

Байесовский выбор архитектуры нейросетевой модели

Сотников А. Д., группа М05-0046
Научный руководитель: к. ф-м н. Бахтеев О.Ю.

Московский Физико-Технический институт
Кафедра интеллектуальных систем

21 июня, 2022

Введение

Задача

Поиск архитектуры нейросетевой модели (ПАНМ) — метод автоматического проектирования архитектуры нейронной сети на заданных задаче и наборе данных.

Мотивация

Методы ПАНМ не обладают достаточной устойчивостью, то есть уязвимы к внешним воздействиям.

Гипотеза

Проведение байесовского вывода повышает устойчивость базового метода.

Предложение

Реализовать модификацию базового метода поиска архитектуры с помощью байесовского вывода распределений параметров и структуры модели.
Предоставить теоретическую интерпретацию предлагаемого метода.

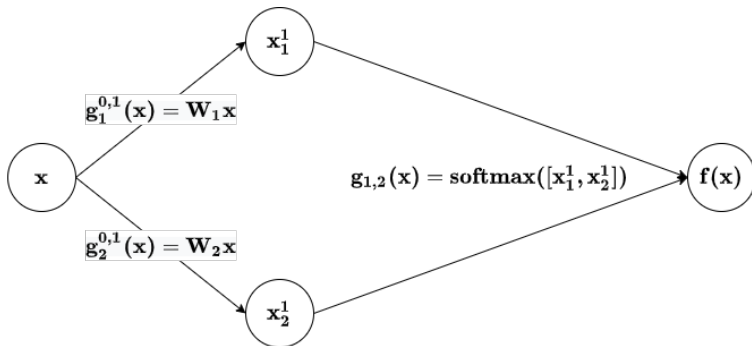
Основные определения

Определение 1

Моделью называется дифференцируемая по параметрам функция

$$\mathbf{f}(\mathbf{w}, \mathbf{x}) = y : \mathbb{W} \times \mathbb{X} \rightarrow \mathbb{Y},$$

где $\mathbf{w} \in \mathbb{W}$ - параметры модели, $\mathbf{x} \in \mathbb{X}$ - признаковое описание входного объекта, $y \in \mathbb{Y}$ - метка входного объекта.



Определение 2

Пусть на ребре (j, k) задан вектор операций $\mathbf{g}^{j,k}$, $|\mathbf{g}^{j,k}| = N^{j,k}$. Структурными параметрами назовём вектор $\gamma^{j,k} = [0, 1]^{N^{j,k}}$.

Структурой модели называется конкатенация её структурных параметров $\Gamma = \{\gamma^{j,k} | (j, k) \in E\}$.

Определение 3

Архитектурой модели называется совокупность её параметров и структуры.

Определение 4

Гиперапараметрами $\mathbf{h} \in \mathbb{H}$ модели назовём параметры распределения $p(\mathbf{w}, \Gamma | \mathbf{h})$.

Существующие методы

Методы ПАНМ

- Методы, основанные на обучении с подкреплением
- Градиентные методы
- Вероятностный подход

Методы, повышающие устойчивость ПАНМ

- "Отравление" пространства поиска
- Оптимизация мер устойчивости
- Регуляризация структурных параметров

Постановка задачи

- Набор данных $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, $\mathbf{x}_i \in \mathbb{R}^m$, $y_i \in \mathbb{N}$
- Параметры $\mathbf{w} \sim p(\mathbf{w} | \mathbf{h})$ и структура $\Gamma \sim p(\Gamma | \mathbf{h})$ модели являются независимыми случайными величинами
- Вероятностная модель:

$$p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{y}, \mathbf{h}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \mathbf{h}) \cdot p(\Gamma | \mathbf{X}, \mathbf{y}, \mathbf{h})$$

- Вводится вариационное распределение

$$q(\mathbf{w}, \Gamma | \theta) = q_{\mathbf{w}}(\mathbf{w} | \theta_{\mathbf{w}}) \cdot q_{\Gamma}(\Gamma | \theta_{\Gamma})$$

- Функция ошибки:

$$\mathcal{L}(\mathbf{w}, \Gamma) = \mathcal{D}_{KL}(q(\mathbf{w}, \Gamma | \theta) || p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{y}, \mathbf{h}))$$

- Ставится двухуровневая оптимизационная задача:

$$\begin{aligned} \min_{\Gamma} \quad & \mathcal{L}_{val}(\hat{\mathbf{w}}, \Gamma) \\ \text{s.t.} \quad & \hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \mathcal{L}_{train}(\mathbf{w}, \Gamma) \end{aligned}$$

Вариационная нижняя оценка обоснованности

По теореме Байеса

$$p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{y}, \mathbf{h}) = \frac{p(\mathbf{y} | \mathbf{w}, \Gamma, \mathbf{X}) \cdot p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{h})}{p(\mathbf{y} | \mathbf{X}, \mathbf{h})}$$

Выразим обоснованность через вариационное распределение

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{X}, \mathbf{h}) &= \iint_{\mathbf{w}, \Gamma} q(\mathbf{w}, \Gamma | \theta) \log p(\mathbf{y} | \mathbf{X}, \mathbf{h}) d\Gamma d\mathbf{w} = \\ \mathbb{E}_{q(\mathbf{w}, \Gamma | \theta)} \left[\log \frac{p(\mathbf{y}, \mathbf{w}, \Gamma | \mathbf{X}, \mathbf{h})}{q(\mathbf{w}, \Gamma | \theta)} \right] &+ \mathcal{D}_{KL}(q(\mathbf{w}, \Gamma | \theta) || p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{y}, \mathbf{h})) \end{aligned}$$

Утверждение 1

$$\min_{\theta} \mathcal{D}_{KL}(q(\mathbf{w}, \Gamma | \theta) || p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{y}, \mathbf{h})) \iff \max_{\theta} \mathbb{E}_{q(\mathbf{w}, \Gamma | \theta)} \left[\log \frac{p(\mathbf{y}, \mathbf{w}, \Gamma | \mathbf{X}, \mathbf{h})}{q(\mathbf{w}, \Gamma | \theta)} \right]$$

- Матожидание можно представить в виде

$$\begin{aligned} & \mathbb{E}_{q(\mathbf{w}, \Gamma | \theta)} \left[\log \frac{p(\mathbf{y}, \mathbf{w}, \Gamma | \mathbf{X}, \mathbf{h})}{q(\mathbf{w}, \Gamma | \theta)} \right] = \\ & = \mathbb{E}_{q(\mathbf{w}, \Gamma | \theta)} \left[\log p(\mathbf{y} | \mathbf{w}, \Gamma, \mathbf{X}) \right] - \mathcal{D}_{KL}(q(\mathbf{w}, \Gamma | \theta) || p(\mathbf{w}, \Gamma | \mathbf{X}, \mathbf{h})) \end{aligned}$$

- KL-дивергенция раскладывается на

$$\mathcal{D}_{KL}(q_{\Gamma}(\Gamma | \theta_{\Gamma}) || p(\Gamma | \mathbf{h})) + \mathcal{D}_{KL}(q_{\mathbf{w}}(\mathbf{w} | \theta_{\mathbf{w}}) || p(\mathbf{w} | \mathbf{h}))$$

Окончательно, функция ошибки принимает вид

$$\begin{aligned} \mathcal{L}(\mathbf{w}, \Gamma) = & -\mathbb{E}_{q(\mathbf{w}, \Gamma | \theta)} \log p(\mathbf{y} | \mathbf{w}, \Gamma, \mathbf{X}) + \\ & + \mathcal{D}_{KL}(q_{\Gamma}(\Gamma | \theta_{\Gamma}) || p(\Gamma | \mathbf{h})) + \mathcal{D}_{KL}(q_{\mathbf{w}}(\mathbf{w} | \theta_{\mathbf{w}}) || p(\mathbf{w} | \mathbf{h})) \end{aligned}$$

Вариационный вывод ПАНМ

Теорема 1 (Сотников, 2022)

Пусть заданы две структуры $\hat{\Gamma}_1$ и $\hat{\Gamma}_2$, функция ошибки на которых принимает одинаковое значение, то есть $\mathcal{L}(\mathbf{w}, \hat{\Gamma}_1) = \mathcal{L}(\mathbf{w}, \hat{\Gamma}_2)$. Пусть также $\|\nabla_{\Gamma}^2 \mathcal{L}(\mathbf{w}, \hat{\Gamma}_1)\| < \|\nabla_{\Gamma}^2 \mathcal{L}(\mathbf{w}, \hat{\Gamma}_2)\|$. Тогда справедливо

$$\|\hat{\Gamma}_1 - \Gamma^*\| > \|\hat{\Gamma}_2 - \Gamma^*\|,$$

где $\Gamma^* = \arg \min_{\Gamma} \mathcal{L}_{val}(\mathbf{w}, \Gamma)$ – оптимальная структура.

Теорема 2 (Сотников, 2022)

Пусть задана функция $G(\mathbf{w}, \Gamma)$. Тогда для любого распределения $q_{\Gamma}(\Gamma | \delta)$ такого, что компоненты Γ – независимые случайные величины, справедливо

$$\mathbb{E}_{q_{\Gamma}(\Gamma | \delta)} [G(\mathbf{w}, \Gamma)] \approx G(\mathbf{w}, \mu) + \frac{\sigma^2}{2} \text{Tr}(\nabla_{\Gamma}^2 G(\mathbf{w}, \mu)),$$

где $\mu = \mathbb{E}_{q_{\Gamma}(\Gamma | \delta)} [\Gamma]$ и $\sigma^2 = \mathbb{D}_{q_{\Gamma}(\Gamma | \delta)} [\Gamma]$.

Параметры эксперимента

Размер изображения	Число классов	Размер \mathcal{D}_{train}	Размер \mathcal{D}_{val}
$1 \times 28 \times 28$	10	60000	10000

Априорные распределения

- $q_{\mathbf{w}}(\mathbf{w} | \theta_{\mathbf{w}}) \sim \mathcal{N}(\mathbf{m}_{\mathbf{w}}, \mathbf{A}_{\mathbf{w}}^{-1})$
- $q_{\Gamma}(\Gamma | \theta_{\Gamma}) \sim$
Gumbel-Softmax($\alpha_1, \dots, \alpha_N$)

T-Shirt/Top

Trouser

Pullover

Dress

Coat

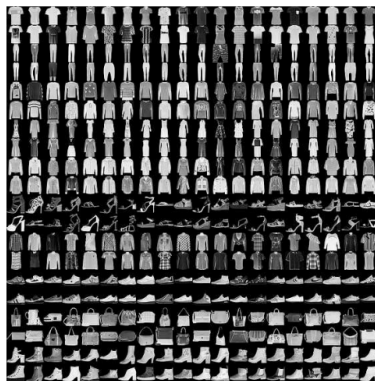
Sandals

Shirt

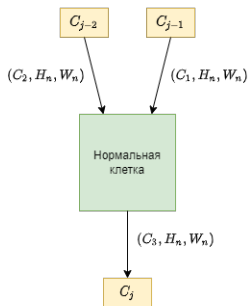
Sneaker

Bag

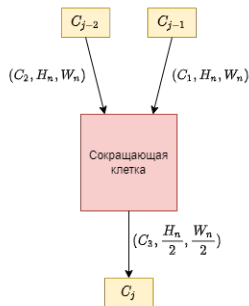
Ankle boots



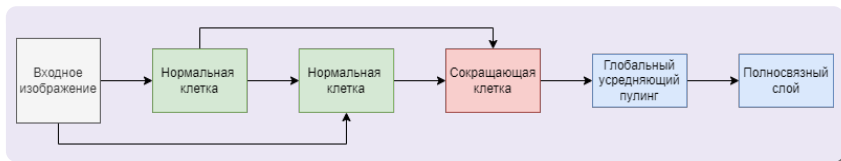
Архитектура поиска



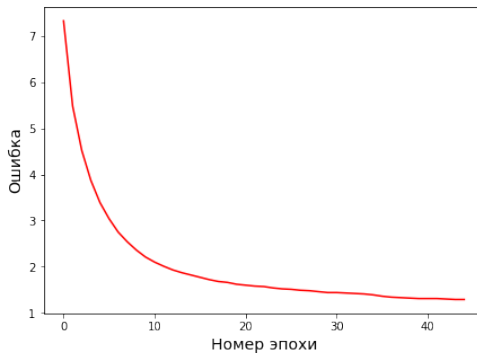
(а)



(б)

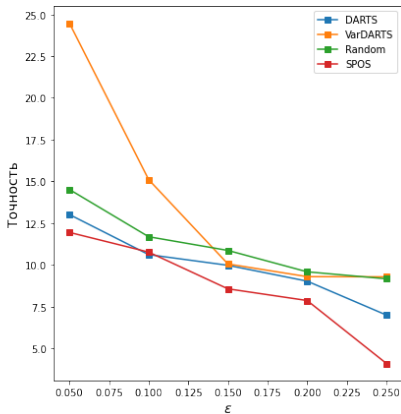
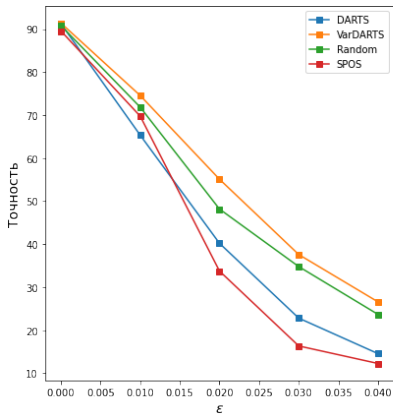


Сравнение с существующими методами

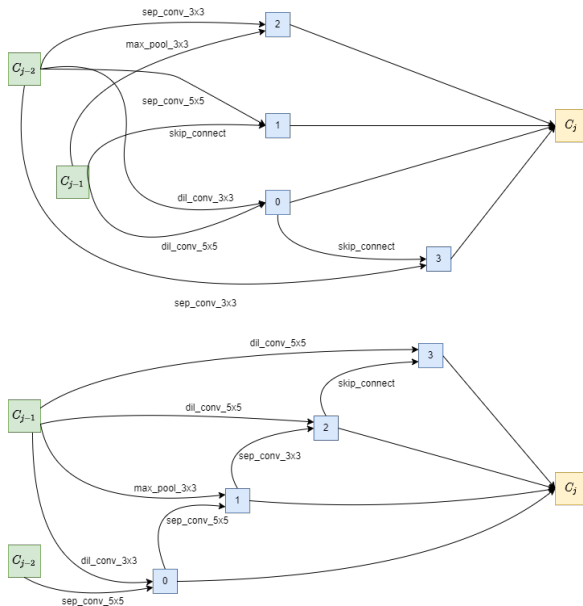


Метод	Accuracy top-1, %	Время обучения, ч
DARTS	91.22	2
Random	90.79	1.5
SPOS	89.54	0.5
VarDARTS	91.36	5

Сравнение с существующими методами



Обученные клетки



Выносятся на защиту

Полученные результаты

- Предложен метод, повышающий устойчивость градиентного ПАНМ, с помощью байесовского вывода
- Предложена теоретическая интерпретация реализованного метода

Дальнейшие исследования

- добавить зависимость распределения параметров модели от её структуры;
- добавить распределения на гиперпараметры модели;
- заменить априорные предположения на другие распределения.