

Поиск согласованных нейросетевых моделей в задаче мультидоменного обучения

К. Д. Яковлев¹ О. Ю. Бахтеев^{1,2} В. В. Стрижов^{1,2}
{iakovlev.kd, bakhteev, strijov}@phystech.edu

¹Москва, Московский физико-технический институт

²Москва, Вычислительный центр им. А.А. Дородницына ФИЦ ИУ РАН

2023

Цель исследования

Цель

Предложить градиентный метод оптимизации гиперпараметров с линейной по количеству параметров и гиперпараметров сложностью итерации и затратами памяти.

Проблема

Существующие методы не гарантируют выполнение следующих условий одновременно: 1) онлайн оптимизация, 2) отсутствие смещения из-за короткого горизонта, 3) линейная сложность итерации и затраты памяти.

Метод решения

Предлагаемый метод основан на агрегации жадных гиперградиентов без дополнительных вычислительных затрат.

Агрегация жадных гиперградиентов

Пусть задано $\gamma \in (0, 1)$. Тогда аппроксимация гиперградиента запишется как:

$$\hat{\nabla}_{\alpha} = \frac{\partial}{\partial \alpha} \mathcal{L}_2(\mathbf{w}_T(\alpha), \alpha) + \sum_{t=1}^T \mathbf{B}_t \frac{\partial \mathcal{L}_2(\mathbf{w}_t, \alpha)}{\partial \mathbf{w}_t} \gamma^{T-t}.$$

| | IFT | RMAD | DrMAD | TruncBP | Proposed |
|--------------------|-----|------|-------|---------|----------|
| Онлайн оптимизация | ✗ | ✓ | ✗ | ✓ | ✓ |
| Длинный горизонт | ✓ | ✓ | ✓ | ✗ | ✓ |
| линейная сложность | ✓ | ✗ | ✓ | ✓ | ✓ |

Постановка задачи оптимизации гиперпараметров

- ▶ Пусть задан вектор параметров модели $\mathbf{w} \in \mathbb{R}^P$ и вектор гиперпараметров $\alpha \in \mathbb{R}^h$. задача оптимизации:

$$\begin{aligned} \alpha^* &= \arg \min_{\alpha} \mathcal{L}_2(\mathbf{w}^*, \alpha), \\ \text{s.t. } \mathbf{w}^* &= \arg \min_{\mathbf{w}} \mathcal{L}_1(\mathbf{w}, \alpha). \end{aligned}$$

- ▶ Пусть внутренняя задача решается с помощью оптимизатора $\Phi(., .)$:

$$\mathbf{w}_{t+1}(\alpha) = \Phi(\mathbf{w}_t, \alpha), \quad t = \overline{1, T}.$$

- ▶ Гиперградиент запишется как:

$$\begin{aligned} \nabla_{\alpha} \mathcal{L}_2(\mathbf{w}_T(\alpha), \alpha) &= \frac{\partial}{\partial \alpha} \mathcal{L}_2(\mathbf{w}_T(\alpha), \alpha) + \sum_{t=1}^T \mathbf{B}_t \mathbf{A}_{t+1} \dots \mathbf{A}_T \frac{\partial \mathcal{L}_2(\mathbf{w}_T(\alpha), \alpha)}{\partial \mathbf{w}}, \\ \mathbf{B}_t &= \frac{\partial \Phi(\mathbf{w}_{t-1}, \alpha)}{\partial \alpha}, \quad \mathbf{A}_t = \frac{\partial \Phi(\mathbf{w}_{t-1}, \alpha)}{\partial \mathbf{w}_{t-1}}. \end{aligned}$$

Аппроксимация гиперградиента

- Пусть задано $\gamma \in (0, 1)$. Тогда аппроксимация гиперградиента запишется как:

$$\hat{\nabla}_{\alpha} = \frac{\partial}{\partial \alpha} \mathcal{L}_2(\mathbf{w}_T(\alpha), \alpha) + \sum_{t=1}^T \mathbf{B}_t \frac{\partial \mathcal{L}_2(\mathbf{w}_t, \alpha)}{\partial \mathbf{w}_t} \gamma^{T-t}.$$

- предположения:

1. $\mathcal{L}_1(\cdot, \alpha)$, $\mathcal{L}_2(\cdot, \alpha)$ являются L -гладкими и μ -сильно выпуклыми.
2. $\frac{\partial^2 \mathcal{L}_1(\cdot, \alpha)}{\partial \mathbf{w} \partial \mathbf{w}^\top}$ является H_w -липшицева.
3. $1 - \eta L \leq \gamma 1 - \eta \mu$
4. $\left\| \frac{\partial \mathcal{L}_1(\mathbf{w}, \alpha)}{\partial \alpha \partial \mathbf{w}^\top} \right\| \leq B$.
5. $\frac{\partial^2 \mathcal{L}_1(\cdot, \alpha)}{\partial \alpha \partial \mathbf{w}^\top}$ является M_b -липшицевой.
6. $\left(\frac{\partial^2 \mathcal{L}_1(\cdot, \alpha)}{\partial \alpha \partial \mathbf{w}^\top} \right)^\top \left(\frac{\partial^2 \mathcal{L}_1(\cdot, \alpha)}{\partial \alpha \partial \mathbf{w}^\top} \right) \succeq \kappa \mathbf{I}$.

Асимптотическая несмещенность гиперградиента

Теорема (Яковлев, 2023)

Пусть выполнены предположения (1-6). Тогда:

$$\|\hat{\nabla}_{\alpha} - \nabla_{\alpha}\|_2 \leq \frac{2LB\|\mathbf{w}_0 - \mathbf{w}_*\|\sqrt{1-\eta\mu}^T}{\sqrt{1-\eta\mu}^{-1} - 1} + B\left\|\frac{\partial\mathcal{L}_2(\mathbf{w}_T, \alpha)}{\partial\mathbf{w}}\right\| \cdot \left[\frac{1}{\eta}\left(\frac{1}{\mu} - \frac{1}{L} + \frac{1}{L}(1-\eta\mu)^T\right) + 2\eta H_w((T-1)\sqrt{1-\eta\mu}^T - \frac{\sqrt{1-\eta\mu}^{T-1} - (1-\eta\mu)^T}{\sqrt{1-\eta\mu}^{-1} - 1}) \right].$$

Теорема (Яковлев, 2023)

Пусть $\mathcal{L}_2 = \mathcal{L}_2(\mathbf{w})$. Тогда найдется $c > 0$:

$$\hat{\nabla}_{\alpha}^{\top} \nabla_{\alpha} \geq c \|\nabla_{\mathbf{w}} \mathcal{L}_2(\mathbf{w}_T, \alpha)\|_2^2.$$

То есть, выполнено достаточное условие спуска.

Постановка вычислительного эксперимента

- ▶ Цель – сравнение качества предложенного подхода с существующими методами подсчета гиперградиента.
- ▶ Эксперимент проводится на задаче очистки обучающей выборки. Приводится точность предсказания на отложенной выборке.
- ▶ Сравниваются следующие методы: DrMAD, IFT, Truncated Backpropagation.

| Method | Valid. Acc. | #JVPs |
|---------------------------------------|--------------|---------------------|
| Truncated backpropagation (Lukethina) | 72.5 | 1 (1) |
| DrMAD | 69.8 | 99 ($2T - 1$) |
| IFT(9, 5) | 70.3 | 50 ($((N + 1)K)$) |
| IFT(4, 10) | 70.7 | 50 ($((N + 1)K)$) |
| Proposed ($\gamma = 0.99$) | 73.5* | 50 (T) |

Из таблицы видно, что предложенный метод превосходит существующие методы оптимизации гиперпараметров в терминах точности предсказаний на отложенной выборке, имея сопоставимые вычислительные затраты.

- ▶ Рассмотрена задача оптимизации гиперпараметров.
- ▶ Предложен метод оптимизации гиперпараметров, удовлетворяющий одновременно трем условиям: 1) онлайн оптимизация, 2) отсутствие смещения из-за короткого горизонта, 3) линейная сложность итерации и затраты памяти.
- ▶ Продемонстрирована работоспособность предлагаемого решения.