

Дифференцируемый алгоритм поиска архитектуры с контролем сложности

К. Д. Яковлев¹ О. С. Гребенькова¹ О. Ю. Бахтеев^{1,2} В. В. Стрижов^{1,2}
{iakovlev.kd, grebenkova.os, bakhteev, strijov}@phystech.edu

¹Москва, Московский физико-технический институт

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

2021

Цель исследования

Цель

Предложить метод поиска архитектуры модели глубокого обучения с контролем сложности модели.

Проблема

Модели глубокого обучения имеет избыточное число параметров. Поиск архитектуры на дискретном множестве является вычислительно сложной задачей.

Метод решения

Предлагаемый метод основан на непрерывной релаксации. Структурные параметры задаются гиперсетью, зависящей от коэффициента, задающего сложность архитектуры. Под гиперсетью понимается модель, порождающая параметры оптимизируемой модели.

Основная литература

-  Hanxiao Liu and Karen Simonyan and Yiming Yang. *DARTS: Differentiable Architecture Search*. CoRR, 2018.
-  David Ha and Andrew M. Dai and Quoc V. Le. *HyperNetworks*. CoRR, 2016.
-  Grebenkova, O., Bakhteev, O.Y., Strijov, V. *Variational deep learning model optimization with complexity control* 2021
-  Jang, E., Gu, S., Poole, B. *Categorical reparameterization with gumbel-softmax*. CoRR, 2016.

Постановка задачи поиска архитектуры

- Архитектура модели представляет собой ориентированный ациклический граф. Каждому ребру ставится в соответствие отображение $\mathbf{g}^{(i,j)}$, причем

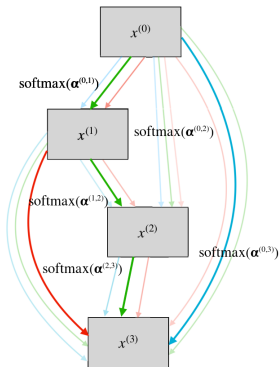
$$\mathbf{x}^{(j)} = \sum_{i < j} \mathbf{g}^{(i,j)}(\mathbf{x}^{(i)}).$$

- Пусть вектор $\vec{\mathbf{g}}^{(i,j)}$ – вектор, составленный из доступных для ребра (i,j) отображений. Пусть вектор $\alpha^{(i,j)}$ – вектор структурных параметров. Смешанная операция

$$\hat{\mathbf{g}}^{(i,j)}(\mathbf{x}^{(i)}) = \langle \mathbf{softmax}(\alpha^{(i,j)}), \vec{\mathbf{g}}^{(i,j)}(\mathbf{x}^{(i)}) \rangle.$$

- Задана выборка $\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{val}}$. Задана функция потерь $\mathcal{L}_{\text{train}}$, \mathcal{L}_{val} . Пусть $\alpha = [\alpha^{(i,j)}]$. Пусть \mathbf{w} – параметры модели. Двухуровневая задача оптимизации

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

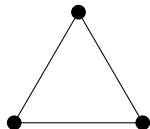


► Смешанная операция:

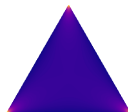
$$\hat{\mathbf{g}}^{(i,j)} = \text{softmax}(\alpha^{(i,j)})_1 \mathbf{g}_1^{(i,j)}(\mathbf{x}^{(i)}) + \\ \text{softmax}(\alpha^{(i,j)})_2 \mathbf{g}_2^{(i,j)}(\mathbf{x}^{(i)}) + \\ \text{softmax}(\alpha^{(i,j)})_3 \mathbf{g}_3^{(i,j)}(\mathbf{x}^{(i)})$$

Распределение гумбель-софтмакс

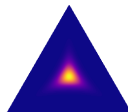
Распределение гумбель-софтмакс определено на симплексе. Пусть $\mathbf{X} \sim \mathcal{GS}(\alpha, t)$, где $\alpha \in \mathbb{R}_{++}^n$, $t > 0$.



$t \rightarrow 0$



$t = 0.995$



$t = 5.0$

Контроль сложности с помощью гиперсети

- ▶ Смешанная операция

$$\hat{\mathbf{g}}^{(i,j)}(\mathbf{x}^{(i)}) = \langle \gamma^{(i,j)}, \vec{\mathbf{g}}^{(i,j)} \rangle, \quad \gamma^{(i,j)} \sim \mathcal{GS}(\exp(\boldsymbol{\alpha}^{(i,j)}), t).$$

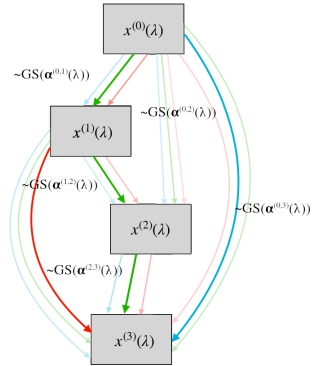
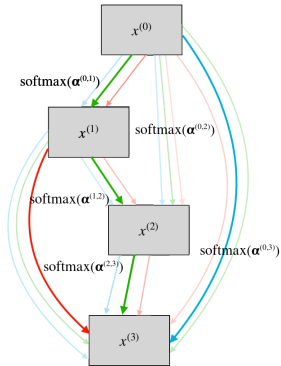
- ▶ Пусть $\Lambda \subset \mathbb{R}$ – множество параметров, задающих сложность. Гиперсеть – это параметрическое отображение из множества Λ во множество структурных параметров модели

$$\boldsymbol{\alpha}^{(i,j)} = \boldsymbol{\alpha}^{(i,j)}(\lambda, \mathbf{a}^{(i,j)}), \quad \lambda \in \Lambda.$$

- ▶ В работе используется кусочно-линейная гиперсеть

$$\boldsymbol{\alpha}^{(i,j)}(\lambda, \mathbf{a}^{(i,j)}) = \sum_{k=0}^{N-1} \left(\frac{\lambda - t_k}{t_{k+1} - t_k} \mathbf{a}_k^{(i,j)} + \left(1 - \frac{\lambda - t_k}{t_{k+1} - t_k} \right) \mathbf{a}_{k+1}^{(i,j)} \right) I[\lambda \in [t_k, t_{k+1}]]$$

DARTS с использованием гиперсети



Структурные параметры порождаются гиперсетью, зависящей от коэффициента, задающего сложность архитектуры. Структурные параметры подчинены распределению Gumbel-Softmax.

Задача оптимизации

- ▶ Пусть вектор $\mathbf{n}(\vec{\mathbf{g}}^{(i,j)})$ хранит количество параметров каждого отображения. Регуляризатор, контролирующий сложность

$$\lambda \sum_{(i,j)} \langle \mathbf{softmax} \left(\alpha^{(i,j)}(\lambda, \mathbf{a}^{(i,j)}) \right), \mathbf{n}(\vec{\mathbf{g}}^{(i,j)}) \rangle.$$

- ▶ Пусть задано распределение $p(\lambda)$ на Λ . Пусть $\gamma = [\gamma^{(i,j)}]$. Параметры $\mathbf{a} = [\mathbf{a}^{(i,j)}]$ гиперсети находятся из задачи оптимизации

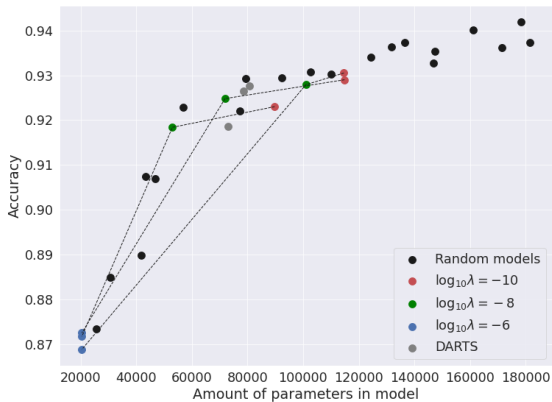
$$\min_{\mathbf{a}} \mathbb{E}_{\lambda \sim p(\lambda)} \left(\mathbb{E}_{\gamma} \mathcal{L}_{\text{val}}(\mathbf{w}^*, \gamma) + \lambda \sum_{(i,j)} \langle \mathbf{softmax} \left(\alpha^{(i,j)}(\lambda, \mathbf{a}^{(i,j)}) \right), \mathbf{n}(\vec{\mathbf{g}}^{(i,j)}) \rangle \right),$$

$$\text{s.t. } \mathbf{w}^* = \arg \min_{\mathbf{w}} \mathbb{E}_{\lambda \sim p(\lambda)} \mathbb{E}_{\gamma} \mathcal{L}_{\text{train}}(\mathbf{w}, \gamma).$$

Постановка вычислительного эксперимента

- ▶ Цель эксперимента – получение зависимости обобщающей способности модели от количества её параметров.
- ▶ Вычислительный эксперимент проводится на выборке Fashion-MNIST. Сравниваются архитектуры, полученные с помощью DARTS, предлагаемого метода и случайные архитектуры.
- ▶ Модель состоит из трех ячеек. Коэффициент $\lambda \sim \mathcal{U}[10^{-10}, 10^{-6}]$. Во время обучения температура распределения гумбель-софтмакс понижалась от 1 до 0.2.

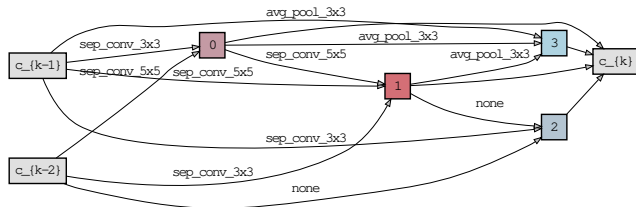
Результаты вычислительного эксперимента



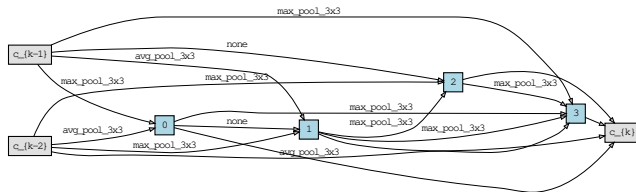
Зависимость качества классификации от количества параметров модели.

Предложенный метод позволяет контролировать сложность архитектуры, изменяя коэффициент регуляризации λ .

Полученные архитектуры



(a) Архитектура, полученная при $\lambda = 10^{-10}$.



(b) Архитектура, полученная при $\lambda = 10^{-6}$.

Чем больше коэффициент регуляризации λ , тем проще получаемая архитектура.

- ▶ Предложен метод, позволяющий контролировать сложность модели в процессе поиска архитектуры.
- ▶ Метод обладает тем свойством, что изменение сложности итоговой модели происходит изменением коэффициента, задающего сложность архитектуры, без дополнительного обучения.
- ▶ Также результаты показывают, что данный метод сопоставим по качеству с DARTS.