

# Анализ сходимости поверхности функции потерь сверточных нейросетевых моделей на основе Гессиана

---

Владислав Мешков, Никита Киселев, Андрей Грабовой

Московский Физико-Технический Институт

## Проблема

Число данных для обучения нейросетей, а также число параметров растет. Необходимо изучать связь между сложностью моделей и необходимым числом объектов в обучающей выборке.

## Цель

Изучить сходимость поверхности функции потерь в пространстве параметров при изменении размера выборки для сетей, представленных в виде произведения зависимых от входа матриц.

## Решение

1. Рассмотреть взаимосвязь гессиана со сходимостью функции потерь.
2. Посчитать норму гессиана для сетей, представленных в виде произведения матриц, и применить данные результаты к сверточным сетям.

## Постановка задачи

Пусть  $f_{\theta}$  — нейросеть,  $\mathbf{x}_k$  — входы,  $\mathbf{y}_k$  — one-hot метки классов,  $\mathcal{L}_k$  — функция потерь на  $k$  первых объектах.

**Изменение значения при добавлении одного объекта**

$$\mathcal{L}_{k+1}(\theta) - \mathcal{L}_k(\theta) = \frac{1}{k+1} (\ell(f_{\theta}(\mathbf{x}_{k+1}), \mathbf{y}_{k+1}) - \mathcal{L}_k(\theta))$$

### Предположение 1

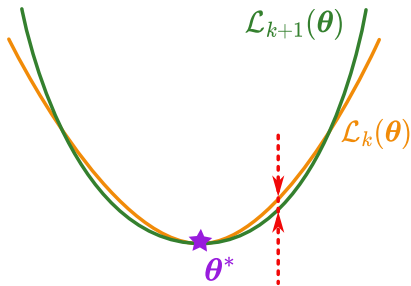
Пусть  $\theta^*$  является точкой локального минимума обеих функций  $\mathcal{L}_k(\theta)$  и  $\mathcal{L}_{k+1}(\theta)$ :

$$\nabla \mathcal{L}_k(\theta^*) = \nabla \mathcal{L}_{k+1}(\theta^*) = 0.$$

### Аппроксимация второго порядка

$$\mathbf{H}^{(k)} = \frac{1}{k} \sum_{i=1}^k \mathbf{H}_i(\theta) = \frac{1}{k} \sum_{i=1}^k \nabla_{\theta}^2 \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

$$\mathcal{L}_k(\theta) \approx \mathcal{L}_k(\theta^*) + \frac{1}{2}(\theta - \theta^*)^{\top} \mathbf{H}^{(k)}(\theta^*)(\theta - \theta^*)$$



# Связь изменения функции потерь с Гессианом

## Абсолютное изменение функции потерь

$$|\mathcal{L}_{k+1}(\boldsymbol{\theta}) - \mathcal{L}_k(\boldsymbol{\theta})| \leq \frac{2}{k+1} \max_{i=1, k+1} |\ell(f_{\boldsymbol{\theta}^*}(\mathbf{x}_i, \mathbf{y}_i))| + \\ + \frac{2}{k+1} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 \max_{i=1, k+1} \|\mathbf{H}_i(\boldsymbol{\theta}^*)\|$$

## Декомпозиция Гессиана

$$\mathbf{H}_i(\boldsymbol{\theta}) = \underbrace{\nabla_{\boldsymbol{\theta} \mathbf{z}_i} \frac{\partial^2 \ell(\mathbf{z}_i, \mathbf{y}_i)}{\partial \mathbf{z}_i^2} \nabla_{\boldsymbol{\theta} \mathbf{z}_i^\top}}_{\mathbf{H}_O} + \underbrace{\sum_{k=1}^K \frac{\partial \ell(\mathbf{z}_i, \mathbf{y}_i)}{\partial z_{ik}} \nabla_{\boldsymbol{\theta}}^2 z_{ik}}_{\mathbf{H}_F}$$

## Аппроксимация Гессиана

- Аппроксимируем Гессиан, пренебрегая  $\mathbf{H}_F$ . В задаче  $K$ -классовой классификации  $\|\mathbf{H}_F\| \ll \|\mathbf{H}_O\|$
- Тогда можно оценить  $\|\mathbf{H}\| \approx \|\mathbf{H}_O\|$ .

# Нейросети представимые в виде произведения матриц

---

Рассмотрим нейросети специального вида, а именно представимые в виде произведения матриц, возможно зависящих от входа.

Пусть нейросеть  $f_{\theta}(\mathbf{x}) := \mathbf{T}^{(L+1)} \mathbf{\Lambda}^{(L)} \dots \mathbf{\Lambda}^{(1)} \mathbf{T}^{(1)} \mathbf{x}$ , где

- $\mathbf{x}$  — вход
- $\mathbf{T}^{(l)}$  — матрица линейного слоя  $l$
- $\mathbf{\Lambda}^{(l)}$  — матрица *ReLU*-активации  $l$ -го слоя, зависящая от входа  $\mathbf{x}$

## Структура $\mathbf{H}_O$ компоненты Гессиана

---

Рассмотрим матрицы:

- $\mathbf{F}$  — матрица сложной структуры, зависящая только от всех  $\mathbf{T}^{(i)}, \mathbf{\Lambda}^{(i)}$  и от  $\mathbf{x}$
- $\mathbf{A} = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T$ , где  $\mathbf{p}$  — вектор вероятностей классов для  $\mathbf{x}$
- $\mathbf{Q}^{(p)} := \frac{\partial \mathbf{T}^{(p)}}{\partial \mathbf{W}^{(p)}}$ , где  $\mathbf{W}^{(p)}$  — параметры  $p$ -го слоя.

**Лемма 1**

$$\mathbf{H}_O(\boldsymbol{\theta}) = \mathbf{Q}^T \mathbf{F}^T \mathbf{A} \mathbf{F} \mathbf{Q}.$$

Данная Лемма позволяет представить норму  $\mathbf{H}_O$  компоненты Гессиана как произведение норм более простых блоков.

### Лемма 2

Пусть  $\|\mathbf{Q}^{(p)}\|_2 \leq q$ ,  $\|\mathbf{T}^{(p)}\|^2 \leq w_{\mathbf{T}}^2 \quad \forall p$ , тогда  
 $\|\mathbf{H}_O\| \leq \sqrt{2}q^2 \|\mathbf{x}\|^2 (L+1)w_{\mathbf{T}}^{2L}$ .

Оценка нормы  $\mathbf{H}_O$  как функция весов является степенной, а как функция числа слоев — показательной.

**Лемма 2** является основой для оценки нормы гессиана в дальнейшем.

Из Леммы, для того, чтобы оценить Гессиан, достаточно оценить  $\|\mathbf{Q}^{(p)}\|$  и  $\|\mathbf{T}^{(p)}\|$  одновременно для всех слоев, что и будет сделано в будущих результатах.

## Свертки как линейная операция

---

Действие свертки на  $\mathbf{x} \in \mathbb{R}^{C \times d}$  представим линейным оператором, действующим на  $\mathbf{x} \in \mathbb{R}^{C \times d}$ , причем с сохранением обозначений :  $\mathbf{T}^{(l)} * \mathbf{x} \rightarrow \mathbf{T}^{(l)} \mathbf{x}$

$$\begin{pmatrix} t_1 & t_2 & t_3 \end{pmatrix} * \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} \rightarrow \begin{pmatrix} t_1 & t_2 & t_3 & 0 & 0 \\ 0 & t_1 & t_2 & t_3 & 0 \\ 0 & 0 & t_1 & t_2 & t_3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix}$$



## Теорема 1

Пусть  $C_l \leq C$ ,  $k_i \leq k$ ,  $d_i \leq d$ ,  $|W_{i,j,k}^{(p)}|^2 \leq w^2$ , где  $C_l, k_l, d_l$  — число каналов, размер ядра и пространственный размер соответственно, тогда

$$\|\mathbf{H}_O\| \leq \sqrt{2} \|x\|^2 d^2 (L+1) (C^2 w^2 k d)^L$$

Из Теоремы видно, что оценка нормы является степенной функцией от числа каналов, весов, размера ядра и размера последовательности

### Теорема 2

Пусть  $|\mathbf{W}_{i,j,k,t}^{(p)}|^2 \leq w^2$ , где  $\mathbf{W}^{(p)}$  — веса  $p$ -го слоя свертки

$C_l, k_l, (m_l, n_l)$  — число каналов, размер ядра, пространственные размеры карты признаков соответственно на  $l$ -м слое нейросети. Пусть

$C_l \leq C, k_l \leq k, m_l \leq m, n_l \leq n$ , тогда

$$\|\mathbf{H}_O\| \leq \sqrt{2} \|x\|^2 C^2 k^2 m n (L+1) (C^2 k^2 w^2 m n)^L.$$

Результат полученный в **Теореме 1** отличается от данного лучшей оценкой

$\|\mathbf{Q}^{(p)}\|$ , это связано с различием в структуре Теплицевых матриц 1D и 2D сверток.

### Лемма 3

Пусть на месте  $l$ -й нелинейности находится max/avg пулинг, причем пусть ядро:  $k_{\text{pool}} \times k_{\text{pool}}$ ,  $\text{stride} = k_{\text{pool}}$ ,  $\text{padding} = 0$ , при этом верны все ограничения предыдущей теоремы, тогда:

$$\|\mathbf{H}_O\| \leq \sqrt{2} \|x\|^2 q^2 \frac{1}{k_{\text{pool}}^{2(L-l+2)}} (L+1)(C^2 k^2 w^2 mn)^L$$

### Лемма 4

Пусть после сверточных слоев находится полносвязная голова размера  $P$ :

$$f_{\theta}(\mathbf{x}) = \mathbf{T}^{(L+P+1)} \mathbf{\Lambda}^{(L+P)} \dots \mathbf{\Lambda}^{(L+1)} \mathbf{T}^{(L+1)} \mathbf{\Lambda}^{(L)} \dots \mathbf{\Lambda}^{(1)} \mathbf{T}^{(1)} \mathbf{x},$$

где  $\mathbf{T}^{(L+1+i)}$  — Линейный слой с параметрами  $h_i, h_{i+1}$ ,  $\mathbf{T}^{(r)}$ -2D-свертки. Пусть имеют место оценки:  $\|\mathbf{T}_{ij}^{(L+1+i)}\| \leq \tilde{w}$  и  $h_p \leq h$ . Тогда в условиях Теоремы 2 имеем

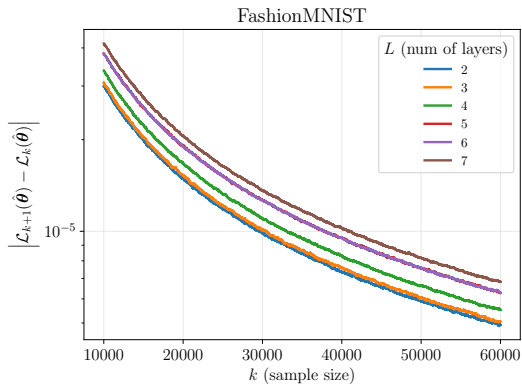
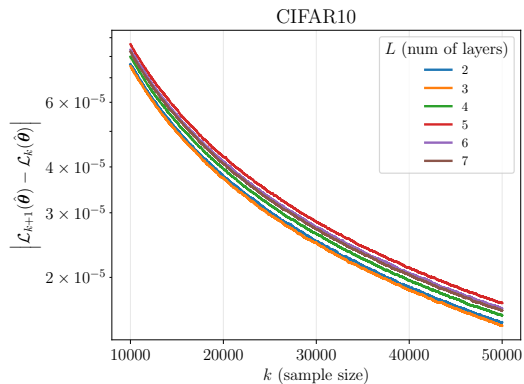
$$\|\mathbf{H}_O\| \leq \sqrt{2} \|\mathbf{x}\|^2 C^2 k^2 mn (h^2 \tilde{w}^2)^P (k^2 C^2 w^2 mn)^L \times \left( L + 1 + P \frac{h^2 \tilde{w}^2}{k^2 C^2 w^2 mn} \right). \quad 10/15$$

## Постановка эксперимента

---

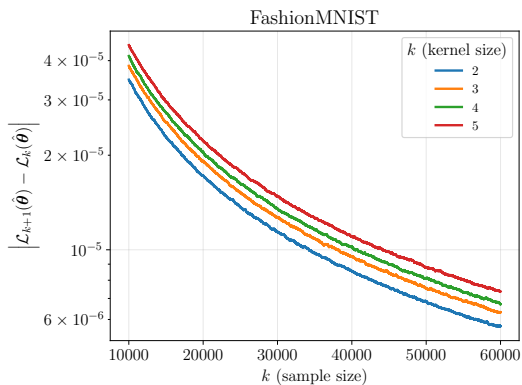
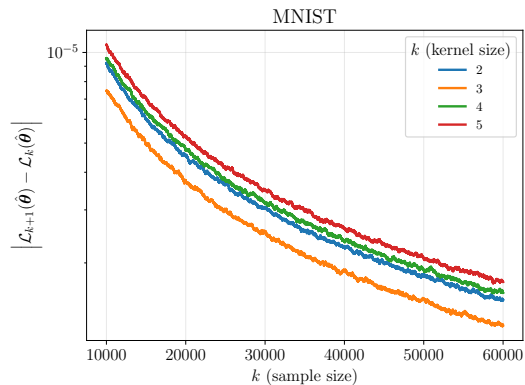
- **Задача:** классификация изображений.
- **Выборка:** MNIST, FashionMNIST, CIFAR10.
- **Архитектура:** L-слойная сверточная нейросеть с ReLU активациями.
- **Проверяемые гипотезы:** Абсолютная разница функции потерь растет в зависимости от
  1. числа сверточных слоев,
  2. размера ядра,
  3. числа фильтров.
- **Постановка эксперимента**
  1. Обучаем модель на полном наборе данных и получаем близкое к оптимальному  $\hat{\theta}$ ,
  2. Находим для всех  $k$ :  $\left| \mathcal{L}_{k+1}(\hat{\theta}) - \mathcal{L}_k(\hat{\theta}) \right|$  для всех  $k = 1, \dots, m$ ,
  3. Предыдущий пункт повторяем меняя порядок элементов в выборке.

## Число слоев



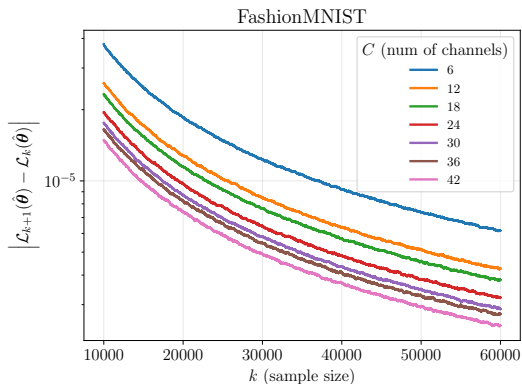
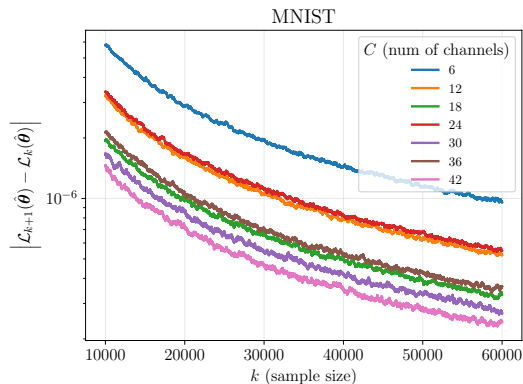
Графики демонстрируют шумное поведение, но видна тенденция: при увеличении числа слоев абсолютная разница функции потерь также увеличивается.

# Размер ядра



При увеличении размера ядра, абсолютная разница функции потерь также растёт.

# Число фильтров



Графики, как видно, демонстрируют обратный результат: при большем числе каналов разница явно меньше. Предположительно, результат таков, так как сети при большем числе фильтров лучше обучились и компонента с Гессианом влияла значительно меньше.

## Основные результаты

1. Предложен способ оценки нормы гессiana для сетей, представленных в виде произведения матриц.
2. Данный результат применен для 2D/1D сверток, а также для пулингов и fc-головы.
3. Предложен способ оценки абсолютной разницы функции потерь для сверточных сетей, основываясь на гессiane.

## Будущее работы

1. Улучшить теоретические оценки, пользуясь разреженностью матриц  $\Lambda^{(i)}$ .
2. Применить данные оценки к другим видам нейронных сетей.
3. Проанализировать другие способы оценки изменения ландшафта функции потерь.