

Регуляризация траектории оптимизации 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. Предлагается обобщение методов дистилляции, заключающееся в градиентной оптимизации метопараметров. На первом уровне оптимизируются параметры модели, на втором — метопараметры, задающие вид оптимизационной задачи. Исследуются свойства оптимизационной задачи и методы предсказания траектории оптимизации метопараметров модели. Под метопараметрами модели понимаются параметры оптимизационной задачи дистилляции. Предложенное обобщение позволяет производить дистилляцию модели с лучшими эксплуатационными характеристиками и за меньшее число итераций оптимизации. Проиллюстрирована комбинация данных подходов с помощью вычислительного эксперимента на выборке CIFAR-10 и на синтетической выборке.

Ключевые слова: *machine learning; knowledge distillation; hyperparameter optimization*

DOI:

1 Введение

В работе исследуется проблема оптимизации моделей глубоких нейросетей. Данная задача требует значительных вычислительных мощностей и является затратной по времени. В данной работе предлагается метод оптимизации, позволяющий улучшить эксплуатационные характеристики модели, а также ускорить ее сходимость к точке оптимума.

Предлагается обобщение метода оптимизации на основе дистилляции знаний. 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, переносащий знания на модель с архитектурой, значительно отличающейся от архитектуры модели-учителя.

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

В работе рассматривается подход к прогнозированию метопараметров, полученных методом градиентной оптимизации. Под метопараметрами понимаются параметры задачи оптимизации. Сложность градиентной оптимизации для метопараметров является квадратичной по числу параметров, и потому вычислительно затратна. Предлагается аппроксимация траектории оптимизации метопараметров 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 метопараметров.

В качестве внешнего критерия качества рассматривается доля правильных ответов:

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

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

Определение 1. 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).

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

$$\mathcal{L}_{\text{train}}(\mathbf{w}, \boldsymbol{\lambda}) = -\beta_1 \underbrace{\sum_{(\mathbf{x}, y) \in \mathfrak{D}_{\text{train}}} \sum_{k=1}^K y^k \log \mathbf{g}(\mathbf{x}, \mathbf{w})|_{T=1}}_{\text{исходная функция потерь}} - \beta_2 \underbrace{\sum_{(\mathbf{x}, y) \in \mathfrak{D}_{\text{train}}} \sum_{k=1}^K \mathbf{f}(\mathbf{x})|_{T=T_0} \log \mathbf{g}(\mathbf{x}, \mathbf{w})|_{T=T_0}}_{\text{слагаемое дистилляции}}, \quad (4)$$

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

Зададим множество метопараметров $\boldsymbol{\lambda}$ как вектор, состоящий из температуры и коэффициента перед слагаемым дистилляции:

$$\boldsymbol{\lambda} = [\beta_1, \beta_2, T].$$

Итоговая оптимизационная задача:

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

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

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

$$56 \quad \mathcal{L}_{\text{val}}(\mathbf{w}, \boldsymbol{\lambda}) = \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{val}}} \sum_{k=1}^K y^k \log \mathbf{g}(\mathbf{x}, \mathbf{w})|_{T=1}. \quad (7)$$

57 **Определение 2.** Назовем *оператором оптимизации* алгоритм U выбора вектора па-
58 раметров \mathbf{w}' по параметрам предыдущего шага \mathbf{w} :

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

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

$$60 \quad \hat{\mathbf{w}} = U \circ U \circ \dots \circ U(\mathbf{w}_0, \boldsymbol{\lambda}) = U^\eta(\mathbf{w}_0, \boldsymbol{\lambda}), \quad (8)$$

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

63 Переопределим задачу минимизации согласно определению оператора U :

$$64 \quad \hat{\boldsymbol{\lambda}} = \arg \max_{\boldsymbol{\lambda} \in \mathbb{R}^2} \mathcal{L}_{\text{val}}(U^\eta(\mathbf{w}_0, \boldsymbol{\lambda})). \quad (9)$$

65 Схема оптимизации метапараметров:

- 66 1. Для каждого $i = \overline{0, l}$, где l — количество итераций, используемых для оптимизации
67 метапараметров.
- 68 2. Решим задачу (9) и получим новое значение метапараметров $\boldsymbol{\lambda}'$.
- 69 3. Положим $\boldsymbol{\lambda} = \boldsymbol{\lambda}'$.

70 3 Градиентные методы оптимизации

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

$$72 \quad U(\mathbf{w}, \boldsymbol{\lambda}) = \mathbf{w} - \gamma \nabla \mathcal{L}_{\text{train}}(\mathbf{w}, \boldsymbol{\lambda}), \quad (10)$$

73 где γ — длина шага градиентного спуска.

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

$$77 \quad \boldsymbol{\lambda}' = \boldsymbol{\lambda} - \gamma_{\boldsymbol{\lambda}} \nabla_{\boldsymbol{\lambda}} \mathcal{L}_{\text{val}}(U(\mathbf{w}, \boldsymbol{\lambda}), \boldsymbol{\lambda}) = \boldsymbol{\lambda} - \gamma_{\boldsymbol{\lambda}} \nabla_{\boldsymbol{\lambda}} \mathcal{L}_{\text{val}}(\mathbf{w} - \gamma \nabla \mathcal{L}_{\text{train}}(\mathbf{w}, \boldsymbol{\lambda}), \boldsymbol{\lambda}). \quad (11)$$

78 Градиентная оптимизация является вычислительно затратной, поэтому предлагается
79 аппроксимировать траекторию оптимизации модели.

80 Предлагается предсказывать траекторию изменения метапараметров модели (а кон-
81 кретно, их градиенты) с помощью линейных сплайнов через определенное число итераций,
82 а в остальное время использовать градиентные методы.

83 4 Вычислительный эксперимент

84 Целью эксперимента является проверка работоспособности предложенного метода ди-
85 стилляции моделей, а также анализ полученных моделей и их метапараметров. Экспери-
86 мент проводится на двух выборках: синтетической модели и выборке CIFAR-10.

87 4.1 Эксперимент на синтетической выборке

88 В эксперименте используется синтетическая выборка с тремя признаками у каждого
 89 объекта. Первые два признака сэмплируются из стандартного нормального распределе-
 90 ния, третий признак — это индикатор того, что первые два признака больше 0. Ответами
 91 для выборки модели-учителя являются значения $\text{sign}(x_1 + x_2)$, где x_1, x_2 — первые два
 92 признака. Ответами для выборки модели-ученика являются $\text{sign}(x_1 + x_2) + \delta$, где δ — это
 93 шум.

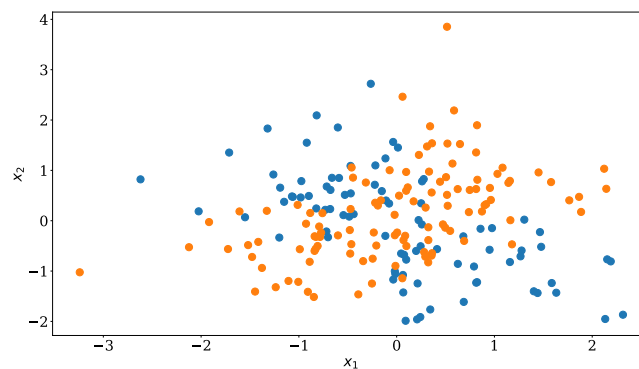


Рис. 1 Выборка для обучения учителя

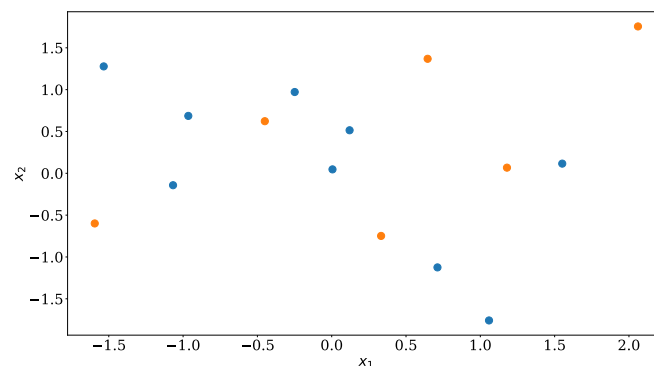


Рис. 2 Выборка для обучения ученика

94 Обучение модели-ученика проводилось несколькими методами: с использованием ди-
 95 стиляции и оптимизации метопараметров градиентными методами, дистилляции с пред-
 96 сказанием траектории оптимизации модели, дистилляции со случайными метопарамет-
 97 рами. При этом для обучения модели с использованием сплайнов дополнительно прово-
 98 дились серии экспериментов для определения наилучшего размера эпохи и наилучшего
 99 количества эпох между предсказаниями траектории с помощью сплайнов.

100 На рис. 4 показан график зависимости точности от номера итерации при различных
 101 размерах эпохи. Согласно данному графику размер эпохи был выбран равным 100.

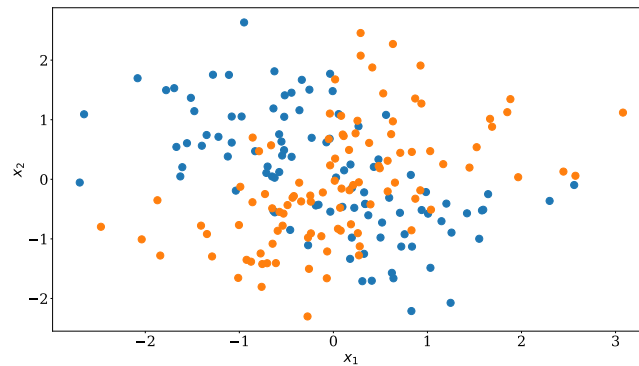


Рис. 3 Тестовая выборка

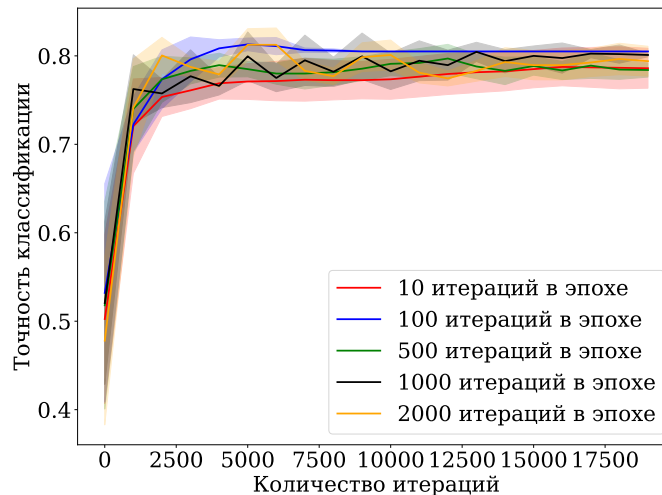


Рис. 4 График зависимости точности классификации от номера итерации при различных значениях размера эпохи

Пусть n — количество эпох между использованием сплайнов. На рис. 5 показан график зависимости точности от номера итерации различных n . Наилучшие результаты достигнуты при $n = 2$.

На рис. 13 показан график зависимости точности от номера итерации при различных подходах к обучению модели. Наилучшие результаты достигнуты при использовании оптимизированных гиперпараметров, но предсказание траектории с помощью сплайнов показало результат не намного хуже предыдущего, причем с увеличением количества итераций точность этих двух методов становилась одинаковой.

На рис. 7, 8 и 9 показаны графики обновления метопараметров β_1 , β_2 и T . Из графика видно, что при $n = 2$, предсказанная траектория наиболее близка к траектории, полученной с помощью только градиентных методов.

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

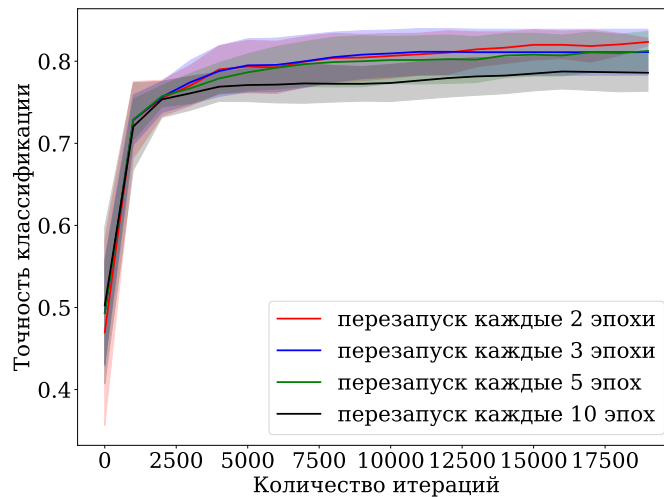


Рис. 5 График зависимости точности классификации от номера итерации при различных n

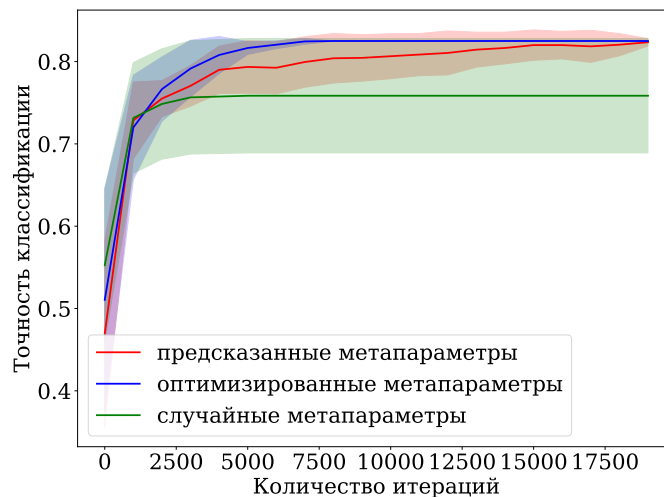


Рис. 6 График зависимости точности классификации от номера итерации

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 моделей-учителей рассматриваются модели из [2], а именно, ResNet-18 и 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 β лежит в пределах от 0 до 1, значение температуры — от 0.1 до 10.

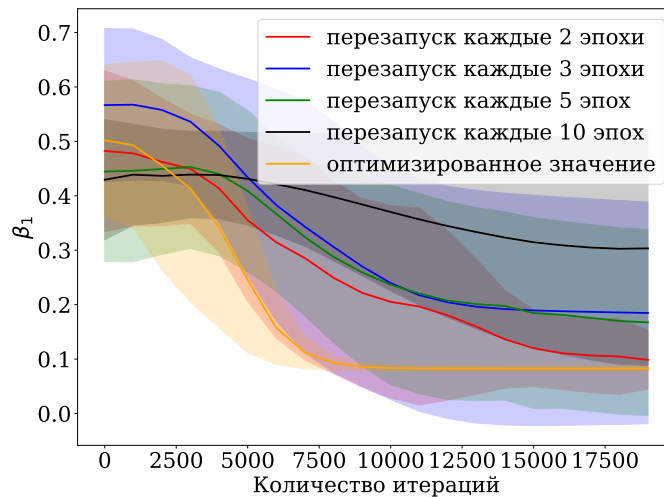


Рис. 7 График зависимости значения β_1 от номера итерации

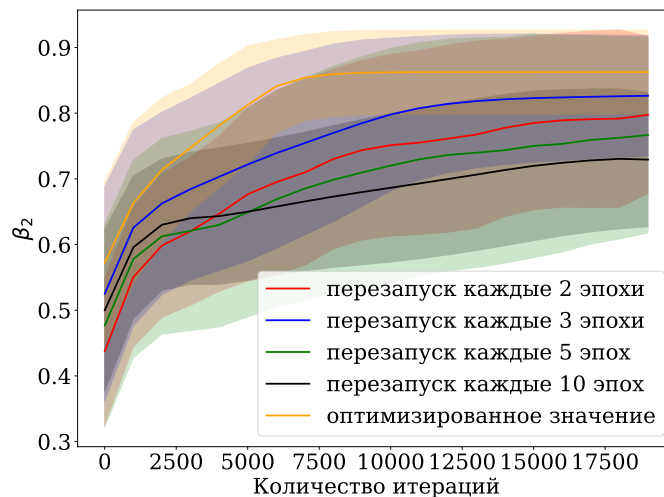


Рис. 8 График зависимости значения β_2 от номера итерации

На рис. 10 изображена зависимость точности от величины коэффициента β . Различные точки отвечают за точность модели без дистилляции, с дистилляцией ResNet и CNN. Можно заметить, что с уменьшением значения коэффициента β значение точности увеличивается.

На рис. 11 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. 12 изображена зависимость β 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 β от 0 до 0.5 и значениях $\log(T)$ от -1

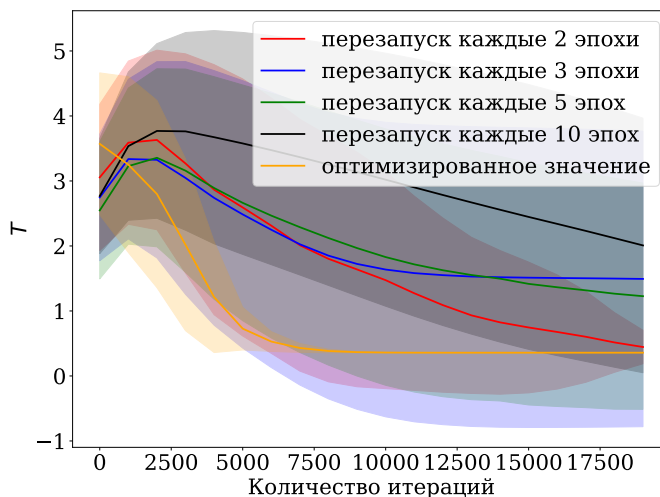


Рис. 9 График зависимости значения T от номера итерации

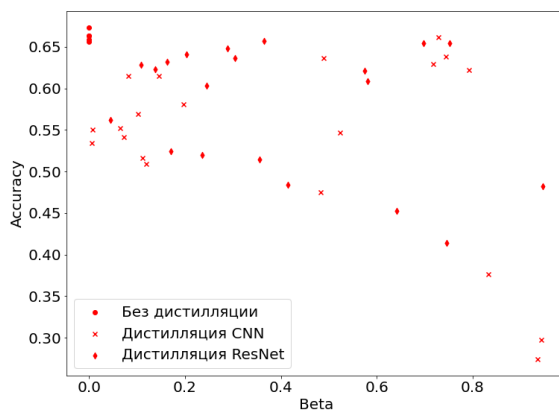


Рис. 10 График зависимости точности от β

137 до 0. Наоборот, точки с низким значением точности, расположены в правом верхнем углу
138 графика.

139 На рис. 13 изображена зависимость точности от количества эпох для обучения модели-
140 ученика без дистилляции, обучения с дистилляцией и случайными метапараметрами, обу-
141 чения с дистилляцией и оптимизацией метапараметров, а также обучения с дистилляцией
142 и оптимальными метапараметрами, полученными в ходе их оптимизации. Можно заме-
143 тить, что точность обучения с дистилляцией гораздо выше, чем без дистилляции. Также
144 наибольшая точность достигается при обучении с дистилляцией и оптимальными метапа-
145 раметрами.

Таблица 1 Результаты эксперимента

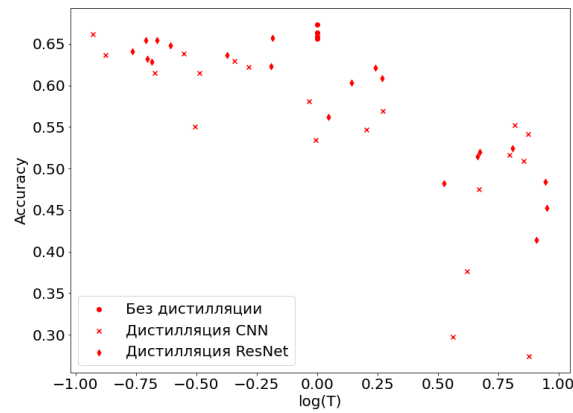


Рис. 11 График зависимости точности от температуры

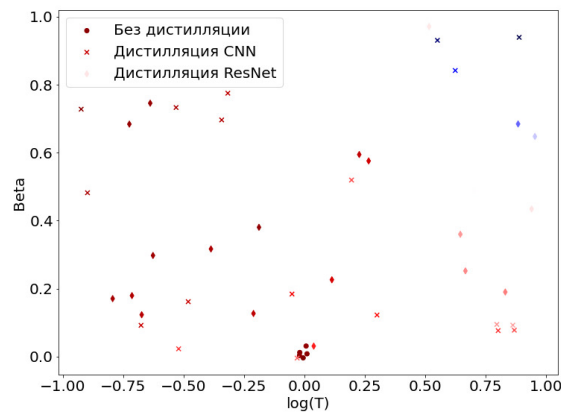


Рис. 12 График зависимости β от температуры с выделенной цветом ассигасу

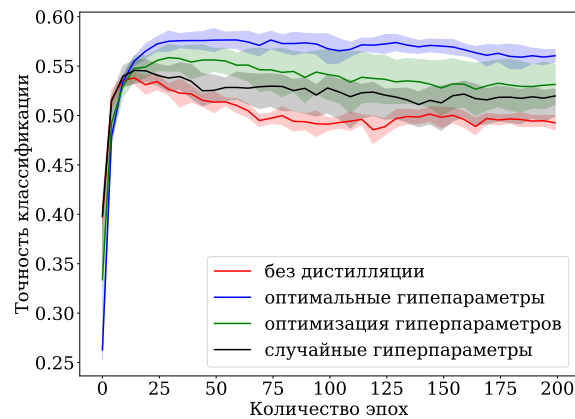


Рис. 13 График зависимости точности от количества эпох

Рис. 14 График зависимости β от количества итераций дистилляции

Рис. 15 График зависимости температуры от количества итераций дистилляции

На рис. 14 изображена зависимость β от количества итераций дистилляции.

На рис. 15 изображена зависимость T от количества итераций дистилляции.

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.

Received