

Регуляризация траектории параметров модели глубокого обучения на основе дистилляции знаний

М. Горпинич, О. Ю. Бахтеев, В. В. Стрижов

gorpinich4@gmail.com; bakhteev@phystech.edu; strijov@ccas.ru

Исследуется задача оптимизации параметров модели глубокого обучения. Во время оптимизации учитывается информация, содержащаяся в модели с более сложной структурой, то есть применяется дистилляция. Предлагается обобщение методов дистилляции, заключающееся в градиентной оптимизации метапараметров. Под метапараметрами модели понимаются параметры оптимизационной задачи дистилляции, а именно, коэффициенты перед слагаемыми в функции ошибки и температура. Функция ошибки состоит из двух слагаемых: правдоподобия исходной выборки и правдоподобия выборки дистилляции. Исследуются свойства оптимизационной задачи и методы прогнозирования траектории оптимизации метапараметров модели. Предложенное обобщение позволяет получить модель с лучшими эксплуатационными характеристиками и за меньшее число итераций оптимизации. Проиллюстрирован данный подход с помощью вычислительного эксперимента на выборке CIFAR-10 и на синтетической выборке.

Ключевые слова: машинное обучение; дистилляция знаний; оптимизация метапараметров; градиентные методы оптимизации; прогнозирование метапараметров

DOI:

1 Введение

В работе исследуется задача оптимизации моделей глубоких нейросетей. Проведение оптимизации требует значительных вычислительных мощностей и является затратной по времени. В данной работе предлагается метод оптимизации, позволяющий улучшить точность предсказаний модели, а также ускорить сходимость траектории оптимизации параметров к точке оптимума.

Предлагается обобщение метода оптимизации на основе дистилляции знаний. Назовем *дистилляцией знаний* задачу оптимизации параметров модели прогнозирования, при которой учитывается не только информация, содержащаяся в выборке, но также и информация, содержащаяся в сторонней модели (модели-учителе). Рассматривается *модель-учитель* более сложной структуры, которая была обучена на выборке. Модель более простой структуры предлагается оптимизировать путем переноса знаний модели учителя на более простую модель, называемую *моделью-учеником*. При этом ее качество будет выше по сравнению с качеством, полученным после оптимизации на той же выборке. Данный подход описан в [1]. В [2] предложен подход к дистилляции знаний, переносящий знания на модель с архитектурой, значительно отличающейся от архитектуры модели-учителя.

Предлагается формулировка задачи в виде двухуровневой оптимизации. На первом уровне оптимизируются параметры модели, на втором уровне — ее метапараметры. Данный подход описан в [3–5]. В [3] рассматривается жадный градиентный метод оптимизации метапараметров. В [4] сравниваются различные градиентные методы оптимизации метапараметров, а также метод случайного поиска.

В работе рассматривается подход к прогнозированию значений метапараметров, полученных методом градиентной оптимизации. Под метапараметрами понимаются параметры

задачи оптимизации. Сложность градиентной оптимизации для метапараметров является квадратичной по числу параметров, и потому вычислительно затратна. Предлагается аппроксимация траектории оптимизации метапараметров на основе приближения траектории линейной моделью. Вычислительный эксперимент проводится на выборке изображений CIFAR-10 [6], а также синтетической выборке.

2 Постановка задачи

Решается задача классификации вида:

$$\mathfrak{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m, \quad \mathbf{x}_i \in \mathbb{R}^n, \quad y_i \in \mathbb{Y} = \{1, \dots, K\}, \quad (1)$$

где y_i — это класс объекта, также будем обозначать \mathbf{y}_i вектором вероятности для класса y_i . Разобьем выборку следующим образом:

$$\mathfrak{D} = \mathfrak{D}_{\text{train}} \sqcup \mathfrak{D}_{\text{val}}. \quad (2)$$

Подвыборку $\mathfrak{D}_{\text{train}}$ будем использовать для оптимизации параметров модели, а подвыборку $\mathfrak{D}_{\text{val}}$ — для оптимизации метапараметров.

Внешним критерием качества назначена доля правильных ответов:

$$\text{accuracy} = \frac{1}{m} \sum_{i=1}^m [\mathbf{g}(\mathbf{x}_i, \mathbf{w}) = y_i], \quad (3)$$

где \mathbf{g} — параметрическая модель классификации с параметрами \mathbf{w} .

Определение 1. Траекторией параметров $\boldsymbol{\nu}$ назовем последовательность $\mathbf{w}_1, \dots, \mathbf{w}_t, \dots$ обновления параметров в ходе оптимизации, где t — количество шагов оптимизации.

Определение 2. Пусть задана функция $D : \mathbb{R}^s \rightarrow \mathbb{R}_+$, задающая схожесть модели-ученика \mathbf{g} и фиксированной модели-учителя \mathbf{f} . D -дистилляцией модели-ученика назовем оптимизацию параметров модели-ученика с траекторией $\boldsymbol{\nu}$, такую что $\lim_{t \rightarrow \infty} D(\mathbf{w}_t) = \min_{\mathbf{w}'_t \in \mathbb{R}^s} D(\mathbf{w}'_t)$.

Функция потерь $\mathcal{L}_{\text{train}}$, в которой учитывается перенос информации от модели учителя \mathbf{f} к модели ученика \mathbf{g} , имеет вид:

$$\begin{aligned} \mathcal{L}_{\text{train}}(\mathbf{w}, \boldsymbol{\lambda}) = & -\lambda_1 \sum_{(\mathbf{x}, y) \in \mathfrak{D}_{\text{train}}} \underbrace{\sum_{k=1}^K y^k \log \frac{e^{\mathbf{g}(\mathbf{x}, \mathbf{w})_i/T}}{\sum_j e^{\mathbf{g}(\mathbf{x}, \mathbf{w})_j/T}} \Big|_{T=1}}_{\text{исходная функция потерь}} \\ & - \lambda_2 \sum_{(\mathbf{x}, y) \in \mathfrak{D}_{\text{train}}} \underbrace{\sum_{k=1}^K \frac{e^{\mathbf{f}(\mathbf{x})_i/T}}{\sum_j e^{\mathbf{f}(\mathbf{x})_j/T}} \Big|_{T=T_0} \log \frac{e^{\mathbf{g}(\mathbf{x}, \mathbf{w})_i/T}}{\sum_j e^{\mathbf{g}(\mathbf{x}, \mathbf{w})_j/T}} \Big|_{T=T_0}}_{\text{слагаемое дистилляции}}, \quad (4) \end{aligned}$$

где y^k — k -я компонента целевого вектора (ответ на k -м классе), T — параметр температуры. Параметр температуры T имеет следующие свойства:

- 1) при $T \rightarrow 0$ получаем вектор, в котором один из классов имеет единичную вероятность;
- 2) при $T \rightarrow \infty$ получаем равновероятные классы.

Выражение $\cdot|_{T=t}$ означает, что параметр температуры T в предыдущей функции равняется t .

Зададим множество метапараметров λ как вектор, состоящий из температуры и коэффициента перед слагаемым дистилляции:

$$\lambda = [\lambda_1, \lambda_2, T].$$

Утверждение 1. При $\lambda_1 = 0$ происходит минимизация функции потерь, которая является D -дистилляцией модели с $D = D_{KL}(\sigma(\mathbf{f}/T|_{T=T_0})_i, \sigma(\mathbf{g}/T|_{T=T_0})_i)$.

Доказательство При $\lambda_1 = 0$:

$$\mathcal{L}_{\text{train}}(\mathbf{w}, \lambda) = \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{train}}} \sum_{k=1}^K \frac{e^{\mathbf{f}(\mathbf{x})_i/T}}{\sum_j e^{\mathbf{f}(\mathbf{x})_j/T}|_{T=T_0}} \log \frac{e^{\mathbf{g}(\mathbf{x}, \mathbf{w})_i/T}}{\sum_j e^{\mathbf{g}(\mathbf{x}, \mathbf{w})_j/T}|_{T=T_0}}.$$

По определению D -дистилляции получим, что должно выполняться условие:

$$\lim_{t \rightarrow \infty} D_{KL}(\sigma(\mathbf{f}(\mathbf{x})/T|_{T=T_0})_i, \sigma(\mathbf{g}(\mathbf{x}, \mathbf{w}_t)/T|_{T=T_0})_i) = \min_{\mathbf{w}'_t \in \mathbb{R}^s} D_{KL}(\sigma(\mathbf{f}(\mathbf{x})/T|_{T=T_0})_i, \sigma(\mathbf{g}(\mathbf{x}, \mathbf{w}'_t)/T|_{T=T_0})_i)$$

Используем следующее неравенство:

$$D_{KL}(\sigma(\mathbf{f}(\mathbf{x})/T|_{T=T_0})_i, \sigma(\mathbf{g}(\mathbf{x}, \mathbf{w}'_{k+1}))) \leq D_{KL}(\sigma(\mathbf{f}(\mathbf{x})/T|_{T=T_0})_i, \sigma(\mathbf{g}(\mathbf{x}, \mathbf{w}'_k)/T|_{T=T_0})_i)$$

Также выполняется $D_{KL}(\sigma(\mathbf{f}(\mathbf{x})/T|_{T=T_0})_i, \sigma(\mathbf{g}(\mathbf{x}, \mathbf{w})/T|_{T=T_0})_i) \geq 0$.

Тогда получаем, что

$$\lim_{t \rightarrow \infty} D_{KL}(\sigma(\mathbf{f}(\mathbf{x})/T|_{T=T_0})_i, \sigma(\mathbf{g}(\mathbf{x}, \mathbf{w}_t)/T|_{T=T_0})_i) = \min_{\mathbf{w}'_t \in \mathbb{R}^s} D_{KL}(\sigma(\mathbf{f}(\mathbf{x})/T|_{T=T_0})_i, \sigma(\mathbf{g}(\mathbf{x}, \mathbf{w}'_t)/T|_{T=T_0})_i)$$

□

Итоговая оптимизационная задача выглядит следующим образом:

$$\hat{\lambda} = \arg \max_{\lambda \in \mathbb{R}^3} \mathcal{L}_{\text{val}}(\hat{\mathbf{w}}, \lambda), \quad (5)$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^s} \mathcal{L}_{\text{train}}(\mathbf{w}, \lambda), \quad (6)$$

где функция \mathcal{L}_{val} определяется как:

$$\mathcal{L}_{\text{val}}(\mathbf{w}, \lambda) = \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{val}}} \sum_{k=1}^K y^k \log \mathbf{g}(\mathbf{x}, \mathbf{w})|_{T=1}. \quad (7)$$

3 Градиентные методы оптимизации

В данном разделе рассматриваются детали оптимизации метапараметров градиентными методами.

71 **Определение 3.** Назовем *оператором оптимизации* алгоритм U выбора вектора па-
 72 раметров \mathbf{w}' по параметрам предыдущего шага \mathbf{w} :

$$\mathbf{w}' = U(\mathbf{w}).$$

73 Оптимизируем параметры \mathbf{w} при помощи η шагов оптимизации:

$$\hat{\mathbf{w}} = U \circ U \circ \dots \circ U(\mathbf{w}_0, \boldsymbol{\lambda}) = U^\eta(\mathbf{w}_0, \boldsymbol{\lambda}), \quad (8)$$

75 где \mathbf{w}_0 — начальное значение вектора параметров \mathbf{w} , $\boldsymbol{\lambda}$ — совокупность метапараметров
 76 модели.

77 Переопределим задачу минимизации согласно определению оператора U :

$$\hat{\boldsymbol{\lambda}} = \arg \max_{\boldsymbol{\lambda} \in \mathbb{R}^3} \mathcal{L}_{\text{val}}(U^\eta(\mathbf{w}_0, \boldsymbol{\lambda})). \quad (9)$$

79 Схема оптимизации метапараметров:

- 80 1. Для каждого $i = \overline{0, l}$, где l — число итераций, используемых для оптимизации мета-
 81 параметров.
- 82 2. Решим задачу (9) и получим новое значение метапараметров $\boldsymbol{\lambda}'$.
- 83 3. Положим $\boldsymbol{\lambda} = \boldsymbol{\lambda}'$.

84 Оптимизационную задачу (5) и (6) решает оператор градиентного спуска:

$$U(\mathbf{w}, \boldsymbol{\lambda}) = \mathbf{w} - \gamma \nabla \mathcal{L}_{\text{train}}(\mathbf{w}, \boldsymbol{\lambda}), \quad (10)$$

86 где γ — длина шага градиентного спуска.

87 Используем метод градиентного спуска, который зависит только от значений парамет-
 88 ров \mathbf{w} на предыдущем шаге. На каждой итерации получим следующее значение метапа-
 89 раметров:

$$\boldsymbol{\lambda}' = \boldsymbol{\lambda} - \gamma_{\boldsymbol{\lambda}} \nabla_{\boldsymbol{\lambda}} \mathcal{L}_{\text{val}}(U(\mathbf{w}, \boldsymbol{\lambda}), \boldsymbol{\lambda}) = \boldsymbol{\lambda} - \gamma_{\boldsymbol{\lambda}} \nabla_{\boldsymbol{\lambda}} \mathcal{L}_{\text{val}}(\mathbf{w} - \gamma \nabla \mathcal{L}_{\text{train}}(\mathbf{w}, \boldsymbol{\lambda}), \boldsymbol{\lambda}). \quad (11)$$

91 Градиентная оптимизация является вычислительно затратной, поэтому предлагается
 92 аппроксимировать траекторию оптимизации модели. Предлагается предсказывать траек-
 93 торию изменения метапараметров модели с помощью линейных моделей через определен-
 94 ное число итераций:

$$\boldsymbol{\lambda}' = \boldsymbol{\lambda} + \mathbf{c}^\top \begin{pmatrix} z \\ 1 \end{pmatrix} \quad (12)$$

96 где σ — сигмоида, z — номер итерации по модулю периодичности обучения линейной
 97 модели, \mathbf{c} — коэффициенты линейного многочлена, которые находятся при помощи МНК.
 98 В остальное время используются градиентные методы.

99 4 Вычислительный эксперимент

100 Целью эксперимента является проверка работоспособности предложенного метода ди-
 101 стилляции моделей, а также анализ полученных моделей и их метапараметров. Экспе-
 102 римент проводится на двух выборках: синтетической модели и выборке CIFAR-10. Ре-
 103 зультаты данной работы и исходный код эксперимента опубликованы в [7] и могут быть
 104 проверены или использованы в дальнейшей работе.

4.1 Эксперимент на синтетической выборке

В эксперименте используется синтетическая выборка:

$$\mathfrak{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m, \quad x_{ij} \in (0, 1), \quad j = 1, 2, \quad x_{i3} = [\text{sign}(x_{i1}) + \text{sign}(x_{i2}) > 0]$$

$$y_i = \text{sign}(x_{i1} \cdot x_{i2} + \delta) \in \mathbb{Y},$$

где $\delta \in N(0, 0.5)$ — это шум. При этом размер выборки модели-ученика намного меньше размера выборки модели-учителя и тестовая выборка совпадает с валидационной.

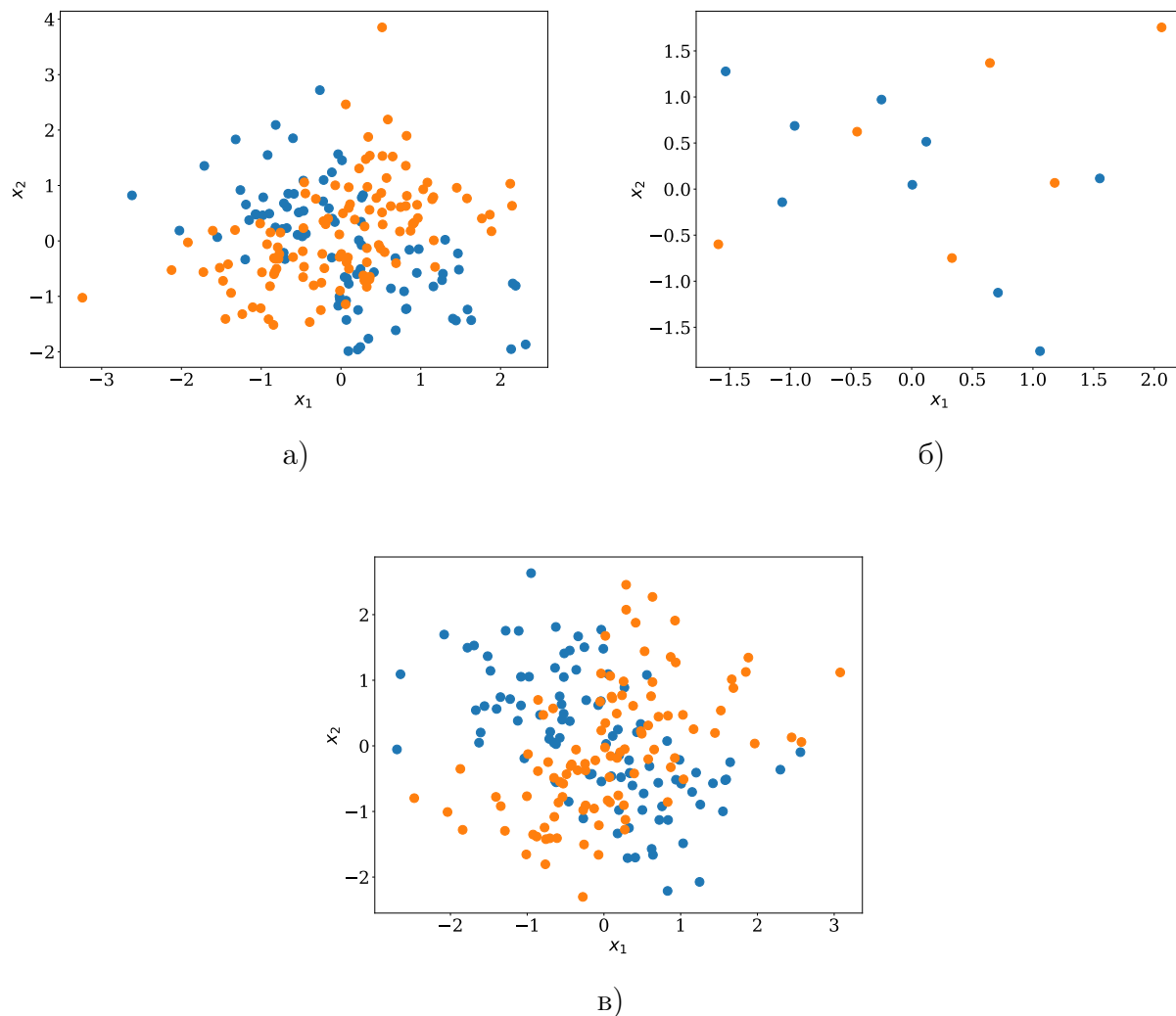


Рис. 1 Визуализация выборки а) для обучения учителя; б) для обучения ученика; в) тестовой выборки

Обучение модели-ученика проводилось несколькими методами: с использованием дистилляции и оптимизации метопараметров градиентными методами, дистилляции с предсказанием траектории оптимизации модели, дистилляции со случайными метопараметрами. При этом для обучения модели с использованием линейной модели дополнительно проводились серии экспериментов для определения наилучшего размера эпохи и наилучшего числа эпох между предсказаниями траектории с помощью линейной модели.

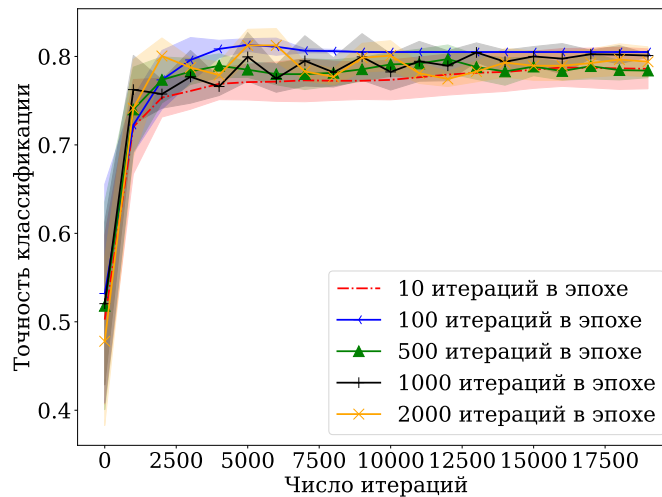


Рис. 2 График зависимости точности классификации от номера итерации при различных значениях размера эпохи

115 На рис. 2 показан график зависимости точности от номера итерации при различных
116 размерах эпохи. Согласно данному графику размер эпохи был выбран равным 100.

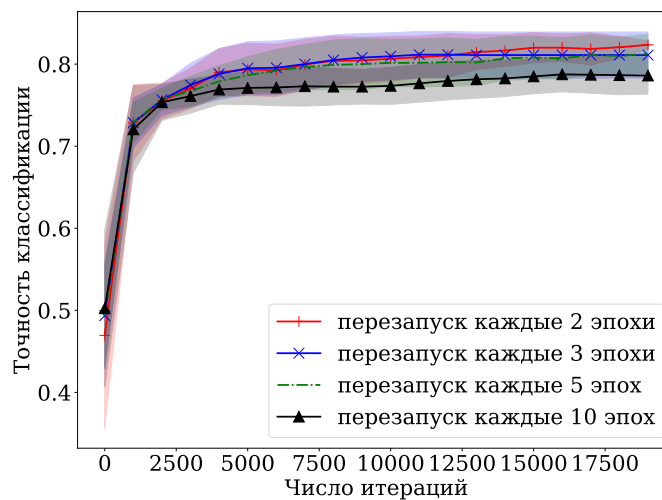


Рис. 3 График зависимости точности классификации от номера итерации при различных n

117 Пусть n — число эпох между использованием линейной модели. На рис. 3 показан
118 график зависимости точности от номера итерации различных n . Наилучшие результаты
119 достигнуты при $n = 2$.

120 На рис. 10 показан график зависимости точности от номера итерации при различных
121 подходах к обучению модели. Наилучшие результаты достигнуты при использовании опти-
122 мизированных метапараметров, но предсказание траектории с помощью линейной модели

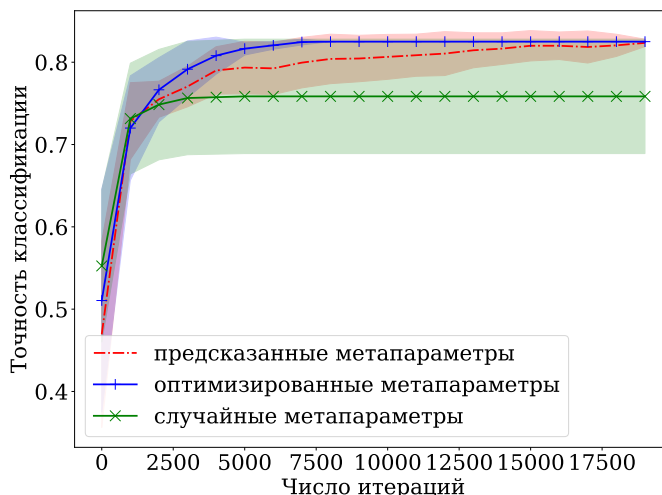


Рис. 4 График зависимости точности классификации от номера итерации

показало результат не намного хуже предыдущего, причем с увеличением числа итераций точность этих двух методов становилась одинаковой.

4.2 Эксперимент на выборке CIFAR-10

В эксперименте используется выборка CIFAR-10, которая состоит из 60000 цветных изображений размера 32×32 пикселя, разделенных на 10 непересекающихся классов. К каждому классу относится 6000 изображений. Выборка делится на обучающую (50000 изображений) и тестовую (10000 изображений) подвыборки. В тестовой выборке содержится 1000 изображений каждого класса.

Внешним критерием качества модели является точность (3). В качестве модели-учителя рассматривается модель из [2], а именно, сверточная нейросеть с тремя сверточными слоями и двумя слоями полносвязной нейросети.

Проведено сравнение среднего качества обучения модели-ученика без дистилляции после 5 запусков, с дистилляцией с моделью-учителем ResNet и сверточной нейросетью после 20 запусков. Значение коэффициента λ_1 лежит в пределах от 0 до 1, значение температуры — от 0.1 до 10.

На рис. 6 изображена зависимость точности от величины коэффициента λ_1 . Различные точки отвечают за точность модели без дистилляции, с дистилляцией ResNet и CNN. Можно заметить, что с уменьшением значения коэффициента λ_1 значение точности увеличивается.

На рис. 7 изображена зависимость точности от T . Для изображения значений температуры используется логарифмическая шкала. По графику видно, что значение температуры уменьшается при увеличении логарифма температуры, но при значениях логарифма от 0.5 до 1 наблюдается резкое уменьшение точности.

На рис. 8 изображена зависимость λ_1 от величины T с выделенной цветом ассигасу. Заметим, что точки с большим значением точности в основном расположены в правом нижнем углу графика, а именно, при значениях λ_1 от 0 до 0.5 и значениях $\log(T)$ от -1 до 0. Наоборот, точки с низким значением точности, расположены в правом верхнем углу графика.

На рис. 9 изображена зависимость метапараметров от числа итераций.

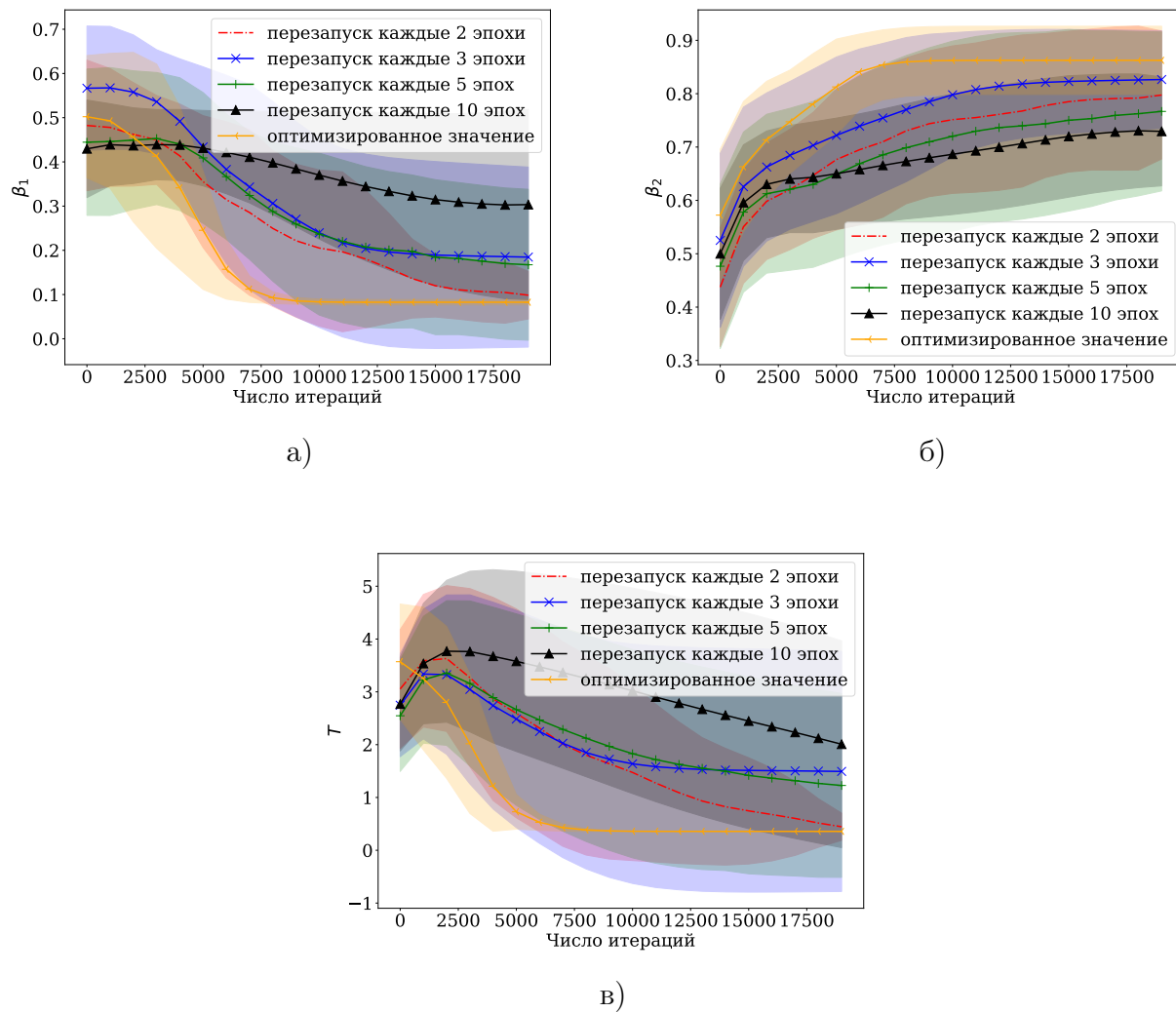


Рис. 5 График зависимости а) λ_1 ; б) λ_2 ; в) температуры от номера итерации

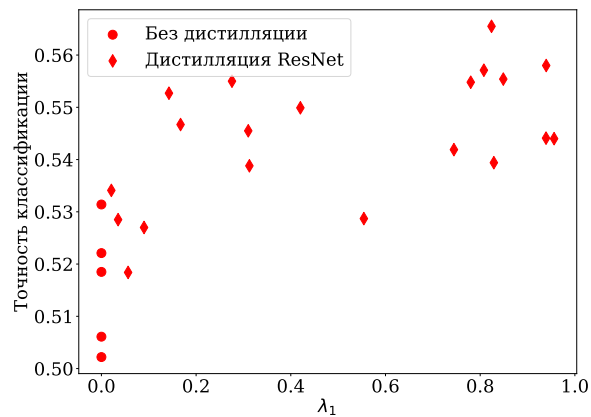


Рис. 6 График зависимости точности от λ_1

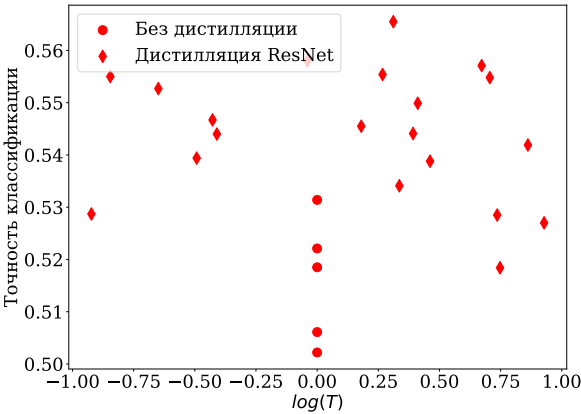


Рис. 7 График зависимости точности от температуры

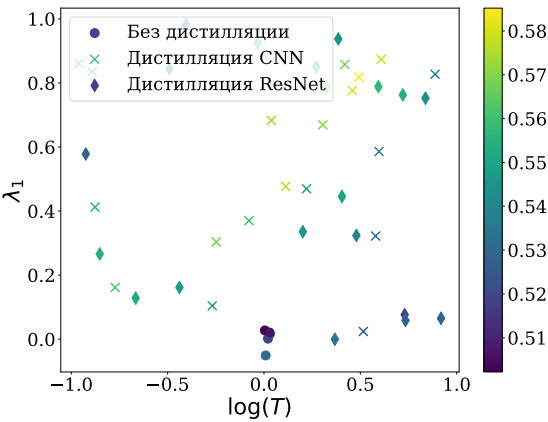


Рис. 8 График зависимости λ_1 от температуры с выделенной цветом ассигасу

На рис. 10 изображена зависимость точности от числа эпох для обучения модели-ученика без дистилляции, обучения с дистилляцией и случайными метапараметрами, обучения с дистилляцией и оптимизацией метапараметров, а также обучения с дистилляцией и оптимальными метапараметрами, полученными в ходе их оптимизации. Можно заметить, что точность обучения с дистилляцией гораздо выше, чем без дистилляции. Также наибольшая точность достигается при обучении с дистилляцией и прогнозированными метапараметрами.

Эксперимент	Точность без дистилляции	Точность с предсказанием метапараметров
Base	0.5465	—
Synthetic	0.76	0.83
CIFAR-10	0.5465	0.5961

Таблица 1 Результаты эксперимента

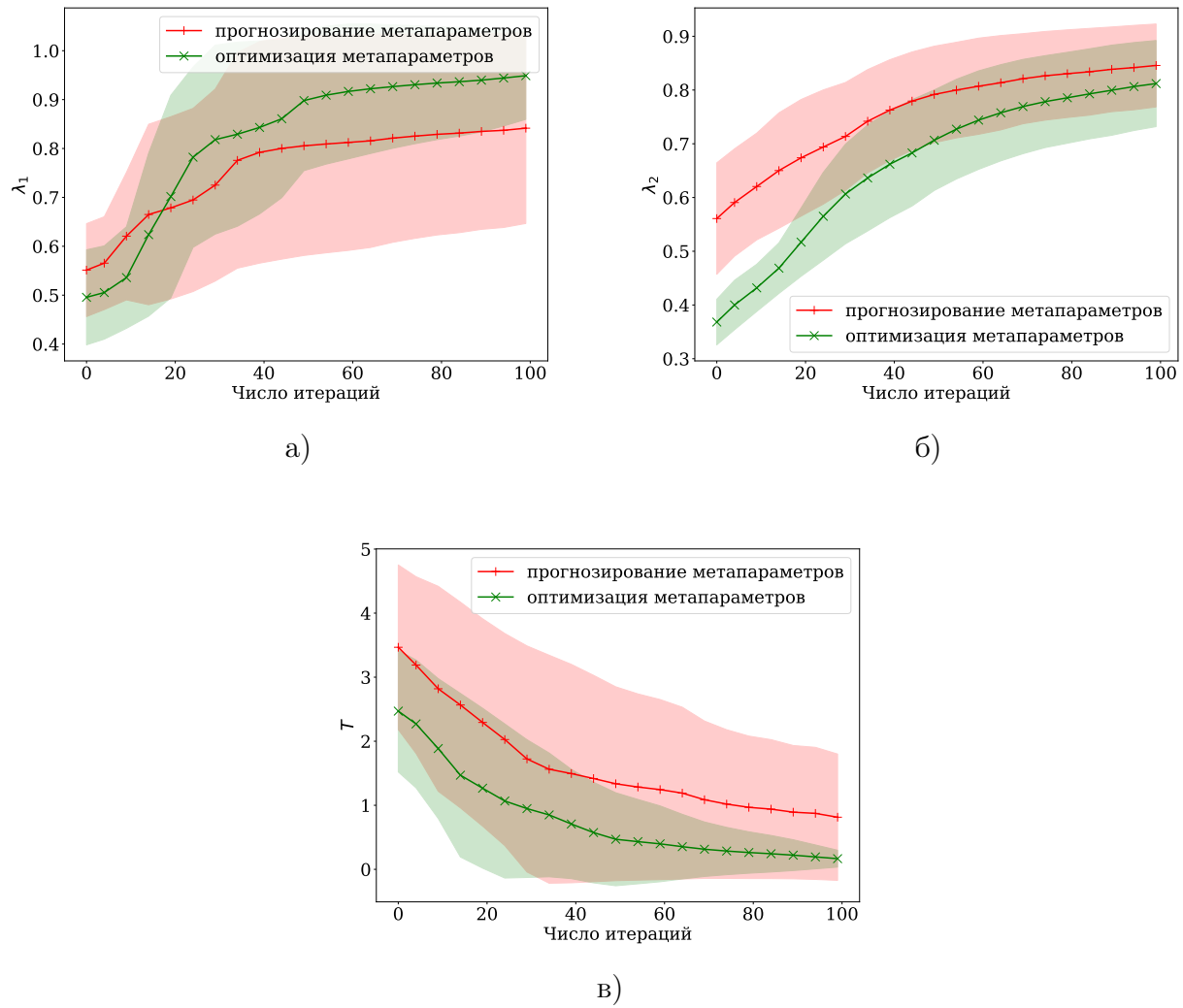


Рис. 9 График зависимости а) λ_1 ; б) λ_2 ; в) температуры от номера итерации

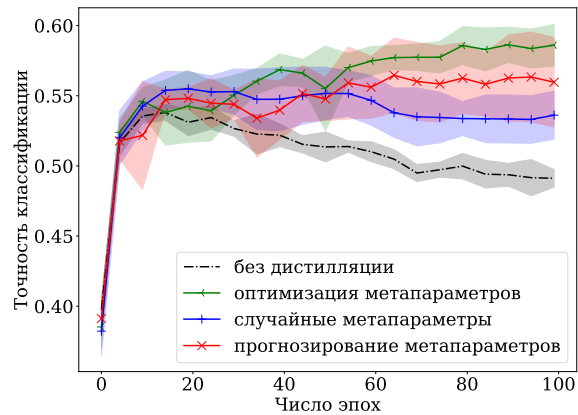


Рис. 10 График зависимости точности от числа эпох

На основании полученных результатов можно сделать вывод о применимости данного метода. Наблюдается выигрыш в точности предсказаний модели при использовании предложенного метода. Также метод является менее вычислительно затратным по сравнению с использованием только градиентной оптимизации.

5 Заключение

Была исследована задача оптимизации параметров модели глубокого обучения. Было предложено обобщение методов дистилляции, заключающееся в градиентной оптимизации метапараметров. На первом уровне оптимизируются параметры модели, на втором — метапараметры, задающие вид оптимизационной задачи. Были исследованы свойства оптимизационной задачи и методы предсказания траектории оптимизации метапараметров модели. Под метапараметрами модели понимаются параметры оптимизационной задачи дистилляции. Предложенное обобщение позволило производить дистилляцию модели с лучшими эксплуатационными характеристиками и за меньшее число итераций оптимизации. Комбинация данных подходов была проиллюстрирована с помощью вычислительного эксперимента на выборке CIFAR-10 и на синтетической выборке. Вычислительный эксперимент показал эффективность градиентной оптимизации для задачи выбора метапараметров дистилляционной функции потерь. Проанализирована возможность аппроксимировать траекторию оптимизации метапараметров локально-линейной моделью. Планируется дальнейшее исследование оптимизационной задачи и анализ качества аппроксимации траектории оптимизации метапараметров более сложными прогностическими моделями.

Литература

- [1] *Hinton Geoffrey E., Vinyals Oriol, Dean Jeffrey.* Distilling the knowledge in a neural network // CoRR, 2015. Vol. abs/1503.02531. URL: <http://arxiv.org/abs/1503.02531>.
- [2] *Passalis Nikolaos, Tzelepi Maria, Tefas Anastasios.* Heterogeneous knowledge distillation using information flow modeling // CVPR. — 2020. P. 2336–2345. URL: <https://ieeexplore.ieee.org/xpl/conhome/9142308/proceeding>.
- [3] *Luketina Jelena, Berglund Mathias, Greff Klaus, Raiko Tapani.* Scalable gradient-based tuning of continuous regularization hyperparameters // CoRR, 2015. Vol. abs/1511.06727. URL: <http://arxiv.org/abs/1511.06727>.
- [4] *Bakhteev Oleg Yu., Strijov Vadim V.* Comprehensive analysis of gradient-based hyperparameter optimization algorithms // Ann. Oper. Res, 2020. Vol. 289. No. 1. P. 51–65.
- [5] *Maclaurin Dougal, Duvenaud David, Adams Ryan P.* Gradient-based hyperparameter optimization through reversible learning // CoRR, 2015. Vol. abs/1502.03492. URL: <http://arxiv.org/abs/1502.03492>.
- [6] *Krizhevsky Alex et al.* Learning multiple layers of features from tiny images, 2009.
- [7] URL: <https://github.com/Intelligent-Systems-Phystech/2021-Project-84>.

Received

Regularizing optimization trajectory of deep learning model parameters with knowledge distillation

M. Gorpinich, O. Yu. Bakhteev, V. V. Strijov

gorpinich4@gmail.com; bakhteev@phystech.edu; strijov@ccas.ru

The paper investigates parameter optimization problem for deep learning neural networks. The knowledge of a cumbersome model is considered during optimization, i.e. the knowledge distillation is used. The paper proposes generalization of knowledge distillation method to optimize meta-parameters by gradient descent. Meta-parameters are the parameters of knowledge distillation optimization problem, namely, the coefficients before terms in error function and the temperature factor. The error function is a sum of likelihood of the initial dataset and the one of distillation dataset. Temperature is a factor of logits of models in softmax function. The authors investigate the properties of optimization problem and methods to predict the optimization path of meta-parameters. Generalized method produces models with higher performance and uses less number of iterations. The algorithm is evaluated on CIFAR-10 dataset and synthetic data.

Keywords: *machine learning; knowledge distillation; metaparameter optimization; gradient-based optimization; metaparameter selection*

DOI:

References

- [1] *Hinton Geoffrey E., Vinyals Oriol, Dean Jeffrey.* Distilling the knowledge in a neural network // CoRR, 2015. Vol. abs/1503.02531. URL: <http://arxiv.org/abs/1503.02531>.
- [2] *Passalis Nikolaos, Tzelepi Maria, Tefas Anastasios.* Heterogeneous knowledge distillation using information flow modeling // CVPR. — 2020. P. 2336–2345. URL: <https://ieeexplore.ieee.org/xpl/conhome/9142308/proceeding>.
- [3] *Luketina Jelena, Berglund Mathias, Greff Klaus, Raiko Tapani.* Scalable gradient-based tuning of continuous regularization hyperparameters // CoRR, 2015. Vol. abs/1511.06727. URL: <http://arxiv.org/abs/1511.06727>.
- [4] *Bakhteev Oleg Yu., Strijov Vadim V.* Comprehensive analysis of gradient-based hyperparameter optimization algorithms // Ann. Oper. Res, 2020. Vol. 289. No. 1. P. 51–65.
- [5] *Maclaurin Dougal, Duvenaud David, Adams Ryan P.* Gradient-based hyperparameter optimization through reversible learning // CoRR, 2015. Vol. abs/1502.03492. URL: <http://arxiv.org/abs/1502.03492>.
- [6] *Krizhevsky Alex et al.* Learning multiple layers of features from tiny images, 2009.
- [7] URL: <https://github.com/Intelligent-Systems-Phystech/2021-Project-84>.

Received