

# 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. Проиллюстрирован данный подход 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:

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

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  $\mathcal{L}_{\text{val}}$  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  $\eta$  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  $\gamma$  — 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. Резуль-  
87 таты данной работы опубликованы в [7] и могут быть проверены или использованы в  
88 дальнейшей работе.

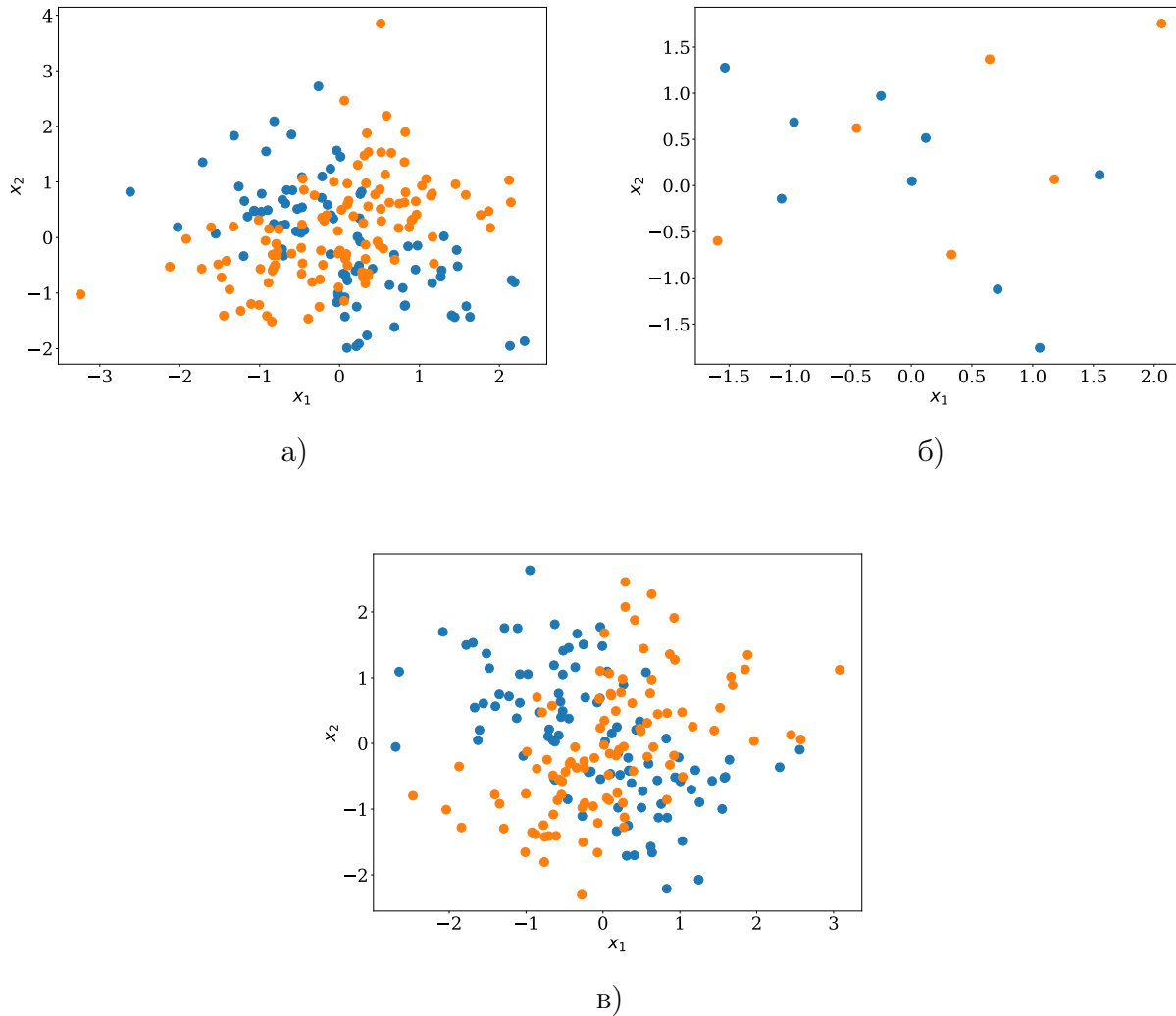
#### 4.1 Эксперимент на синтетической выборке

В эксперименте используется синтетическая выборка:

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

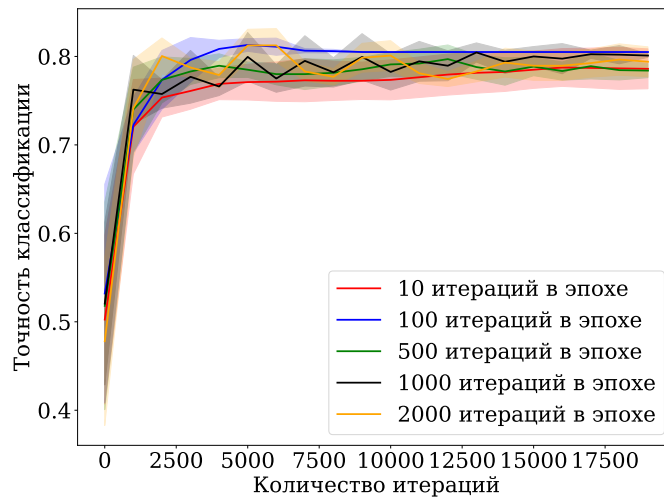
$$y_i = \text{sign}(x_{i1} * x_{i2} + \delta) \in \mathbb{Y},$$

где  $\delta$  — это шум. При этом размер выборки модели-ученика намного меньше размера выборки модели-учителя.



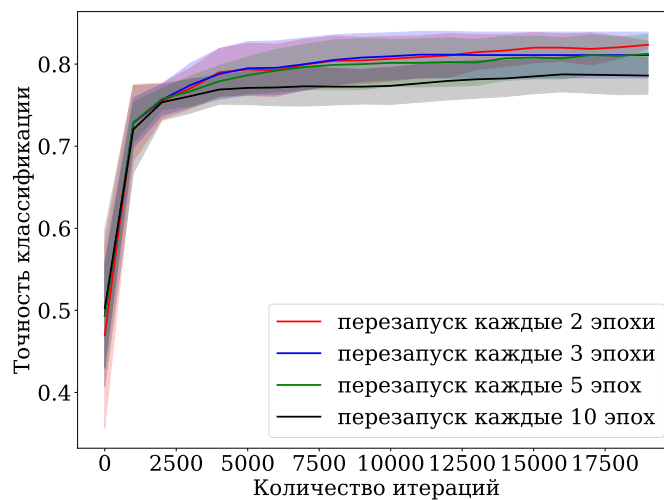
**Рис. 1** Визуализация выборки а) для обучения учителя; б) для обучения ученика; в) тестовой выборки

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



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

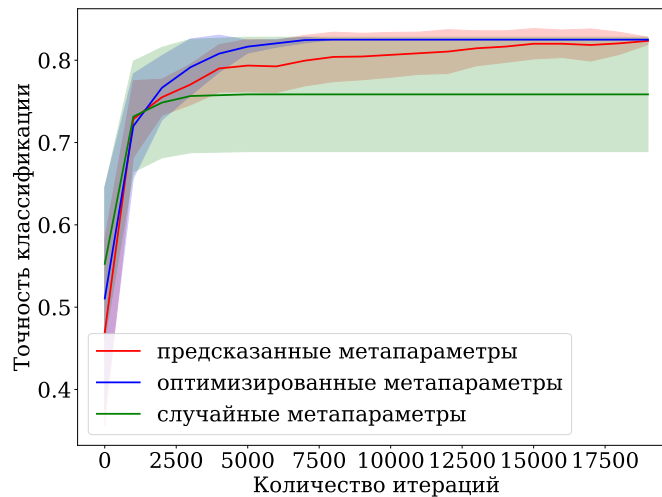
99 На рис. 2 показан график зависимости точности от номера итерации при различных  
 100 размерах эпохи. Согласно данному графику размер эпохи был выбран равным 100.



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

101 Пусть  $n$  — количество эпох между использованием сплайнов. На рис. 3 показан график  
 102 зависимости точности от номера итерации различных  $n$ . Наилучшие результаты достиг-  
 103 нуты при  $n = 2$ .

104 На рис. 9 показан график зависимости точности от номера итерации при различных  
 105 подходах к обучению модели. Наилучшие результаты достигнуты при использовании опти-  
 106 мизированных гиперпараметров, но предсказание траектории с помощью сплайнов пока-



**Рис. 4** График зависимости точности классификации от номера итерации

107 зало результат не намного хуже предыдущего, причем с увеличением количества итераций  
 108 точность этих двух методов становилась одинаковой.

## 109 4.2 Эксперимент на выборке CIFAR-10

110 В эксперименте используется выборка CIFAR-10, которая состоит из 60000 цветных  
 111 изображений размера  $32 \times 32$  пикселя, разделенных на 10 непересекающихся классов. К  
 112 каждому классу относится 6000 изображений. Выборка делится на обучающую (50000  
 113 изображений) и тестовую (10000 изображений) подвыборки. В тестовой выборке содер-  
 114 жится 1000 изображений каждого класса.

115 Внешним критерием качества модели является точность (3). В качестве моделей-  
 116 учителей рассматриваются модели из [2], а именно, ResNet-18 и сверточная нейросеть  
 117 с тремя сверточными слоями и двумя слоями полносвязной нейросети.

118 Проведено сравнение среднего качества обучения модели-ученика без дистилляции по-  
 119 сле 5 запусков, с дистилляцией с моделью-учителем ResNet и сверточной нейросетью после  
 120 20 запусков. Значение коэффициента  $\beta$  лежит в пределах от 0 до 1, значение температу-  
 121 ры — от 0.1 до 10.

122 На рис. 6 изображена зависимость точности от величины коэффициента  $\beta$ . Различ-  
 123 ные точки отвечают за точность модели без дистилляции, с дистилляцией ResNet и CNN.  
 124 Можно заметить, что с уменьшением значения коэффициента  $\beta$  значение точности уве-  
 125 личивается.

126 На рис. 7 изображена зависимость точности от  $T$ . Для изображения значений темпера-  
 127 туры используется логарифмическая шкала. По графику видно, что значение температу-  
 128 ры уменьшается при увеличении логарифма температуры, но при значениях логарифма  
 129 от 0.5 до 1 наблюдается резкое уменьшение точности.

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

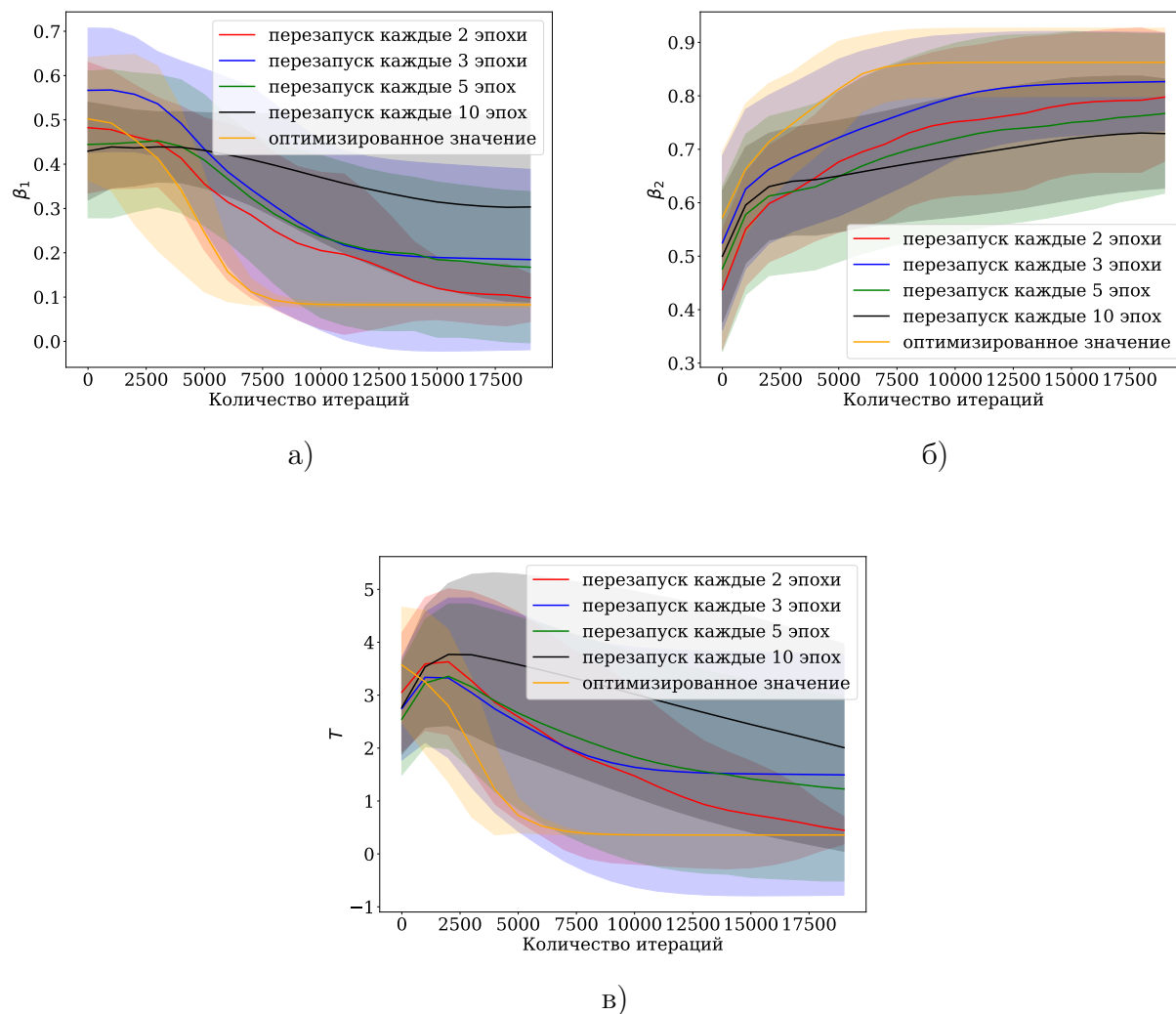


Рис. 5 График зависимости а)  $\beta_1$ ; б)  $\beta_2$ ; в) температуры от номера итерации

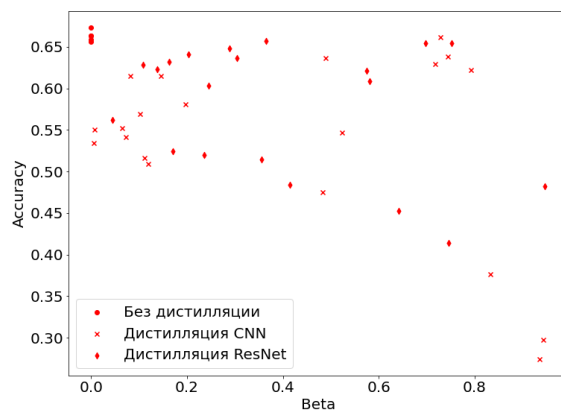
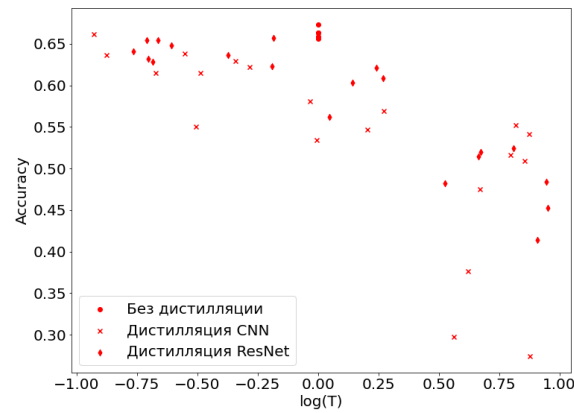
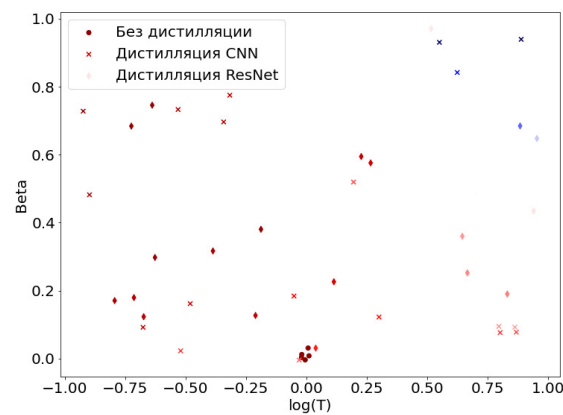


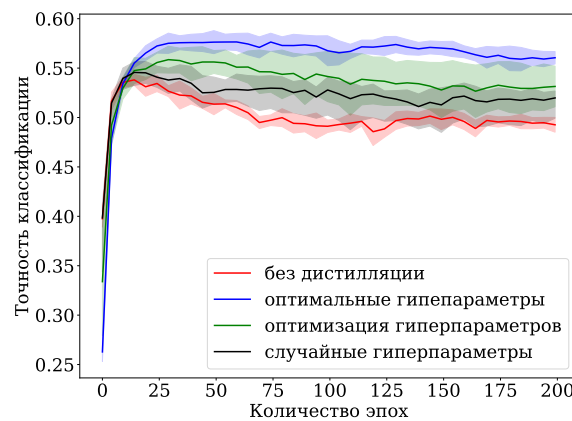
Рис. 6 График зависимости точности от  $\beta$



**Рис. 7** График зависимости точности от температуры



**Рис. 8** График зависимости  $\beta$  от температуры с выделенной цветом ассигасу



**Рис. 9** График зависимости точности от количества эпох

135 На рис. 9 изображена зависимость точности от количества эпох для обучения модели-  
 136 ученика без дистилляции, обучения с дистилляцией и случайными метопараметрами, обу-



137 чения с дистилляцией и оптимизацией метапараметров, а также обучения с дистилляцией  
 138 и оптимальными метапараметрами, полученными в ходе их оптимизации. Можно заме-  
 139 тить, что точность обучения с дистилляцией гораздо выше, чем без дистилляции. Также  
 140 наибольшая точность достигается при обучении с дистилляцией и оптимальными метапа-  
 141 раметрами.

Таблица 1 Результаты эксперимента

142 В таблице 1 приведены результаты эксперимента.

Рис. 10 График зависимости  $\beta$  от количества итераций дистилляции

Рис. 11 График зависимости температуры от количества итераций дистилляции

143 На рис. 10 изображена зависимость  $\beta$  от количества итераций дистилляции.

144 На рис. 11 изображена зависимость  $T$  от количества итераций дистилляции.

## 145 5 Заключение

146 Была исследована задача оптимизации параметров модели глубокого обучения. Было  
 147 предложено. обобщение методов дистилляции, заключающееся в градиентной оптимиза-  
 148 ции метапараметров. На первом уровне оптимизируются параметры модели, на втором —  
 149 метапараметры, задающие вид оптимизационной задачи. Были исследованы свойства оп-  
 150 тимизационной задачи и методы предсказания траектории оптимизации метапараметров  
 151 модели. Под метапараметрами модели понимаются параметры оптимизационной задачи  
 152 дистилляции. Предложенное обобщение позволило производить дистилляцию модели с  
 153 лучшими эксплуатационными характеристиками и за меньшее число итераций оптимиза-  
 154 ции. Комбинация данных подходов была проиллюстрирована с помощью вычислительного  
 155 эксперимента на выборке CIFAR-10 и на синтетической выборке. Вычислительный экспе-  
 156 римент показал эффективность градиентной оптимизации для задачи выбора метарапара-  
 157 метров дистилляционной функции потерь. Проанализирована возможность аппроксими-  
 158 ровать траекторию оптимизации метапараметров локально-линейной моделью. Планиру-  
 159 ется дальнейшее исследование оптимизационной задачи и анализ качества аппроксимации  
 160 траектории оптимизации метапараметров более сложными прогностическими моделями.

## 161 Литература

- 162 [1] *Hinton Geoffrey E., Vinyals Oriol, Dean Jeffrey*. Distilling the knowledge in a neural network //  
 163 CoRR, 2015. Vol. abs/1503.02531. URL: <http://arxiv.org/abs/1503.02531>.
- 164 [2] *Passalis Nikolaos, Tzelepi Maria, Tefas Anastasios*. Heterogeneous knowledge distillation using  
 165 information flow modeling // CVPR. — 2020. P. 2336–2345. URL: <https://ieeexplore.ieee.org/xpl/conhome/9142308/proceeding>.
- 166 [3] *Luketina Jelena, Berglund Mathias, Greff Klaus, Raiko Tapani*. Scalable gradient-based tuning  
 167 of continuous regularization hyperparameters // CoRR, 2015. Vol. abs/1511.06727. URL: <http://arxiv.org/abs/1511.06727>.

- 170 [4] *Bakhteev Oleg Yu., Strijov Vadim V.* Comprehensive analysis of gradient-based hyperparameter  
171 optimization algorithms // *Ann. Oper. Res.*, 2020. Vol. 289. No. 1. P. 51–65.
- 172 [5] *Maclaurin Dougal, Duvenaud David, Adams Ryan P.* Gradient-based hyperparameter optimization  
173 through reversible learning // *CoRR*, 2015. Vol. abs/1502.03492. URL: [http://arxiv.org/abs/](http://arxiv.org/abs/1502.03492)  
174 [1502.03492](http://arxiv.org/abs/1502.03492).
- 175 [6] *Krizhevsky Alex et al.* Learning multiple layers of features from tiny images, 2009.
- 176 [7] URL: <https://github.com/Intelligent-Systems-Phystech/2021-Project-84>.

177 *Received*