

Выбор структуры модели глубокого обучения

Бахтеев Олег

МФТИ

13.02.2019

Графовое представление модели глубокого обучения

Определение

Задан граф (V, E) . Для каждого ребра $(j, k) \in E$ определен вектор базовых функций мощности $K^{j,k}$:

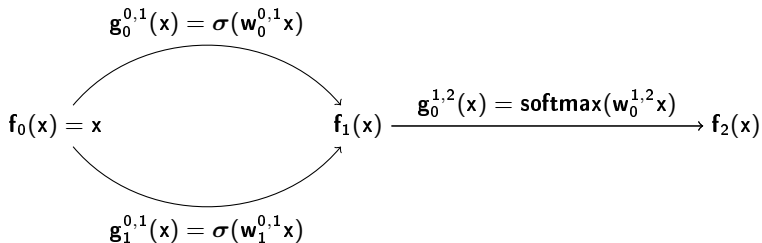
$$\mathbf{g}^{j,k} = [\mathbf{g}_0^{j,k}, \dots, \mathbf{g}_{K^{j,k}}^{j,k}]$$

. Пусть для каждой вершины $v \in V$ определена функция агрегации \mathbf{agg}_v . Граф (V, E) в совокупности со множеством векторов базовых функций $\{\mathbf{g}^{j,k}, (j, k) \in E\}$ и множеством функций агрегаций $\{\mathbf{agg}_v, v \in V\}$ называется *параметрическим семейством моделей* \mathfrak{F} , если функция, задаваемая как

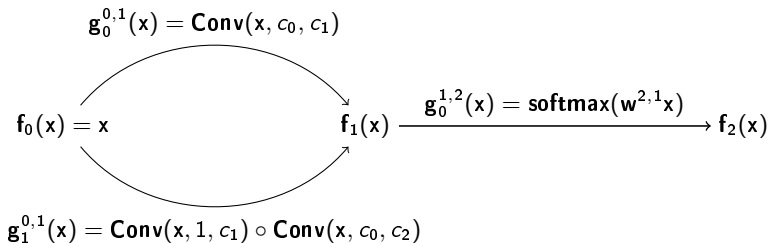
$$\mathbf{f}_k(\mathbf{x}) = \mathbf{agg}_k (\{ \langle \gamma^{j,k}, \mathbf{g}^{j,k} \rangle (\mathbf{f}_j(\mathbf{x})) \mid j \in \text{Adj}(v_k) \}), \quad \mathbf{f}_0(\mathbf{x}) = \mathbf{x} \quad (1)$$

является моделью при любых значениях векторов, $\gamma^{j,k} \in [0, 1]^{K^{j,k}}$.

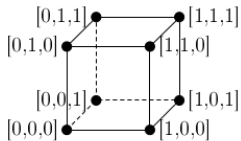
Пример: двуслойная нейросеть



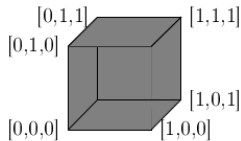
Пример: сверточная сеть



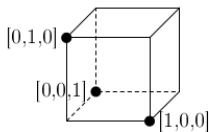
Ограничения на структурные параметры



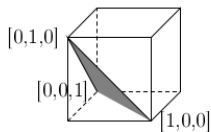
(a)



(б)



(в)



(г)

Статистические критерии качества модели

Параметрическая сложность — наименьшая дивергенция между априорным распределением параметров и апостериорным распределением параметров:

$$C_{\text{param}} = \min_{\mathbf{h}} D_{\text{KL}}(p(\mathbf{W}, \Gamma | \mathbf{y}, \mathbf{X}) || p(\mathbf{W}, \Gamma | \mathbf{h})).$$

Структурная сложность модели — энтропия апостериорного распределения структуры модели:

$$C_{\text{struct}} = -E_p \log p(\Gamma | \mathbf{y}, \mathbf{X}).$$

Выбор оптимальной модели

Основные проблемы выбора оптимальной модели

- Интеграл правдоподобия $p(\mathbf{y}|\mathbf{X}, \mathbf{h})$ невычислим аналитически.
- Задача его оптимизации многоэкстремальна и невыпукла.

Требуется

Предложить метод поиска субоптимального решения задачи оптимизации, обобщающего различные алгоритмы оптимизации:

- Оптимизация правдоподобия.
- Последовательное увеличение и снижение сложности модели.
- Полный перебор вариантов структуры модели.

Распределение на структуре

Пусть для каждого ребра (j, k) задан нормированный положительный вектор $\gamma_{j,k} \in \mathbb{R}_+^{|K_{j,k}|}$, определяющий веса базовых функций из $\mathbf{g}(j, k)$. Перечислим основные свойства, которыми должно обладать распределение такого вектора:

- 1 $p(\gamma^{j,k})$ является непрерывным на симплексе $\Delta^{K^{j,k}-1}$.
- 2 При устремлении температуры к бесконечности распределение сходится к равномерному: $\lim_{c_{\text{temp}} \rightarrow \infty} p(\gamma^{j,k} | c_{\text{temp}}) = \mathcal{U}(\Delta^{K^{j,k}-1})$.
- 3 При устремлении температуры к нулю распределение сходится к сингулярному распределению следующего вида:

$$\lim_{c_{\text{temp}} \rightarrow 0} p(\gamma_k^{j,k}) = m_k^{j,k},$$

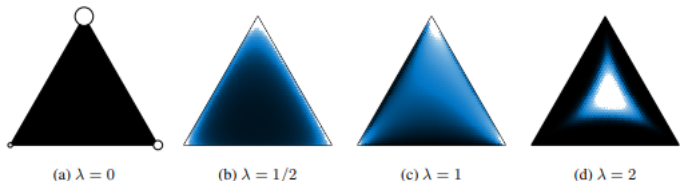
где $m^{j,k}$ — параметр распределения.

Варианты распределений

- 1 Дирихле;
- 2 Гумбель-Софтмакс:

$$\hat{\gamma}_h = \exp(\log(m_h + \text{Gum}_h) c_{\text{temp}}^{-1}) \sum_{l=1}^{K_{j,k}} \exp(\log(m_l + \text{Gum}_l) c_{\text{temp}}^{-1}),$$

где $\text{Gum} \sim -\log(-\log \mathcal{U}(0, 1))$.



Maddison et al., 2017.

Оптимизация параметров вариационного распределения

Параметры вариационного распределения $q(\mathbf{W}, \Gamma) = q_{\mathbf{W}}(\mathbf{W})q_{\Gamma}(\Gamma)$ оптимизируем:

$$L = E_q \log p(\mathbf{y} | \mathbf{X}, \mathbf{W}, \Gamma, \mathbf{A}^{-1}, c_{\text{temp}}) - c_{\text{reg}} D_{KL}(p(\mathbf{w}, \Gamma | \mathbf{A}^{-1}, \mathbf{m}, c_{\text{temp}}) || q(\mathbf{W}), q(\Gamma)) \rightarrow \max_{\mathbf{A}_q, \mu_q, \mathbf{m}_q} .$$

Теорема

Пусть $c_{\text{reg}} > 0$. Тогда $\frac{1}{m} L(c_{\text{reg}})$ сходится п.н. к той же функции, что и $\frac{c_{\text{reg}}}{m_0} L(c_{\text{reg}} = 1)$.

Интерпретация: для достаточно большого m и $c_{\text{reg}} \neq 1$ оптимизация параметров и гиперпараметров эквивалентна оптимизации ELBO для выборки другой мощности.

Теорема [Бахтеев, 2018].

Пусть Γ_1 и Γ_2 — реализации Γ , такие что:

- $\Gamma_1 \in \bar{\Delta}(\Gamma)$.
- $\Gamma_2 \notin \bar{\Delta}(\Gamma)$.

Тогда для любых положительно определенных матриц \mathbf{A}_1 и \mathbf{A}_2 и векторов $\mathbf{m}_1, \mathbf{m}_2, \min(\mathbf{m}_1) > 0$ справедлива следующее отношение апостериорных вероятностей:

$$\lim_{c_{\text{temp}} \rightarrow 0} \frac{p(\Gamma_2, \mathbf{W}_2 | \mathbf{y}, \mathbf{X}, \mathbf{A}_1, \mathbf{m}_2, c_{\text{temp}})}{p(\Gamma_1, \mathbf{W} | \mathbf{y}, \mathbf{X}, \mathbf{A}_1, \mathbf{m}_1, c_{\text{temp}})} = \infty.$$

Оптимизация параметров априорного распределения

Гиперпараметры \mathbf{A} , \mathbf{m} оптимизируем:

$$Q = c_{\text{train}} E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{W}, \mathbf{\Gamma}, \mathbf{A}^{-1}, c_{\text{prior}}) - c_{\text{prior}} D_{\text{KL}}(p(\mathbf{W}, \mathbf{\Gamma}|\mathbf{A}^{-1}, \mathbf{m}, c_{\text{emp}})||q(\mathbf{W}, \mathbf{\Gamma})) - c_{\text{comb}} \sum_{p' \in \mathcal{P}} D_{\text{KL}}(\mathbf{\Gamma}|p') \rightarrow \max,$$

где \mathcal{P} — множество (возможно пустое) распределений на структуре модели.

- c_{train} — коэффициент правдоподобия выборки;
- c_{prior} — коэффициент регуляризации модели;
- c_{comb} — коэффициент перебора структуры.

Общая задача оптимизации

Общая задача оптимизации — двухуровневая:

$$\begin{aligned} \hat{\mathbf{A}}, \hat{\mathbf{m}} &= \arg \max_{\mathbf{A}, \mathbf{m}} Q = \\ &= c_{\text{train}} E_{\hat{q}} \log p(y|\mathbf{X}, \mathbf{W}, \mathbf{\Gamma}, \mathbf{A}^{-1}, c_{\text{prior}}) - c_{\text{prior}} D_{\text{KL}}(p(\mathbf{W}, \mathbf{\Gamma} | \mathbf{A}^{-1}, \mathbf{m}, c_{\text{temp}}) || \hat{q}(\mathbf{W}, \mathbf{\Gamma})) - \\ &\quad - c_{\text{comb}} \sum_{\rho' \in \mathcal{P}} D_{\text{KL}}(\mathbf{\Gamma} | \rho'), \end{aligned}$$

где

$$\hat{q} = \arg \max_q L = E_q \log p(y|\mathbf{X}, \mathbf{W}, \mathbf{\Gamma}, \mathbf{A}^{-1}, c_{\text{temp}}) - c_{\text{reg}} D_{\text{KL}}(p(\mathbf{w}, \mathbf{\Gamma} | \mathbf{A}^{-1}, \mathbf{m}, c_{\text{temp}}) || q(\mathbf{W}), q(\mathbf{\Gamma}))$$

Параметрическая сложность

Обозначим за $F(c_{\text{reg}}, c_{\text{train}}, c_{\text{prior}}, c_{\text{comb}}, \mathbf{P}, c_{\text{temp}})$ множество экстремумов функции L при решении задачи двухуровневой оптимизации.

Утверждение

Пусть $\mathbf{f} \in F(1, 1, c_{\text{prior}}, 0, \emptyset, c_{\text{temp}})$. При устремлении c_{prior} к бесконечности параметрическая сложность модели \mathbf{f} устремляется к нулю.

$$\lim_{c_{\text{prior}} \rightarrow \infty} C_{\text{param}}(\mathbf{f}) = 0.$$

Параметрическая сложность

Обозначим за $F(c_{\text{reg}}, c_{\text{train}}, c_{\text{prior}}, c_{\text{comb}}, \mathbf{P}, c_{\text{temp}})$ множество экстремумов функции L при решении задачи двухуровневой оптимизации.

Утверждение

Пусть $\mathbf{f}_1 \in F(1, 1, c_{\text{prior}}^1, 0, \emptyset, c_{\text{temp}})$, $\mathbf{h}_2 \in F(1, 1, c_{\text{prior}}^2, 0, \emptyset, c_{\text{temp}})$, $c_{\text{prior}}^1 < c_{\text{prior}}^2$. Пусть вариационные параметры моделей \mathbf{f}_1 и \mathbf{f}_2 лежат в области U , в которой соответствующие функции L и Q являются локально-выпуклыми. Тогда модель \mathbf{f}_1 имеет параметрическую сложность, не меньшую чем у \mathbf{f}_2 .

$$C_{\text{param}}(\mathbf{f}_1) \geq C_{\text{param}}(\mathbf{f}_2).$$

Структурная сложность

Утверждение

Пусть для каждого ребра (i, j) семейства моделей \mathfrak{F} априорное распределение

$$p(\gamma_{i,j}) = \lim_{c_{\text{temp}} \rightarrow 0} \mathcal{GS}(c_{\text{temp}}).$$

Пусть $c_{\text{reg}} > 0$, $c_{\text{train}} > 0$, $c_{\text{prior}} > 0$. Пусть $\mathbf{f} \in F(c_{\text{reg}}, c_{\text{train}}, c_{\text{prior}}, 0, \emptyset, c_{\text{temp}})$. Тогда структурная сложность модели \mathbf{f} равняется нулю.

$$C_{\text{struct}}(\mathbf{f}) = 0.$$

Структурная сложность

Гипотеза

Пусть $\mathbf{f}_1 \in F(c_{\text{reg}}, c_{\text{train}}, c_{\text{prior}}, 0, \emptyset, c_{\text{temp}}^1)$, $\mathbf{h}_2 \in \lim_{c_{\text{temp}}^2 \rightarrow \infty} F(c_{\text{reg}}, c_{\text{train}}, c_{\text{prior}}, 0, \emptyset, c_{\text{temp}}^2)$. Пусть вариационные параметры моделей f_1 и f_2 лежат в области U , в которой соответствующие функции L и Q являются локально-выпуклыми. Тогда разница структурных сложностей моделей ограничена выражением:

$$C_{\text{struct}}(\mathbf{f}_1) - C_{\text{struct}}(\mathbf{f}_2) \leq E_q^1 \log p(y|\mathbf{X}, \mathbf{W}, \mathbf{\Gamma}, \mathbf{A}^{-1}, c_{\text{temp}}^1) - E_q^2 \log p(y|\mathbf{X}, \mathbf{W}, \mathbf{\Gamma}, \mathbf{A}^{-1}).$$

Полный перебор

Пусть для каждого ребра (i, j) семейства моделей \mathfrak{F} априорное распределение

$$p(\gamma_{i,j}) = \lim_{c_{\text{temp}} \rightarrow 0} \mathcal{GS}(c_{\text{temp}}).$$

Рассмотрим последовательность \mathbf{P} , состоящую из $N = \prod_{(j,k) \in E} K_{j,k}$ моделей, полученных в ходе оптимизаций вида:

$$f_1 \in F(c_{\text{reg}}, 0, 0, \emptyset, c_{\text{comb}}, c_{\text{temp}}),$$

$$f_2 \in F(c_{\text{reg}}, 0, 0, \{q_1(\Gamma)\}, c_{\text{comb}}, c_{\text{temp}}),$$

$$f_3 \in F(c_{\text{reg}}, 0, 0, \{q_1(\Gamma), q_2(\Gamma)\}, c_{\text{comb}}, c_{\text{temp}}),$$

где $c_{\text{reg}} > 0, c_{\text{comb}} > 0$.

Гипотеза

Вариационные распределения q_{Γ} структур последовательности \mathbf{P} вырождаются в распределения вида $\delta(\hat{\mathbf{m}})$, где $\hat{\mathbf{m}}$ — точка на декартовом произведении вершин симплексов структуры модели.

Последовательность соответствует полному перебору структуры Γ .