

# 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:** к.ф.-м.н. Бахтеев Олег Юрьевич

Весна 2021 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.

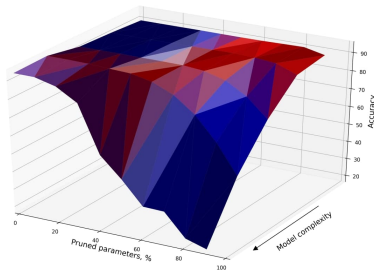


График зависимости точности классификации от процента удалённых параметров и параметра сложности модели

- ① выборка

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

- ② 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;

- ③ 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;

- ④ 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:

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

$$\mathcal{L}_{\mathbf{w}}(\mathcal{D}, \mathbf{w}) = -D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w})).$$

$$\mathcal{L}_E = \mathbb{E}_{q(\mathbf{w})} \mathcal{L}_{\mathcal{D}}(\mathcal{D}|\mathbf{w})$$

Обобщенная функция обоснованности:

$$\mathfrak{L}(\lambda) = \lambda \mathcal{L}_{\mathbf{w}}(\mathfrak{D}, \mathbf{w}) + \mathcal{L}_E(\mathfrak{D}); \quad (1)$$

$$\mathfrak{L}(\lambda) = \mathcal{L}_{\mathbf{w}}(\mathfrak{D}, p(\lambda \mathbf{w}), q(\mathbf{w})) + \mathcal{L}_E(\mathfrak{D}); \quad (2)$$

$$\mathfrak{L}(\lambda) = \lambda \|\mathbf{w}\|^2 + \mathcal{L}_E(\mathfrak{D}). \quad (3)$$

## Максимизация функционала

$$\mathfrak{G}(\lambda) = \arg \max_{\mathbf{w} \in \mathbb{R}^n} (\log p(\mathfrak{D}|\mathbf{w}) - \lambda D_{KL}(q(\mathbf{w})||p(\mathbf{w}))), \quad (4)$$

$$\mathfrak{G}(\lambda) = \arg \max_{\mathbf{w} \in \mathbb{R}^n} (\log p(\mathfrak{D}|\mathbf{w}) - D_{KL}(q(\mathbf{w})||p(\lambda \mathbf{w}))), \quad (5)$$

$$\mathfrak{G}(\lambda) = \arg \max_{\mathbf{w} \in \mathbb{R}^n} (\log p(\mathfrak{D}|\mathbf{w}) - \lambda \|\mathbf{w}\|^2). \quad (6)$$

## Гиперсеть

Параметрическое отображение из множества  $\Lambda$  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,  $\Lambda$  — множество параметров, контролирующих сложность модели.

## Реализация с линейной аппроксимацией

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

Для аппроксимации оптимизационной задачи (4) предлагается оптимизировать параметры гиперсети  $\mathbf{U} \in \mathbb{R}^u$  по случайно порожденным значениям параметра сложности  $\lambda \in \Lambda$ :

$$\mathbb{E}_{\lambda \sim P(\lambda)} (\log p(\mathcal{D}|\mathbf{w}) - \lambda D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}))) \rightarrow \max_{\mathbf{U} \in \mathbb{R}^u}.$$

**Критерий удаления параметров** — относительная плотность модели:

$$\rho(w_i) \propto \exp \frac{\mu_i^2}{2\sigma_i^2},$$

$$\rho(w_i) \propto \exp(-w_i^2).$$

## Теорема

Пусть выполнены следующие условия:

- 1 существует компакт  $\mathbb{U} \in \mathbb{R}^n$ , который содержит единственный минимум  $\mathbf{f}(\mathbf{w}^*(\lambda))$  для каждого  $\lambda \in \Lambda$ ;
- 2 существует последовательность моделей  $\mathbf{f}(\mathbf{w}_n(\lambda))$  такая, что  $\mathbb{E}\mathcal{L}(\mathbf{f}(\mathbf{w}_n(\lambda))) \xrightarrow{n \rightarrow \infty} \max.$

Тогда  $g(\mathbf{f}(\mathbf{w}_n(\lambda))) \xrightarrow{L_1} g(\mathbf{f}(\mathbf{w}^*(\lambda)))$ , где  $g$  это непрерывный критерий для удаления параметров.



## Цель

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

Проведено сравнение следующих моделей:

- (а) вариационная сеть;
- (б) сеть с репараметризацией;
- (в) базовая нейросеть с регуляризатором;
- (г) вариационная сеть с линейной гиперсетью;
- (д) сеть с репараметризацией и линейной гиперсетью ;
- (е) базовая нейросеть с регуляризатором и линейной гиперсетью.

Вид используемой нейросети для эксперимента на MNIST:

$$\mathbf{f}(\mathbf{x}, \mathbf{w}) = \text{softmax}(\mathbf{w}_2^\top \text{ReLU}(\mathbf{w}_1^\top \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2),$$

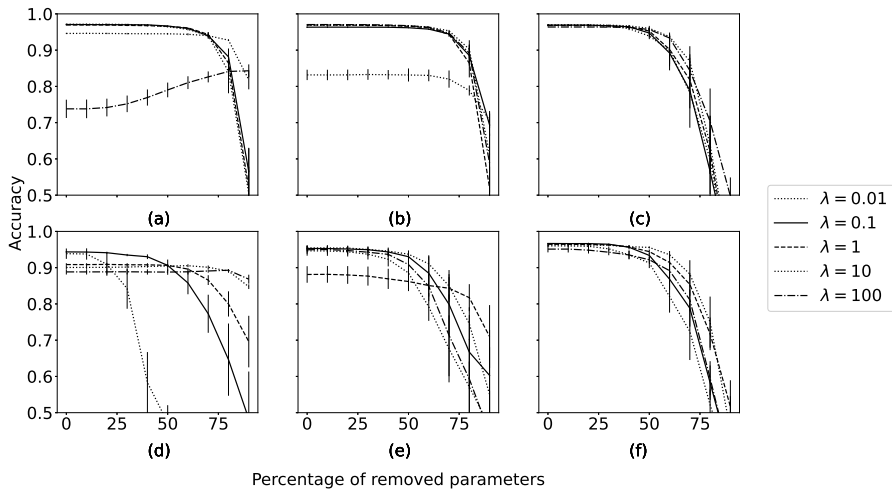
где  $\mathbf{w}_1, \mathbf{b}_1$  — параметры первого слоя,  $\mathbf{w}_2, \mathbf{b}_2$  — параметры второго слоя,

$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^k \exp(x_j)} \quad i = 1, \dots, k,$$

$$\text{ReLU}(\mathbf{x}) = \max(0, \mathbf{x}).$$

Критерий качества модели — точность классификации

$$\text{Accuracy} = 1 - \frac{1}{m} \sum_{i=1}^m [f(\mathbf{x}_i, \mathbf{w}) \neq y_i]$$



- ❶ Вариационный метод позволяет удалить  $\approx 60\%$  параметров при всех  $\lambda$  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 Вариационная сеть запущена на CIFAR, идет настройка эксперимента. Планируется провести полноценное сравнение на CIFAR и добавить описание эксперимента в дипломную работу.



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.