

# Federated Multi-Task Learning

Арте́м Ага́рков

МФТИ

12 декабря 2023

## 1 Постановка задачи

- Федеративное обучение
- Multi-task обучение
- Федеративное multi-task обучение (FMTL)

## 2 FMTL

- Общая постановка multi-task обучения
- Некоторые наблюдения
- Вывод двойственной задачи

## 3 Двойственная задача

## 4 Локальные квадратичные подзадачи

## 5 Обновление $\Omega$

- Структура с кластерами

## 6 Алгоритм

## 7 Эксперименты

При федеративном обучении цель состоит в том, чтобы обучить модель на основе данных, которые находятся на  $m$  распределенных нодах и были сгенерированы ими. В качестве примера рассмотрим решение задачи классификации активности пользователей мобильных телефонов в сотовой сети на основе их индивидуальных данных датчиков, текстов или изображений.

Каждый узел (телефон),  $t \in [m]$ , может генерировать данные по своему распределению, поэтому естественно подбирать отдельные модели,  $[w_1, w_2, \dots, w_m]$ , для распределенных данных - по одной для каждого локального набора данных.

Multi-task обучение - это область машинного обучения, в которой одновременно решаются несколько задач обучения, при этом используются общие черты и различия между задачами. Это может привести к повышению эффективности обучения и точности предсказания для моделей, специфичных для конкретной задачи, по сравнению с обучением моделей по отдельности.

Пример для задачи классификации: спам-фильтр, который можно рассматривать как отдельные, но связанные между собой задачи классификации для разных пользователей.

Предположим, что разные люди имеют различные распределения признаков, которые отличают спам от допустимых писем, например, англоговорящий человек может решить, что все письма на русском языке являются спамом, а русскоговорящий - нет. Тем не менее, в этой задаче классификации есть определенная общность для разных пользователей, например, одним из общих признаков может быть текст, связанный с переводом денег.

# Федеративное multi-task обучение (FMTL)

Вернемся к примеру федеративного обучения для задачи изучения активности пользователей телефонов. Между моделями часто существует структура (например, люди могут вести себя одинаково при использовании своих телефонов). Тогда моделирование этой структуры с помощью multi-task обучения является естественной стратегией для **повышения производительности и увеличения размера выборки** для каждого узла.

# Общая постановка multi-task обучения

Данные  $X_t \in \mathbb{R}^{d \times n_t}$  для каждой из  $m$  нод в multi-task подходе обучаем параметры  $w_t \in \mathbb{R}^d$  на данных для каждой ноды с произвольной выпуклой функцией потерь  $\ell_t$

$$\min_{\mathbf{W}, \Omega} \left( \sum_{t=1}^m \sum_{i=1}^{n_t} \ell_t(\mathbf{w}_t^T \mathbf{x}_t^i, y_t^i) + \mathcal{R}(\mathbf{W}, \Omega) \right) \quad (1)$$

$\mathbf{W} := [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m] \in \mathbb{R}^{d \times m}$  ( $t$ -й столбец - параметры для  $t$ -й задачи. Матрица  $\Omega \in \mathbb{R}^{m \times m}$  отвечает за взаимосвязь между задачами. Может быть задана изначально или найдена в процессе обучения. Существуют различные подходы к определению  $\mathcal{R}$  и, соответственно, задачи MTL. Наиболее распространенная бивыпуклая формулировка:

$$\mathcal{R}(\mathbf{W}, \Omega) = \lambda_1 \text{tr}(\mathbf{W} \Omega \mathbf{W}^T) + \lambda_2 \|\mathbf{W}\|_F^2 \quad (2)$$

# Некоторые наблюдения

## Наблюдение 1

В общем случае, (1) не является выпуклой функцией по  $\mathbf{W}, \Omega$ . И даже в случае, когда (1) выпуклая, решение для  $\mathbf{W}, \Omega$  одновременно может быть сложным.

## Наблюдение 2

При фиксировании  $\Omega$  обновление  $\mathbf{W}$  зависит как от данных  $\mathbf{X}$ , которые распределены по узлам, так и от структуры  $\Omega$ , которая известна централизованно.

## Наблюдение 3

При фиксированном  $\mathbf{W}$  оптимизация для  $\Omega$  зависит только от  $\mathbf{W}$ , но не от данных  $\mathbf{X}$ .



Исходя из этих наблюдений, естественно предложить подход попеременной оптимизации для решения задачи (1), при котором на каждой итерации мы фиксируем либо  $W$ , либо  $\Omega$  и оптимизируем по другому параметру, чередуя их до достижения сходимости. Отметим, что решение для  $\Omega$  не зависит от данных и поэтому может быть вычислено централизованно; поэтому далее рассматриваем только задачу оптимизации для  $W$ .

# Вывод двойственной задачи

Итак, мы рассматриваем итерацию с оптимизацией  $\mathbf{W}$ :

$$\min_{\mathbf{W}} \left( \sum_{t=1}^m \sum_{i=1}^{n_t} \ell_t(\mathbf{w}_t^T \mathbf{x}_t^i, y_t^i) + \mathcal{R}(\mathbf{W}, \Omega) \right)$$

Более абстрактно, можем записать исходную задачу в виде:

$$\min_{\mathbf{w}, \mathbf{v}} (f(\mathbf{v}) + g(\mathbf{w})), \text{ s.t. } \mathbf{v} = \mathbf{X}^T \mathbf{w},$$

где  $\mathbf{w} \in \mathbb{R}^{dm}$  - "вытянутый" вектор параметров,  
 $\mathbf{X} := \text{Diag}(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m) \in \mathbb{R}^{md \times n}$ . Будем искать двойственную задачу для этой абстрактной формулировки.

# Вывод двойственной задачи

Лагранжиан:  $L(\mathbf{w}, \mathbf{v}; \alpha) = f(\mathbf{v}) + g(\mathbf{w}) + \alpha^T(\mathbf{X}^T \mathbf{w} - \mathbf{v})$ .

Тогда инфимум Лагранжиана может быть преобразован так:

$$\begin{aligned}\inf_{\mathbf{w}, \mathbf{v}} L(\alpha, \mathbf{w}, \mathbf{v}) &= \inf_{\mathbf{v}} \left( f(\mathbf{v}) - \alpha^T \mathbf{v} \right) + \inf_{\mathbf{w}} \left( g(\mathbf{w}) + \alpha^T \mathbf{X}^T \mathbf{w} \right) \\ &= -\sup_{\mathbf{v}} \left( \alpha^T \mathbf{v} - f(\mathbf{v}) \right) - \sup_{\mathbf{w}} \left( -(\alpha^T \mathbf{X}^T) \mathbf{w} - g(\mathbf{w}) \right) \\ &= -f^*(\alpha) - g^*(-\mathbf{X}\alpha)\end{aligned}$$

Далее просто меняем знаки и получаем двойственную задачу оптимизации:

$$\max_{\alpha} (-f^*(\alpha) - g^*(-\mathbf{X}\alpha)) = \min_{\alpha} (f^*(\alpha) + g^*(-\mathbf{X}\alpha))$$

И еще раз поменяем знак для эстетики:

$$\min_{\alpha} (f^*(-\alpha) + g^*(\mathbf{X}\alpha))$$

# Двойственная задача

Для задачи (1) рассмотрим выведенную двойственную задачу, которая позволит перейти к распределенным подзадачам.

Пусть  $n := \sum_{t=1}^m n_t$  и  $X := \text{Diag}(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m) \in \mathbb{R}^{md \times n}$ . При фиксированной  $\Omega$ :

$$\min_{\alpha} \left( \mathcal{D}(\alpha := \sum_{t=1}^m \sum_{i=1}^{n_t} \ell_t^*(-\alpha_t^i) + \mathcal{R}^*(\mathbf{X}\alpha) \right), \quad (3)$$

где  $\ell_t^*$  и  $\mathcal{R}^*$  сопряженные двойственным функции  $\ell_t$  и  $\mathcal{R}$  соответственно, и  $\alpha_t^i$  - двойственная переменная для данных  $(\mathbf{x}_t^i, y_t^i)$

# Двойственная задача

Заметим, что  $\mathcal{R}^*$  зависит от  $\Omega$ , но для простоты обозначений не будем это указывать. Чтобы вывести распределенные подзадачи из этой глобальной двойственной задачи, сделаем предположение, описанное ниже, на регуляризатор  $\mathcal{R}$ .

## Предположение 1

Для  $\Omega$ , мы предполагаем, что существует симметричная положительно определенная матрица  $\mathbf{M} \in \mathbb{R}^{md \times md}$ , зависящая от  $\Omega$ , для которой функция  $\mathcal{R}$  является сильно выпуклой относительно  $\mathbf{M}^{-1}$ . Заметим, что это соответствует предположению, что  $\mathcal{R}^*$  будет гладкой относительно матрицы  $\mathbf{M}$ .

# Локальные квадратичные подзадачи

Для решения (1) на узлах определяем следующие локальные подзадачи, которые являются квадратичной аппроксимацией двойственной задачи (3). Эти подзадачи находят обновления  $\Delta\alpha_t \in \mathbb{R}^{n_t}$  для двойственных переменных в  $\alpha$ , и требуют только локальных данных  $X_t$ .  $t$ -я подзадача задается:

$$\begin{aligned} \min_{\Delta\alpha_t} \mathcal{G}_t^{\sigma'}(\Delta\alpha_t; \mathbf{v}_t, \alpha_t) \\ \mathcal{G}_t^{\sigma'}(\Delta\alpha_t; \mathbf{v}_t, \alpha_t) := \sum_{i=1}^{n_t} \ell_t^*(-\alpha_t^* - \Delta\alpha_t^i) + \langle \mathbf{w}_t(\alpha), \mathbf{X}_t \Delta\alpha_t \rangle + \\ + \frac{\sigma'}{2} \|\mathbf{X}_t \Delta\alpha_t\|_{M_t}^2 + c(\alpha), \end{aligned}$$

где  $c(\alpha) = \frac{1}{m} \mathcal{R}^*(\mathbf{X}\alpha)$ ,  $M_t \in \mathbb{R}^{d \times d}$  -  $t$ -й диагональный блок симметричной положительно определенной матрицы  $M$ ,  $\mathbf{v} = \mathbf{X}\alpha$ .

# Связь primal и dual переменных

После того, как мы нашли обновление  $\Delta\alpha$ , необходимо обновить параметры  $\mathbf{w}_t$  на каждой ноде. Связь primal и dual переменных найдем из условий оптимальности первого порядка.

Заметим, что если  $\mathcal{R}$  выпуклая, замкнутая, собственная функция, мы можем воспользоваться теоремой Фенхеля-Моро:  $\mathcal{R}^{**} = \mathcal{R}$ .

Перепишем по определению:

$$\mathcal{R}(\mathbf{w}) = \sup_{\alpha} \left( \mathbf{w}^T \alpha - \mathcal{R}^*(\mathbf{X}\alpha) \right).$$

Тогда если на  $\alpha^*$  достигается  $\sup$ , выполнены условия оптимальности первого порядка:

$$\nabla_{\alpha} \left( \mathbf{w}^T \alpha^* - \mathcal{R}^*(\mathbf{X}\alpha^*) \right) = 0 \Rightarrow \mathbf{w}(\alpha^*) = \nabla \mathcal{R}^*(\mathbf{X}\alpha^*)$$

Мы выяснили, как делать распределенную оптимизацию по  $W$ . Однако в нашем алгоритме еще есть шаг с централизованным обновлением  $\Omega$ .

Как мы говорили в начале, существуют разные подходы к задаче MTL (определению  $\mathcal{R}$ ). Будем рассматривать задачу MTL со структурой кластеров.



# Структура с кластерами

При такой постановке векторы весов для каждой задачи  $\mathbf{w}_t$  предполагаются "близкими" по некоторой метрике к остальным весам в кластере. То есть все веса в кластере располагаются "близко" к среднему.

Такую структуру описывает матрица  $\Omega = (I_{m \times m} - \frac{1}{m} \mathbf{1}\mathbf{1}^T)^2$ .

В случае кластеров регуляризатор имеет вид:

$$\mathcal{R}(\mathbf{W}, \Omega) = \lambda_1 \text{tr}(\mathbf{W}\Omega\mathbf{W}^T) + \lambda_2 \|\mathbf{W}\|_F^2,$$

где  $\lambda_1, \lambda_2 > 0$  параметры. Здесь считаем матрицу  $\Omega$  известной заранее, количество кластеров = 1.

Однако в случае большего числа кластеров  $\mathcal{R}$  может быть невыпуклой, поэтому производится выпуклая релаксация:

$$\mathcal{R}(\mathbf{W}, \Omega) = \lambda \text{tr}(\mathbf{W}(\eta I + \Omega)^{-1} \mathbf{W}^T), \Omega \in \mathcal{Q} = \{\mathbf{Q} | \mathbf{Q} \succeq 0, \text{tr}(\mathbf{Q}) = k, \mathbf{Q} \preceq I\},$$

где  $\lambda$  и  $\nu$  параметры регуляризации,  $k$  - количество кластеров.

---

**Algorithm 1** MOCHA: Federated Multi-Task Learning Framework

---

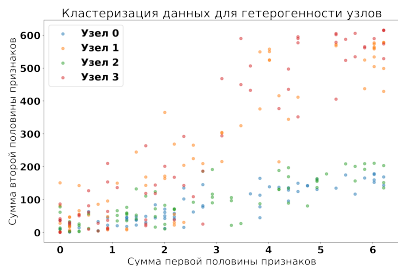
```
1: Input: Data  $\mathbf{X}_t$  from  $t = 1, \dots, m$  tasks, stored on one of  $m$  nodes, and initial matrix  $\mathbf{\Omega}_0$ 
2: Starting point  $\boldsymbol{\alpha}^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{v}^{(0)} := \mathbf{0} \in \mathbb{R}^b$ 
3: for iterations  $i = 0, 1, \dots$  do
4:   Set subproblem parameter  $\sigma'$  and number of federated iterations,  $H_i$ 
5:   for iterations  $h = 0, 1, \dots, H_i$  do
6:     for tasks  $t \in \{1, 2, \dots, m\}$  in parallel over  $m$  nodes do
7:       call local solver, returning  $\theta_t^h$ -approximate solution  $\Delta\boldsymbol{\alpha}_t$  of the local subproblem
8:       update local variables  $\boldsymbol{\alpha}_t \leftarrow \boldsymbol{\alpha}_t + \Delta\boldsymbol{\alpha}_t$ 
9:       return updates  $\Delta\mathbf{v}_t := \mathbf{X}_t\Delta\boldsymbol{\alpha}_t$ 
10:    reduce:  $\mathbf{v}_t \leftarrow \mathbf{v}_t + \Delta\mathbf{v}_t$ 
11:   Update  $\mathbf{\Omega}$  centrally based on  $\mathbf{w}(\boldsymbol{\alpha})$  for latest  $\boldsymbol{\alpha}$ 
12: Central node computes  $\mathbf{w} = \mathbf{w}(\boldsymbol{\alpha})$  based on the latest  $\boldsymbol{\alpha}$ 
13: return:  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_m]$ 
```

---

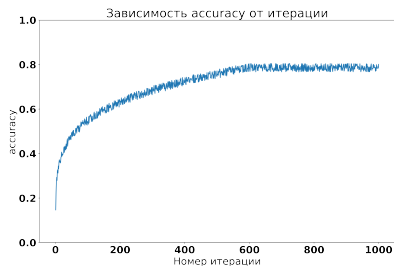
С помощью рассмотренного метода попробуем решить задачу бинарной классификации.

Датасет: phishing из LibSVM

Чтобы обеспечить гетерогенность в данных, выделим кластеры с помощью алгоритма k-means. Полученные кластеры будут нодами.



(a) Гетерогенность



(b) Accuracy

# The End