Регуляризация траектории оптимизации параметров модели глубокого обучения на основе дистилляции знаний

M. 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 гиперпараметров. На первом уровне оптимизируются параметры модели, на втором — гиперпараметры, задающие вид оптимизационной задачи. Исследуются свойства оптимизационной задачи и различные виды оператора оптимизации. 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.

Ключевые слова: machine learning; knowledge distillation; hyperparameter optimization; recurrent neural network

DOI:

₋ 1 Введение

4

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

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, переносящий 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 гиперпараметры. Данный подход описан в [3–5]. В [3] рассматривается жадный градиентный метод оптимизации гиперпараметров. В [4] сравниваются различные градиентные методы оптимизации гиперпараметров, а также метод случайного поиска.

В работе рассматривается вид задачи оптимизации, а также различные виды оператора оптимизации. Данный подход с использованием нейросети LSTM описан в [6]. Вычислительный эксперимент проводится на выборке изображений CIFAR-10 [7].

2 Постановка задачи

Решается задача классификации вида:

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

М. Горпинич и др.

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

Разобьем выборку следующим образом:

27

31

32

36

44

45

46

47

49

$$\mathfrak{D} = \mathfrak{D}_{\text{train}} \sqcup \mathfrak{D}_{\text{val}}. \tag{2}$$

Подвыборку \mathfrak{D}_{train} будем использовать для оптимизации параметров модели, а подвызо борку \mathfrak{D}_{val} — для оптимизации гиперпараметров.

В качестве внешнего критерия качества рассматривается доля правильных ответов:

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

 ${f g}$ где ${f g}$ — параметрическая модель классификации с параметрами ${f w}.$

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

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

 $_{
m 37}~$ где T- параметр температуры. Параметр температуры T имеет следующие свойства:

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

⁴⁰ Выражение $\cdot|_{T=t}$ обозначает, что параметр температуры T в предыдущей функции равня⁴¹ ется t.

Зададим множество гиперпараметров **h** как вектор, состоящий из температуры и коэффициента перед слагаемым дистилляции:

$$\mathbf{h} = [\beta, T].$$

Итоговая оптимизационная задача:

$$\hat{\mathbf{h}} = \arg\max_{\mathbf{h} \in \mathbb{R}^2} \mathcal{L}_{\text{val}}(\hat{\mathbf{w}}, \mathbf{h}), \tag{5}$$

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^s}{\min} \, \mathcal{L}_{\text{train}}(\mathbf{w}, \mathbf{h}), \tag{6}$$

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

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

Определение 1. Назовем *оператором оптимизации* алгоритм U выбора вектора параметров \mathbf{w}' по параметрам предыдущего шага \mathbf{w} :

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

Оптимизируем параметры **w** при помощи η шагов оптимизации:

$$\hat{\mathbf{w}} = U \circ U \circ \dots \circ U(\mathbf{w}_0, \mathbf{h}) = U^{\eta}(\mathbf{w}_0, \mathbf{h}), \tag{8}$$

54 где ${f w}_0$ — начальное значение вектора параметров ${f w},\,{f h}$ — совокупность гиперпараметров 55 модели.

Переопределим задачу минимизации согласно определению оператора U:

$$\hat{\mathbf{h}} = \arg\max_{\mathbf{h} \in \mathbb{R}^2} \mathcal{L}_{\text{val}} \left(U^{\eta}(\mathbf{w}_0, \mathbf{h}) \right). \tag{9}$$

Схема оптимизации гиперпараметров:

- 1. Для каждого $i = \overline{0,l}$, где l количество итераций, используемых для оптимизации гиперпараметров.
- 61 2. Решим задачу (9) и получим новое значение гиперпараметров \mathbf{h}' .
- $_{62}$ 3. Положим $\mathbf{h} = \mathbf{h}'$.

52

56

57

58

59

60

64

65

67

68

69

70

71

72

73 74

75

76

77

78

79

80

81

82

83

84

85

87

3 Градиентные методы оптимизации

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

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

66 где γ — длина шага градиентного спуска.

Используем метод градиентного спуска, который зависит только от значений параметров **w** на предыдущем шаге. На каждой итерации получим следующее значение гиперпараметров:

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

4 Вычислительный эксперимент

Целью эксперимента является проверка работоспособности предложенного метода дистилляции моделей, а также анализ полученных моделей и их гиперпараметров.

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 моделейучителей 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.

На рис. 1 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.

М. Горпинич и др.

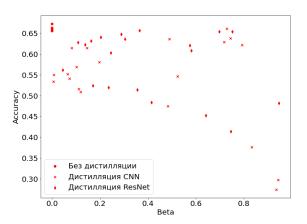


Рис. 1 График зависимости точности от β

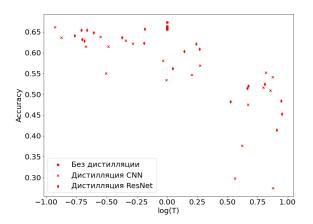


Рис. 2 График зависимости точности от температуры

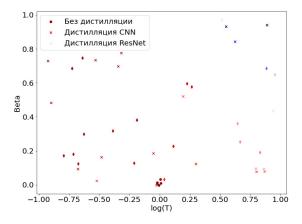


Рис. 3 График зависимости β от температуры с выделенной цветом ассигасу

На рис. 2 изображена зависимость точности от T. Для изображения значений температуры используется логарифмическая шкала. По графику видно, что значение температу-

ры уменьшается при увеличении логарифма температуры, но при значениях логарифма
 от 0.5 до 1 наблюдается резкое уменьшение точности.

⁹⁴ На рис. 3 изображена зависимость β от величины T с выделенной цветом ассигасу. ⁹⁵ Заметим, что точки с большим значением точности в основном расположены в правом ⁹⁶ нижнем углу графика, а именно, при значениях β от 0 до 0.5 и значениях $\log(T)$ от -1 ⁹⁷ до 0. Наоборот, точки с низким значением точности, расположены в правом верхнем углу ⁹⁸ графика.

Таблица 1 Результаты эксперимента

В таблице 1 приведены результаты эксперимента.

Рис. 4 График зависимости β от количества итераций дистилляции

Рис. 5 График зависимости температуры от количества итераций дистилляции

100 На рис. 4 изображена зависимость β от количества итераций дистилляции. 101 На рис. 5 изображена зависимость T от количества итераций дистилляции.

102 4.1 Название параграфа

Разделы и параграфы, за исключением списков литературы, нумеруются.

104 5 Заключение

103

105 Желательно, чтобы этот раздел был, причём он не должен дословно повторять ан-106 нотацию. Обычно здесь отмечают, каких результатов удалось добиться, какие проблемы 107 остались открытыми.

108 Литература

- 109 [1] Hinton Geoffrey E., Vinyals Oriol, Dean Jeffrey. Distilling the knowledge in a neural network //
 110 CoRR, 2015. Vol. abs/1503.02531. URL: http://arxiv.org/abs/1503.02531.
- 111 [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.
- 114 [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.
- 117 [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.
- 119 [5] Maclaurin Dougal, Duvenaud David, Adams Ryan P. Gradient-based hyperparameter optimization 120 through reversible learning // CoRR, 2015. Vol. abs/1502.03492. URL: http://arxiv.org/abs/ 1502.03492.
- 122 [6] Andrychowicz Marcin, Denil Misha, Colmenarejo Sergio Gomez, Hoffman Matthew W.,
 123 Pfau David et al. Learning to learn by gradient descent by gradient descent // CoRR, 2016.
 124 Vol. abs/1606.04474. URL: http://arxiv.org/abs/1606.04474.

М. Горпинич и др.

125 [7] Krizhevsky Alex et al. Learning multiple layers of features from tiny images, 2009.

Received Received