

# Быстрая оптимизация мультизадачных моделей

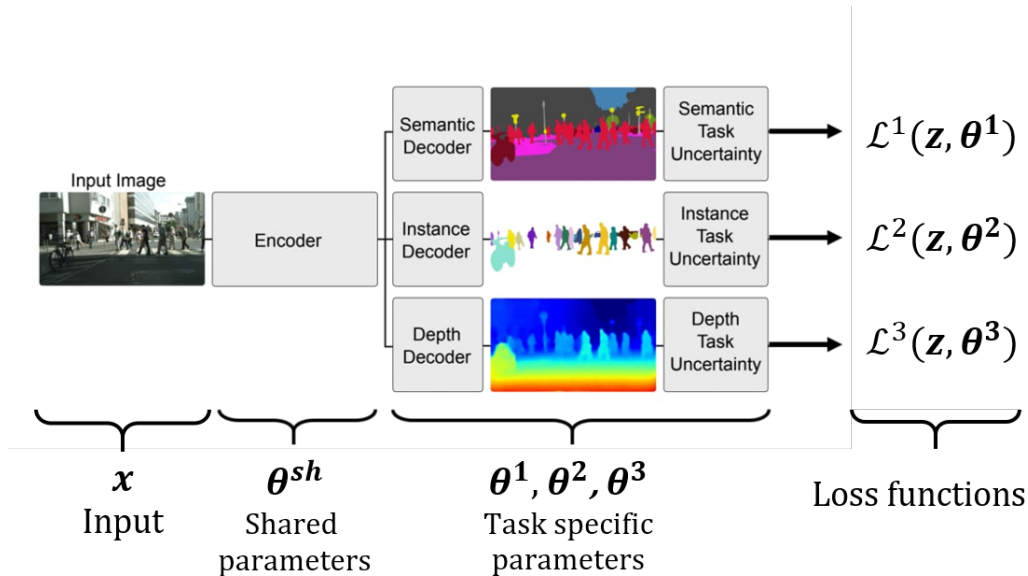
Филатов Андрей

Московский физико-технический институт  
Кафедра интеллектуальных систем

Научный руководитель: д.ф.-м.н. Стрижов В.В.  
Консультант: Меркулов Д.М.

9 Июня, 2021

# Мультизадачные модели



## Многокритериальная оптимизация

Заданы  $T$  функции потерь (задач)  $\mathcal{L}^t$ . Требуется их одновременная оптимизации:

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

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

# Методы решения задачи мультикритериальной оптимизации

1. Взвешивание<sup>1</sup>: сводим мультикритериальную оптимизацию к одномерной оптимизации следующим образом:

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{t=1}^T w_t \mathcal{L}^t(\boldsymbol{\theta}).$$

2. Методы нулевого порядка: эволюционные алгоритмы и мультикритериальная байесовская оптимизация.<sup>2</sup>
3. Градиентные методы<sup>3</sup>: градиентный спуск, метод Ньютона.

---

<sup>1</sup>Chen и др., “Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks”.

<sup>2</sup>Deb и др., “A fast elitist non-dominated sorting genetic algorithm for multi-objective optimization: NSGA-II”.

<sup>3</sup>Fliege и Svaiter, “Steepest descent methods for multicriteria optimization”.

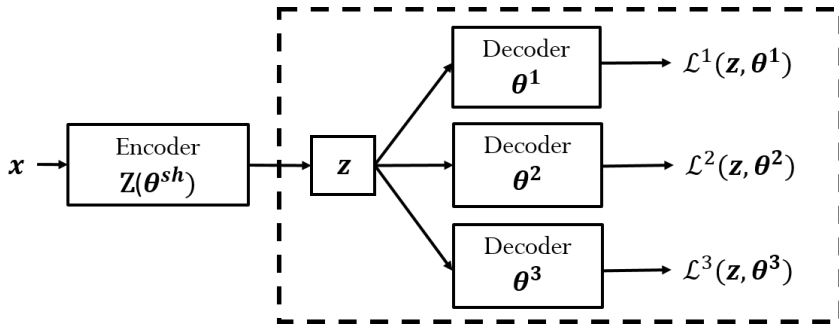
## Мотивация

- ▶ Для градиентных методов оптимизации мультизадачных моделей теоретически обосновано лишь использование линейного поиска для нахождения шага;
- ▶ Линейный поиск неэффективен на практике из высокой вычислительной стоимости;

## Цель

Создать вычислительно эффективный метод линейного поиска.

## Быстрый линейный поиск



В алгоритме быстрого линейного поиска мы оптимизируем в скрытом пространстве  $Z$ .

## Градиентные методы

- Multiple-gradient descent algorithm (MGDA)<sup>4</sup>

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

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

- Min-max подход<sup>5</sup>

$$\min_{\mathbf{d}} \max_t \left( \frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}} \mathbf{d} \right)_t + g^t(\mathbf{d})$$

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

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

---

<sup>4</sup>Désidéri, “Mgda variants for multi-objective optimization”.

<sup>5</sup>Fliege и Svaiter, “Steepest descent methods for multicriteria optimization”.

## Алгоритм линейного поиска

Пусть получено  $\mathbf{d}$  — направление убывания всех функций:  $\forall t \nabla_{\boldsymbol{\theta}} \mathcal{L}^t \mathbf{d} < 0$ .  
Необходимо найти шаг  $\eta$ , чтобы:

$$\mathcal{L}^t(\boldsymbol{\theta} - \eta \mathbf{d}) < \mathcal{L}^t(\boldsymbol{\theta}), \forall t \in \{1 \dots T\}. \quad (\star)$$

### Теорема (Fliege 2000)<sup>6</sup>

Если условие  $(\star)$  будет выполнено на каждой итерации, то для любой сходящейся подпоследовательности  $\{\boldsymbol{\theta}_{k_j}\}_{j=1}^{\infty} : \lim_{j \rightarrow \infty} \boldsymbol{\theta}_{k_j} = \hat{\boldsymbol{\theta}}$ , созданной градиентным спуском, предел этой последовательности  $\hat{\boldsymbol{\theta}}$  — Парето стационарная точка.

---

<sup>6</sup>Fliege и Svaiter, “Steepest descent methods for multicriteria optimization”.



# Алгоритм линейного поиска

## Правило Армихо

На каждом шаге градиентного спуска нам найти шаг  $\eta$ , чтобы выполнялось следующее правило Армихо  $\forall t \in \{1 \dots T\}$ :

$$\mathcal{L}^t(\boldsymbol{\theta}^{sh} - \eta \mathbf{d}_{sh}, \boldsymbol{\theta}^t - \eta \nabla_{\boldsymbol{\theta}^t} \mathcal{L}^t) \leq \mathcal{L}^t - \eta \beta \left\| \frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}^t} \right\|^2 - \eta \beta \left( \frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}^{sh}} \right)^\top \mathbf{d}_{sh}.$$

## Модифицированное правило Армихо

На каждом шаге градиентного спуска нам найти шаг  $\eta$ , чтобы выполнялось следующее правило Армихо  $\forall t \in \{1 \dots T\}$ :

$$\mathcal{L}^t(\mathbf{z} - \eta \mathbf{d}_z, \boldsymbol{\theta}^t - \eta \nabla_{\boldsymbol{\theta}^t} \mathcal{L}^t) \leq \mathcal{L}^t - \eta \beta \left\| \frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}^t} \right\|^2 - \eta \beta \left( \frac{\partial \mathcal{L}^t}{\partial \mathbf{z}} \right)^\top \mathbf{d}_z.$$

# Основной результат

## Модифицированное правило Армихо

На каждом шаге градиентного спуска нам найти шаг  $\eta$ , чтобы выполнялось следующее правило Армихо :

$$\mathcal{L}^t(\mathbf{z} - \eta \mathbf{d}_z, \boldsymbol{\theta}^t - \eta \nabla_{\boldsymbol{\theta}^t} \mathcal{L}^t) \leq \mathcal{L}^t - \eta \beta \left\| \frac{\partial \mathcal{L}^t}{\partial \boldsymbol{\theta}^t} \right\|^2 - \eta \beta \left( \frac{\partial \mathcal{L}^t}{\partial \mathbf{z}} \right)^\top \mathbf{d}_z$$

## Теорема (Филатов 2021)

Каждый частичный предел последовательности, созданной градиентным спуском с модифицированным правилом Армихо, является Парето стационарной точкой.

## Базовый<sup>7</sup> и быстрый алгоритмы

---

### Algorithm 1: Backtracking( $\gamma, lr_{ub}$ )

---

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

```
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:      $\tilde{\theta}^t \leftarrow \theta^t - \eta \cdot \nabla_{\theta^t} \mathcal{L}^t$ 
6:   end for
7: until правило Армихо
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( $\gamma, lr_{ub}$ )

---

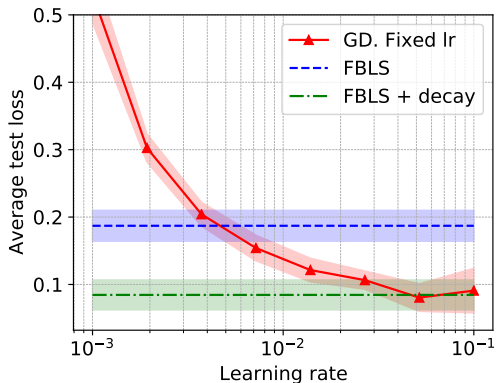
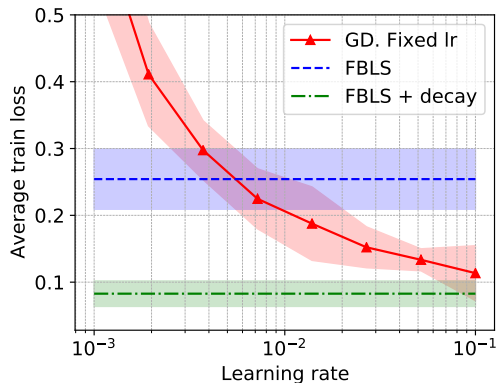
Ensure: Learning rate  $\eta = lr_{ub}/\gamma$

```
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:      $\tilde{\theta}^t \leftarrow \theta^t - \eta \cdot \nabla_{\theta^t} \mathcal{L}^t$ 
6:   end for
7: until правило Армихо
8: for  $t \leftarrow 1$  to  $T$  do
9:    $\theta_{new}^t \leftarrow \tilde{\theta}^t$ 
10: end for
11:  $\theta_{new}^{sh} \leftarrow \theta^{sh} - \eta \cdot \frac{\partial \theta^{sh}}{\partial z} d_z$ 
```

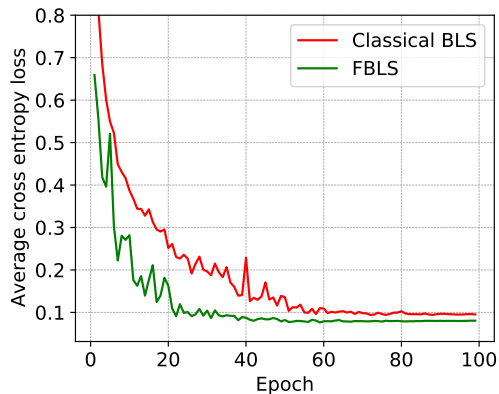
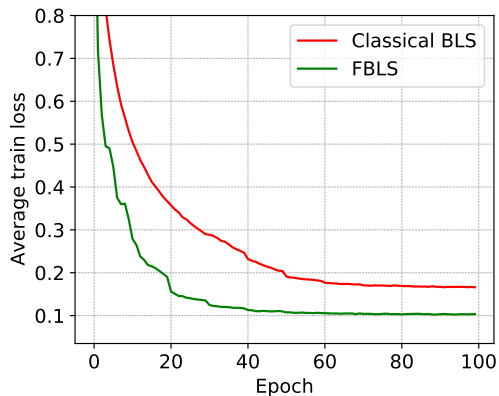
---

<sup>7</sup>Armijo, "Minimization of functions having Lipschitz continuous first partial derivatives".

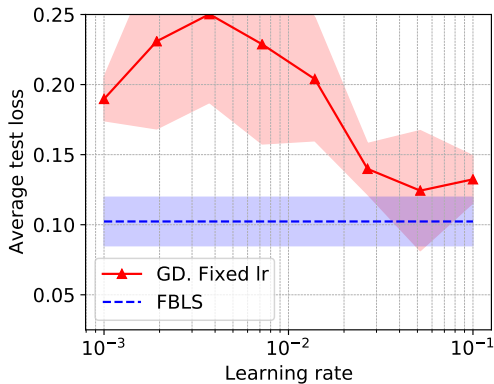
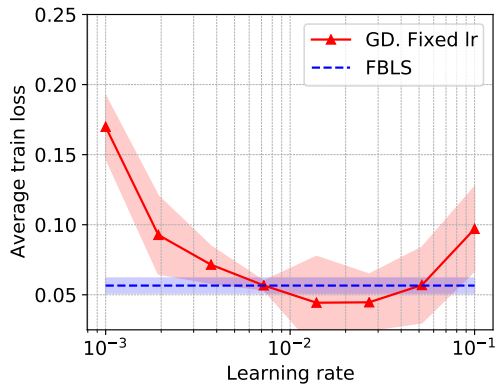
## Сравнение предложенного метода с градиентным спуском на MultiMNIST



## Сравнение предложенного метода с базовым линейным поиском на MultiMNIST



## Сравнение предложенного метода с градиентным спуском на CIFAR-10



## Сравнение времени работы алгоритмов

	MNIST ↓	CIFAR-10 ↓	Cityscapes ↓
FBLS (Ours)	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)	-

## Результаты, выносимые на защиту

1. Предложен алгоритм быстрого линейного поиска для мультизадачных моделей.
2. Подтверждена теоретическая сходимость быстрого линейного поиска к Парето стационарной точке.
3. Проверена практическая эффективность метода на задачах MultiMNIST, CIFAR-10, Cityscapes.