

Порождение моделей заданной сложности с использованием байесовских гиперсетей

О. С. Гребенькова

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем
Научный руководитель: Бахтеев Олег Юрьевич

Отчёт по НИР
Осень 2020 г.

Задача построения модели глубокого обучения

Цель

Предложить метод оптимизации модели глубокого обучения с контролем сложности модели.

Исследуемая проблема

По построению семейство моделей глубокого обучения имеет избыточное число параметров. Поэтому оптимизация и выбор модели с наперед заданной сложностью является вычислительно сложной задачей.

Метод решения

Предлагаемый метод заключается в представлении модели глубокого обучения в виде гиперсети, с использованием байесовского подхода. Гиперсеть — сеть, которая порождает параметры для оптимальной модели.

- ❶ выборка

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}, \quad i = 1, \dots, m, \quad \mathbf{x}_i \in \mathbb{R}^m, \quad y_i \in \{1, \dots, Y\},$$

- ❷ модель

$$\mathbf{f}(\mathbf{x}, \mathbf{w}) : \mathbb{R}^m \times \mathbb{R}^n \longrightarrow \mathbb{R}^Y,$$

где $\mathbf{w} \in \mathbb{R}^n$ — пространство параметров модели;

- ❸ априорное распределение вектора параметров в пространстве \mathbb{R}^n :

$$p(\mathbf{w}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{A}_{\text{pr}}^{-1}),$$

где $\boldsymbol{\mu}, \mathbf{A}_{\text{pr}}^{-1}$ — вектор средних и матрица ковариации априорного распределения;

- ❹ распределение, аппроксимирующее неизвестное апостериорное распределение $p(\mathbf{w}|\mathcal{D})$:

$$q(\mathbf{w}) \sim \mathcal{N}(\mathbf{m}, \mathbf{A}_{\text{ps}}^{-1}).$$

Здесь $\mathbf{m}, \mathbf{A}_{\text{ps}}^{-1}$ — вектор средних и матрица ковариации. Предполагается, что:

$$q(\mathbf{w}) \approx p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

Логарифмическая функция правдоподобия выборки:

$$\mathcal{L}_{\mathcal{D}}(\mathbf{w}|\mathcal{D}) \propto \log p(\mathcal{D}|\mathbf{w}).$$

Логарифм обоснованности модели:

$$\log p(\mathcal{D}) = \log \int_{\mathbf{w} \in \mathbb{R}^n} p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w}.$$

При оценке интеграла получаем:

$$\begin{aligned}\mathcal{L}(\mathcal{D}) = \log p(\mathcal{D}) &\geq \int_{\mathbf{w} \in \mathbb{R}^n} q(\mathbf{w}) \log \frac{p(\mathbf{w})}{q(\mathbf{w})} d\mathbf{w} + \int_{\mathbf{w} \in \mathbb{R}^n} q(\mathbf{w}) \log p(\mathcal{D}|\mathbf{w}) d\mathbf{w} = \\ &= \mathcal{L}_{\mathbf{w}}(\mathcal{D}, \mathbf{w}) + \mathcal{L}_E(\mathcal{D}).\end{aligned}$$

Обобщенная функция обоснованности: $\lambda \mathcal{L}_{\mathbf{w}}(\mathcal{D}, \mathbf{w}) - \mathcal{L}_E(\mathcal{D})$

Максимизация функционала

$$\mathfrak{G}(\lambda) = \arg \max_{\mathbf{w} \in \mathbb{R}^n} (\log p(\mathcal{D}|\mathbf{w}) - \lambda D_{KL}(q(\mathbf{w})||p(\mathbf{w}))).$$

Гиперсеть

Параметрическое отображение из множества Λ во множество параметров модели \mathbb{R}^n :

$$\mathbf{G} : \Lambda \times \mathbb{R}^u \rightarrow \mathbb{R}^n,$$

где \mathbb{R}^u — множество допустимых параметров гиперсети, Λ — множество параметров, контролирующих сложность модели.

Реализация с отображением во множество матриц низкого ранга

$$\mathbf{G}_{\text{lowrank}}(\lambda) = (\mathbf{f}(\lambda)\mathbf{U}_1)^\top (\mathbf{f}(\lambda)\mathbf{U}_2) + \mathbf{B}_1.$$

Реализация с линейной аппроксимацией

$$\mathbf{G}_{\text{linear}}(\lambda) = \lambda \mathbf{b}_2 + \mathbf{b}_3.$$

Дополнительно была испробована локальная репараметризация параметров. Были рассмотрены другие виды гиперсетей, пока линейная аппроксимация показывает себя лучше всего.

Цель

Исследовать поведение обобщенной функции обоснованности модели.
Сравнить методы построения разных моделей. Сравнить с теоретическими результатами.

Проведено сравнение следующих моделей:

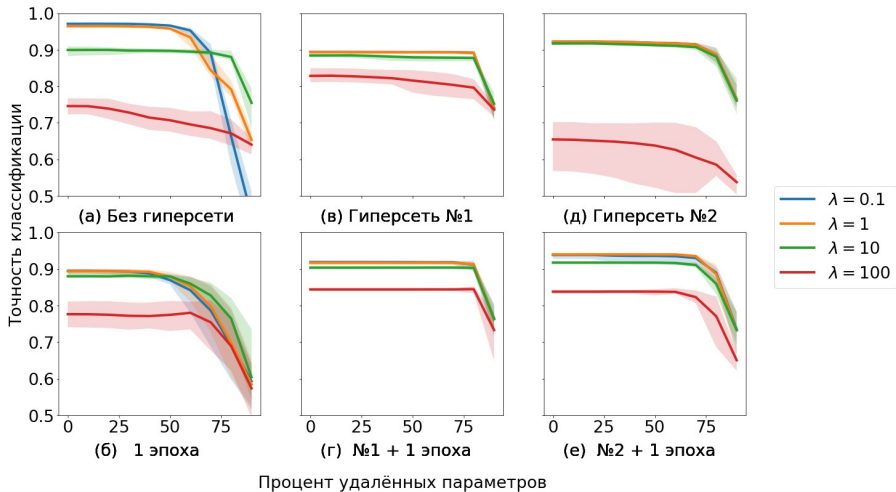
- (а) построения модели напрямую без использования гиперсети;
- (б) построения модели напрямую без использования гиперсети с оптимизацией за одну эпоху;
- (в) построение с использованием гиперсети;
- (г) построение с использованием гиперсети с дообучением итоговой модели за одну эпоху;
- (д) построение с использованием гиперсети;
- (е) построение с использованием гиперсети с дообучением итоговой модели за одну эпоху.
- (ж) построение модели без вариационного вывода
- (з) построение модели без вариационного вывода с использованием гиперсетей

Критерий удаления параметров — относительная плотность модели:

$$\rho(w_i) \approx \exp \frac{\mu_i^2}{2\sigma_i^2}.$$

Критерии качества модели:

- ❶ Точность классификации — Accuracy = $1 - \frac{1}{m} \sum_{i=1}^m [f(\mathbf{x}_i, \mathbf{w}) \neq y_i]$.
- ❷ Количество обновлений параметров модели.
- ❸ Обобщенная обоснованность модели.
- ❹ Стабильность модели S .



- ❶ Вариационный метод позволяет удалить $\approx 60\%$ параметров при всех λ без значительной потери точности классификации.
- ❷ Несмотря на потерю в качестве, гиперсеть получает схожие результаты, что и обычные модели при меньших вычислительных затратах.
- ❸ По графикам видно, что модель сохраняет схожие свойства (к примеру точность классификации) при прореживании.
- ❹ Планируется исследовать свойства вариационных гиперсетей и их применение к задаче выбора модели с контролем сложности.
- ❺ Сейчас ведется работа на базовым экспериментом для сравнения и теоретической частью диплома

ALEX GRAVES

Practical Variational Inference for Neural Networks // Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain

DAVID HA AND ANDREW M. DAI AND QUOC V. LE

HyperNetworks // CoRR, vol. abs/1609.09106, 2018.

TOM VENIAT AND LUDOVIC DENOYER

Learning Time/Memory-Efficient Deep Architectures With Budgeted Super Networks // CVPR, 2018, Pp. 3492–3500.

JONATHAN LORRAINE AND DAVID DUVENAUD

Stochastic Hyperparameter Optimization through Hypernetworks // CoRR, vol. abs/1802.09419, 2018.