

Оптимизация архитектуры нейросети с контролем эксплуатационных характеристик на целевом устройстве

Фирсов Сергей

Научный руководитель: к.ф.-м.н. О.Ю. Бахтеев

Московский физико-технический институт

2025

- ▶ **Neural Architecture Search:** Задача автоматизированного поиска оптимальной архитектуры нейросети.
- ▶ **Цель:** Получать архитектуры решающие поставленную ML задачу, при этом удовлетворяя вычислительным или ресурсным ограничениям.
- ▶ **Проблемы:**
 - ▶ Обширное пространство для поиска
 - ▶ Баланс точности и сложности
 - ▶ Получение семейства решений

Постановка задачи

- ▶ Многоклассовая классификация. Выборка $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{\mathcal{N}}$, $\mathbf{x}_i \in X$, $y_i \in Y$.
- ▶ Модель — нейронная сеть из повторяющихся блоков, каждый из которых представляется в виде DAG. Задача выбирать рёбра-операции $\mathbf{o}^{(m)}$ для каждой пары вершин в блоке.
- ▶ Для релаксации дискретной задачи выбора архитектуры в непрерывную, вводится *смешанная операция*

$$\hat{\mathbf{g}}^{(i,j)}(\mathbf{x}) = \sum_{m=1}^k \gamma_m^{(i,j)} \mathbf{o}^{(m)}(\mathbf{x}), \quad \gamma^{(i,j)} \sim \text{GumbelSoftmax}(\boldsymbol{\alpha}^{(i,j)}, t).$$

Регуляризация сложности

Вводится *векторный параметр сложности* \mathbf{s} , компоненты которого — коэффициенты регуляризации по соответствующим операциям.

$$\mathbf{s} \in \Delta^{k-1}, \quad \Delta^{k-1} = \left\{ \mathbf{s} \in \mathbf{R}^k \mid \sum_{m=1}^k s_m = 1, s_m \geq 0 \right\},$$

где $k = |\mathcal{O}|$, Δ^{k-1} симплекс размерности $k - 1$.

Определим функцию затрат Cost архитектуры γ на основе \mathbf{s} :

$$\text{Cost}(\gamma; \mathbf{s}) = \sum_{(i,j) \in E} \sum_{m=1}^k \gamma_m^{(i,j)} s_m.$$

Оптимизационная задача

Для генерации архитектуры вводится гиперсеть $u_{\mathbf{a}}: \Delta^{k-1} \rightarrow \Gamma$ с параметрами \mathbf{a} , которая по заданному вектору сложности $\mathbf{s} \in \Delta^{k-1}$ порождает логиты $\boldsymbol{\alpha}$, тем самым определяя архитектуру модели с учётом эксплуатационных характеристик.

Theorem (Фирсов, 2025)

$$\gamma(t) \sim \text{GumbelSoftmax}(\boldsymbol{\alpha}, t), \quad \bar{\gamma}(t) = \text{Softmax}(\boldsymbol{\alpha}/t) = \mathbb{E}[\gamma(t)].$$

Тогда

$$\lim_{t \rightarrow 0^+} \mathbb{E}[f(\gamma(t))] = f(\bar{\gamma}(0)).$$

$$\mathbb{E}_{\mathbf{s} \sim U(\Delta^{k-1})} \left[\mathbb{E}_{\gamma \sim \text{GS}(\boldsymbol{\alpha}(\mathbf{s}), t)} \mathcal{L}_{\text{task}}(\mathbf{w}, \gamma) + \kappa \text{Cost}(\bar{\gamma}(\mathbf{s}); \mathbf{s}) \right] \rightarrow \min_{\mathbf{w}, \mathbf{a}}$$

Иллюстрация

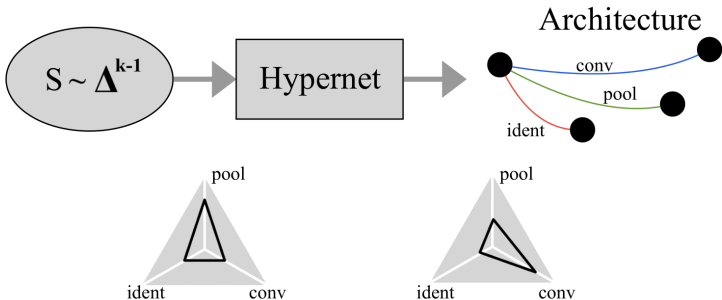


Рис. 1. Иллюстрация предлагаемого метода.

- ▶ Предлагается использовать векторный параметр сложности s , компоненты которого — коэффициенты регуляризации по соответствующим операциям.
- ▶ Гиперсеть на основе s генерирует архитектурные параметры для нейросети.

Постановка эксперимента

- ▶ Эксперименты поставлены для подтверждения работоспособности метода и иллюстрации его возможностей.
- ▶ Представленный эксперимент проводится на выборке Fashion-MNIST.
- ▶ Рассматриваются модели из 3 блоков по 4 вершины.
- ▶ Параметры \mathbf{w} и параметры гиперсети \mathbf{a} тренировались с помощью оптимизаторов Adam.
- ▶ Температура Gumbel–Softmax линейно уменьшалась от $t=1.0$ до $t=0.2$ по эпохам.

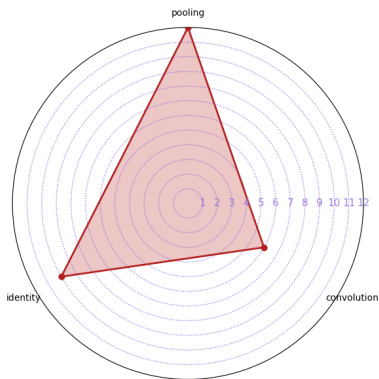
После обучения получено семейство нейронных сетей, из которого для анализа выделяется четыре архитектуры, соответствующие приведённым ниже векторам сложности $s = (pool, ident, conv)$. Компоненты вектора являются штрафами за соответствующие операции.

(A) (0.33, 0.33, 0.33), (B) (0.15, 0.15, 0.70),

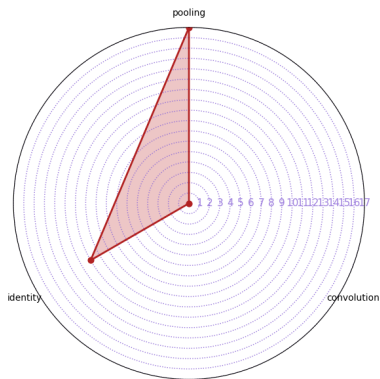
(C) (0.70, 0.15, 0.15), (D) (0.15, 0.70, 0.15).

Представленные архитектуры отражают равномерный штраф за операции (A), увеличенный за свёртки (B), увеличенный за пулинг (C) и увеличенный за пропуск (D).

Распределение операций



(A) равномерный штраф

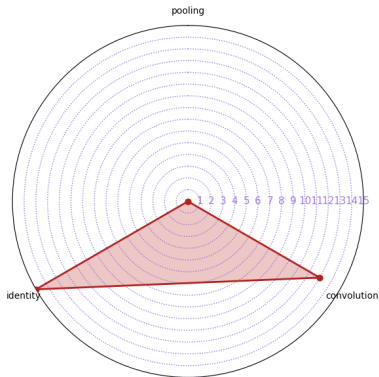


(B) увеличенный за свёртки

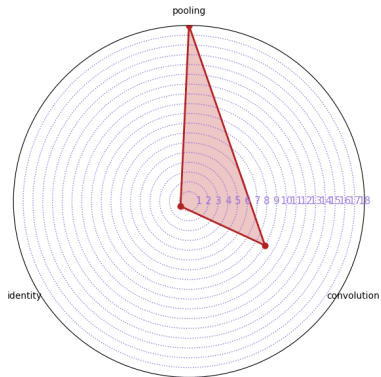
Рис. 2. Статистика по операциям для архитектур эксперимента 1:

(A) равномерный штраф, (B) увеличенный за свёртки.

Распределение операций



(C) увеличенный за пулинг



(D) увеличенный за пропуск

Рис. 3. Статистика по операциям для архитектур эксперимента 1:

(C) увеличенный за пулинг, (D) увеличенный за пропуск.

Результаты

Модели	A	B	C	D
Accuracy	82.5%	79.2%	85.0%	81.4
Количество параметров	38 304	5 120	143 488	12 320
Количество pooling	12	17	0	18
Количество convolution	6	0	13	9
Количество identity	10	11	15	1

Таблица 1. Сравнение параметров полученных архитектур. (A) равномерный штраф, (B) увеличенный за свёртки, (C) увеличенный за пулинг, (D) увеличенный за пропуск.

Результаты экспериментов подтверждают корректность работы предлагаемого метода. Заметим, что при увеличении штрафа за операцию — уменьшается её использование в модели. Такая регуляризация позволяет учитывать различные эксплуатационные ограничения без необходимости дообучения модели на них.

Результаты

- ▶ Предложен метод позволяющий получать семейство моделей с возможностью контроля сложности обучения и аппаратных ограничения.
- ▶ Метод обладает возможностью получать архитектуры моделей за счёт изменения вектора параметра сложности без необходимости дополнительного дообучения.
- ▶ Вычислительные эксперименты подтверждают работоспособность метода и демонстрируют заявленную гибкость получаемого решения.

Результаты работы доложены на конференции МФТИ и в университете Сириус.

В рамках дальнейшего развития метода планируется добавление многокритериального учета сложности и проведение полноформатного сопоставления результатов с классическими NAS-подходами.

Список литературы

- [1] K. Yakovlev и др. «Neural architecture search with structure complexity control». В: *Recent Trends in Analysis of Images, Social Networks and Texts*. 2022, с. 207—219.
- [2] J. Dong, Y. Yang и ... «BigNAS: Scaling Up Neural Architecture Search with Big Single Stage Models». В: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2021, с. 3173—3182.
- [3] T. Elsken, J. H. Metzen и F. Hutter. «Neural Architecture Search: A Survey». В: *Journal of Machine Learning Research* 20.55 (2019), с. 1—21.
- [4] H. Liu, K. Simonyan и Y. Yang. «DARTS: Differentiable architecture search». В: *International Conference on Learning Representations*. 2019.
- [5] H. Cai, L. Zhu и S. Han. *ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware*. 2019.

Примечание. Клетка

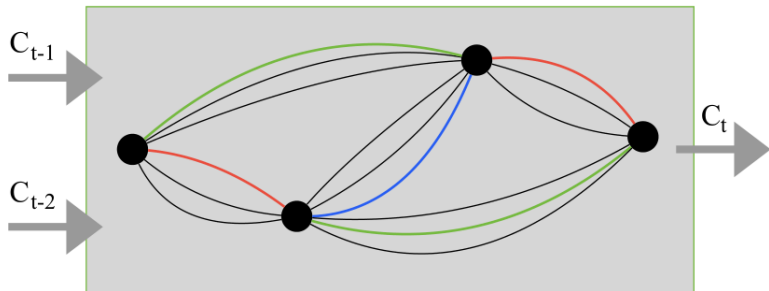


Рис. 4. Иллюстрация клетки.

Примечание. Гиперсеть

Введём гиперсеть \mathbf{u}_a с параметрами \mathbf{a} , которая на основе \mathbf{s} строит логиты $\alpha^{(i,j)}$ для каждого ребра (i,j) как взвешенную комбинацию заранее выученных логитов опорных архитектур $\alpha_{\text{ref},r}^{(i,j)}$, полученных в процессе оптимизации гиперсети. По вектору сложности \mathbf{s} вычисляются веса интерполяции

$$w_r(\mathbf{s}) = \frac{\exp(-\|\mathbf{s} - \mathbf{s}_{\text{ref},r}\|^2)}{\sum_{l=1}^m \exp(-\|\mathbf{s} - \mathbf{s}_{\text{ref},l}\|^2)}, \quad r = 1, \dots, m,$$

где $\mathbf{s}_{\text{ref},r} \in \Delta^{k-1}$ — r -я опорная точка, а m — их общее число. Затем итоговые логиты для ребра (i,j) вычисляются как

$$\alpha^{(i,j)}(\mathbf{s}) = \sum_{r=1}^m w_r(\mathbf{s}) \alpha_{\text{ref},r}^{(i,j)}.$$

Примечание. Теорема

Стохастическое порождение $\gamma \sim GS$ вносит высокую дисперсию градиентов, особенно при малой температуре t . Так как математическое ожидание GS -вектора совпадает с обычным Softmax $\hat{\gamma} = \mathbb{E}[\gamma] = \text{Softmax}(\alpha/t)$, можно безопасно подменить γ на $\hat{\gamma}$ в Cost части функции потерь, где важна стабильность. Формальную корректность такой подмены устанавливает следующая теорема.

Theorem (Фирсов, 2025)

Пусть для каждого набора весов $\mathbf{w} \in \mathbb{R}^n$ и любого входа $\mathbf{x} \in \mathcal{X}$ функция

$$f(\gamma) = \mathcal{L}_{\text{task}}(\mathbf{w}, \gamma) + \kappa \langle \gamma, \mathbf{s} \rangle, \quad \mathbf{s} \in \Delta^{k-1},$$

непрерывна и ограничена на симплексе

$\Delta^{k-1} = \{\gamma \in \mathbb{R}^k \mid \sum_{i=1}^k \gamma_i = 1, \gamma_i \geq 0\}$. Положим

$$\gamma(t) \sim \text{GumbelSoftmax}(\alpha, t), \quad \bar{\gamma}(t) = \text{Softmax}(\alpha/t) = \mathbb{E}[\gamma(t)].$$

Тогда

$$\lim_{t \rightarrow 0^+} \mathbb{E}[f(\gamma(t))] = f(\bar{\gamma}(0)).$$

Примечание. Оптимизационная задача

$$\mathbb{E}_{\mathbf{s} \sim U(\Delta^{k-1})} \left[\mathbb{E}_{\gamma \sim \text{GS}(\alpha(\mathbf{s}), t)} \mathcal{L}_{\text{task}}(\mathbf{w}, \gamma) + \kappa \text{Cost}(\bar{\gamma}(\mathbf{s}); \mathbf{s}) \right] \rightarrow \min_{\mathbf{w}, \mathbf{a}},$$

где \mathbf{w} — все параметры внутри примитивных операций $\mathcal{O}^{(m)}$; \mathbf{a} — параметры гиперсети \mathbf{u}_a , которая генерирует логиты $\alpha^{(i,j)}$ для всех $(i, j) \in E$ при подстановке \mathbf{s} ;
 $\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_k) \in \Delta^{k-1}$ — вектор сложности, порождённый из равномерного распределения на симплексе Δ ; $\mathcal{L}_{\text{task}}$ — кросс-энтропия по обучающей выборке; $\gamma^{(i,j)}(\mathbf{s})$ — веса операций на ребре (i, j) , получаемые через Gumbel-Softmax от логитов $\mathbf{u}_a^{(i,j)}(\mathbf{s})$; γ — вектор состоящий из $\gamma^{(i,j)}$; $\bar{\gamma} = \mathbb{E}[\gamma(t)]$; Cost — штраф за использование дорогих операций; κ — гиперпараметр, управляющий важностью штрафа.