

Быстрая оптимизация мультизадачных моделей

Филатов Андрей

Московский физико-технический институт
Кафедра интеллектуальных систем

Апрель, 2021

Постановка задачи

Задача оптимизации

Рассмотрим T функции (задач) \mathcal{L}^t . Цель заключается в их одновременной оптимизации, то есть:

$$\min_{\substack{\boldsymbol{\theta}^{sh}, \\ \boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^T}} \mathcal{L}(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^T) = \min_{\substack{\boldsymbol{\theta}^{sh}, \\ \boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^T}} \left(\mathcal{L}^1(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^1), \dots, \mathcal{L}^T(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^T) \right)^\top \quad (1)$$

где $\boldsymbol{\theta}^{sh}$ - параметры общие параметры, а $\boldsymbol{\theta}^t$ отдельные параметры для каждой задачи

Обозначения

Парето доминирование

Считаем, что точка θ_1 доминируется точкой θ_2 если:

$$\forall i = 1 \dots T \quad \mathcal{L}^i(\theta_2) \leq \mathcal{L}^i(\theta_1)$$

$$\exists j : \mathcal{L}^j(\theta_2) < \mathcal{L}^j(\theta_1)$$

Парето оптимальность

Точка $\theta = [\theta^{sh}, \theta^t]$ для задачи (1) Парето оптимальная тогда и только тогда, когда она допустимая и не существует другой точки, которая Парето доминирует θ :

- ▶ Для любой допустимой $\hat{\theta}$: $\mathcal{L}^i(\theta) \leq \mathcal{L}^i(\hat{\theta})$ для всех $i \in \{1, \dots, T\}$
- ▶ $\mathcal{L}(\theta^{sh}, \theta^1, \dots, \theta^T) \neq \mathcal{L}(\hat{\theta}^{sh}, \hat{\theta}^1, \dots, \hat{\theta}^T)$

Обозначения

Парето станционарность

Точка $\boldsymbol{\theta}^* = [\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^t]$ Парето станционарная, если существует набор коэффициентов $\alpha^1, \dots, \alpha^T \geq 0$, таких что $\sum_{t=1}^T \alpha^t = 1$ и $\sum_{t=1}^T \alpha^t \nabla_{\boldsymbol{\theta}^{sh}} \hat{\mathcal{L}}^t(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^t) = 0$ и для всех задач t , $\nabla_{\boldsymbol{\theta}^t} \hat{\mathcal{L}}^t(\boldsymbol{\theta}^{sh}, \boldsymbol{\theta}^t) = 0$

Методы решения

- ▶ Взвешивание: сводим мультикритериальную оптимизацию к одномерной оптимизации следующим образом:

$$L(x) = \sum_{t=1}^T w_t \mathcal{L}^t(x)$$

- ▶ Методы нулевого порядка: эволюционные алгоритмы и мультикритериальная байесовская оптимизация
- ▶ Методы более старшего порядка: градиентный спуск, метод Ньютона

Методы первого порядка

► MGDA

$$\min_{\alpha^1, \dots, \alpha^T} \left\{ \left\| \sum_{t=1}^T \alpha^t \nabla_{\mathbf{x}} \hat{\mathcal{L}}^t(\mathbf{x}) \right\|_2^2 \mid \sum_{t=1}^T \alpha^t = 1, \alpha^t \geq 0 \quad \forall t \right\} \quad (2)$$

$$d^* = \sum_{t=1}^T \tilde{\alpha}^t \nabla_{\mathbf{x}} \hat{\mathcal{L}}^t$$

► Минмакс подход

$$\min_u \max_t \left(\frac{\partial \mathcal{L}}{\partial \mathbf{x}} d \right)_t + g^t(d) \quad (3)$$

Рассматривая $g^t(d) = \|d\|^2$ получаем двойственную к MGDA.

Рассматривая $g^t(d) = \frac{1}{2} d^T H^t d$ получаем метод Ньютона (H^t гессиан \mathcal{L}^t).

Модификации

- ▶ PCGrad (Tianhe Yu et al. 2020)
Предложил преобразование искомых градиентов, которое позволяет получить направление минимизирующие всех функции.
- ▶ EDM (Katrutsa et al. 2020)
Добавил нормирование градиентов, чтобы сделать метод более устойчивым в случае, когда градиенты разных задач сильно различаются по норме.

Линейный поиск

Пусть получено d — направление убывания всех функций.
Необходимо найти шаг η , чтобы:

$$\forall t = 1 \dots T : \mathcal{L}^t(x - \eta d) < \mathcal{L}^t(x) \quad (4)$$

Теорема

Если это условие (4) будет выполнено на каждой итерации, то каждый частичный предел последовательности, созданной градиентным спуском, — Парето стационарная точка.

Линейный поиск

Armijo rule

На каждом шаге градиентного спуска нам найти шаг η , чтобы выполнялось следующее Armijo rule:

$$\mathcal{L}^t(\boldsymbol{\theta}^{sh} - \eta d_{sh}, \boldsymbol{\theta}^t - \eta \nabla_{\boldsymbol{\theta}^t} \mathcal{L}^t) \leq \mathcal{L}^t - \eta \beta \left\| \frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}^t} \right\|^2 - \eta \beta \left(\frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}^{sh}} \right)^\top d_{sh} \quad (5)$$

Мотивация

- ▶ Получив направление для градиентного спуска нужно выбрать шаг;
- ▶ Теоретически обосновано лишь использование линейного поиска для нахождения шага;
- ▶ Линейный поиск неэффективен на практике из-за высокой вычислительной стоимости из-за дополнительных forward и backward проходов;

Цель

Предложить вычислительно эффективный метод, сохранив теоретические результаты.

Линейный поиск

Новое правило Армихо

На каждом шаге градиентного спуска нам найти шаг η , чтобы выполнялось следующее правило Армихо:

$$\mathcal{L}^t(\mathbf{Z} - \eta d_z, \boldsymbol{\theta}^t - \eta \nabla_{\boldsymbol{\theta}^t} \mathcal{L}^t) \leq \mathcal{L}^t - \eta \beta \left\| \frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}^t} \right\|^2 - \eta \beta \left(\frac{\partial \mathcal{L}^t}{\partial \mathbf{Z}} \right)^\top d_z \quad (6)$$

Теорема (Филатов 2021)

Каждый частичный предел последовательности, сгенерированной в результате оптимизации с правилом Армихо (6), является Парето стационарной точкой.

Классический и быстрый алгоритмы

Algorithm 1: Backtracking(γ, l_{ub})

Ensure: Learning rate $\eta = l_{ub}/\gamma$

```
1: repeat
2:    $\eta \leftarrow \gamma \cdot \eta$ 
3:    $\tilde{\theta}^{sh} \leftarrow \theta^{sh} - \eta \cdot d_{sh}$ 
4:   for  $t \leftarrow 1$  to  $T$  do
5:      $\tilde{\theta}^t \leftarrow \theta^t - \eta \cdot \nabla_{\theta^t} L^t$ 
6:   end for
7: until Armijo condition 5
8: for  $t \leftarrow 1$  to  $T$  do
9:    $\theta_{new}^t \leftarrow \tilde{\theta}^t$ 
10: end for
11:  $\theta_{new}^{sh} \leftarrow \tilde{\theta}^{sh}$ 
```

Algorithm 2: Fast backtracking(γ, l_{ub})

Ensure: Learning rate $\eta = l_{ub}/\gamma$

```
1: repeat
2:    $\eta \leftarrow \gamma \cdot \eta$ 
3:    $z \leftarrow z - \eta \cdot d_z$ 
4:   for  $t \leftarrow 1$  to  $T$  do
5:      $\tilde{\theta}^t \leftarrow \theta^t - \eta \cdot \nabla_{\theta^t} L^t$ 
6:   end for
7: until Armijo condition 6
8: for  $t \leftarrow 1$  to  $T$  do
9:    $\theta_{new}^t \leftarrow \tilde{\theta}^t$ 
10: end for
11:  $\theta_{new}^{sh} \leftarrow \theta^{sh} - \eta \cdot \frac{\partial \theta^{sh}}{\partial z} d_z$ 
```

Эксперименты

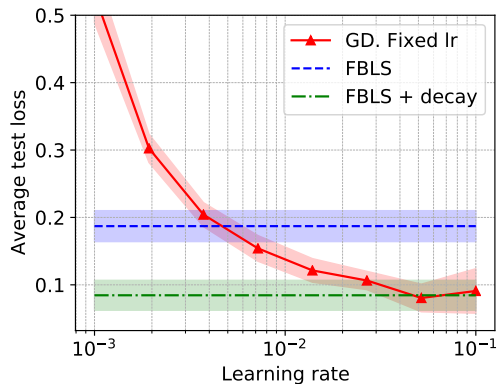
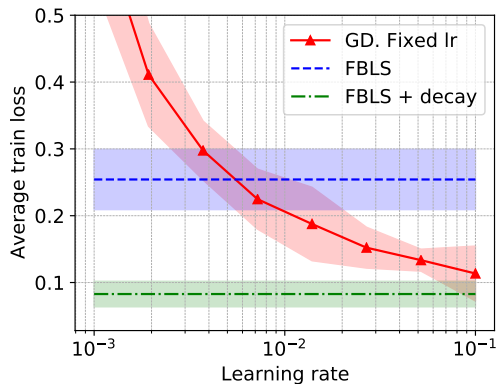


Рис.: Сравнение работы предложенного метода против градиентного спуска с фиксированным шагом после 50 эпох на датасете MultiMNIST

Эксперименты

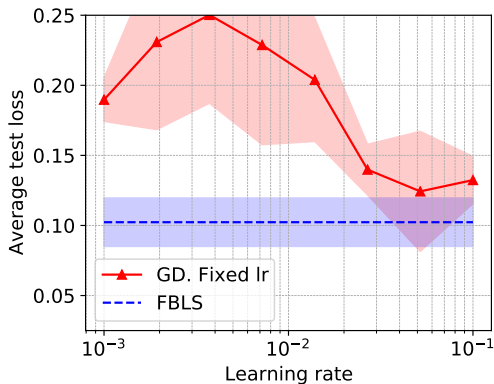
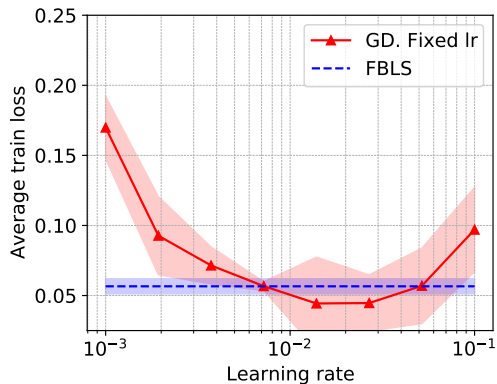


Рис.: Сравнение работы предложенного метода против градиентного спуска с фиксированным шагом после 50 эпох на датасете CIFAR-10

Эксперименты

	MNIST ↓	CIFAR-10 ↓	Cityscapes ↓
Fast backtracking	1.05 (143)	0.15 (85)	1.28 (76800)
Backtracking	1.37 (195)	1.18 (650)	-
Classical SGD	1.0 (143)	1.0 (550)	1.0 (60000)
Latent space SGD	0.95 (136)	0.14 (80)	-

Таблица: Сравнение времени работы (секунд на эпоху), меньше лучше

Результаты, выносимые на защиту

- ▶ Предложен метод оптимизации мультизадачных модели;
- ▶ Подтверждена теоретическая сходимость предложенного метода;
- ▶ Проверена эффективность метода на прикладных задачах при различных условиях.

- ▶ В рамках работы не были исследованы другие методы линейного поиска. Однако, предложенные выводы не ограничены backtracking line search.
- ▶ В рамках работы не были исследованы современные методы оптимизации. Рассмотрение методов более высокого порядка и методов редукции дисперсии будут рассмотрены в дальнейшем.