

Байесовский выбор наиболее правдоподобной структуры модели глубокого обучения

О. Ю. Бахтеев

Научный руководитель: д.ф.-м.н. В.В. Стрижов
Московский физико-технический институт (государственный университет)

Интеллектуализация обработки информации
ИОИ-2018
11.10.2018

Выбор структуры модели глубокого обучения

Цель работы:

Развитие теории байесовского выбора модели и исследование свойств методов выбора моделей глубокого обучения.

Задачи:

- Предложить алгоритм оптимизации параметров, гиперпараметров и структурных параметров моделей глубокого обучения.
- Предложить метод выбора модели наиболее правдоподобной структуры.
- Исследовать свойства оптимизационных алгоритмов выбора модели.

Основные проблемы

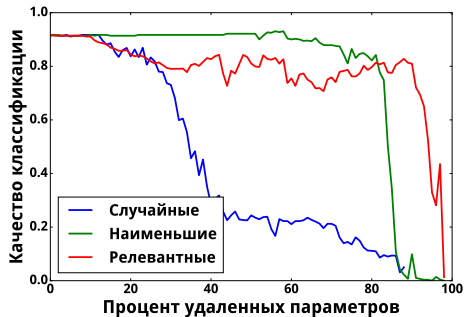
- Многоэкстремальность задачи оптимизации параметров модели.
- Вычислительная сложность оптимизации.
- Большое число параметров и гиперпараметров.

Исследование основывается на следующих работах

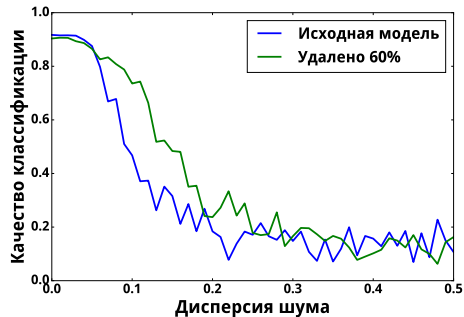
- Graves A. Practical variational inference for neural networks //Advances in Neural Information Processing Systems. – 2011
- Maclaurin D., Duvenaud D., Adams R. Gradient-based hyperparameter optimization through reversible learning //International Conference on Machine Learning. – 2015
- Luketina J. et al. Scalable gradient-based tuning of continuous regularization hyperparameters //International Conference on Machine Learning. - 2016
- J. Fu et al., DrMAD: Distilling Reverse-Mode Automatic Differentiation for Optimizing Hyperparameters of Deep Neural Networks // IJCAI - 2016
- Pedregosa F. Hyperparameter optimization with approximate gradient //International Conference on Machine Learning. – 2016. –

Проблемы оптимизации моделей глубокого обучения

Правдоподобие моделей с избыточным количеством параметров не меняется при удалении параметров.



Избыточность параметров модели



Неустойчивость модели

Глубокое обучение предполагает оптимизацию моделей с заведомо избыточной сложностью.

Задача выбора структуры модели

Однослойная нейросеть:

$$\mathbf{f}(\mathbf{x}) = \text{softmax} \left(\mathbf{W}^T \mathbf{f}_1(\mathbf{x}) \right), \quad \mathbf{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow 2^{|\mathbb{Y}|}, \quad \mathbf{x} \in \mathbb{R}^n.$$

$$\mathbf{f}_1(\mathbf{x}) = \gamma_{0,1}^1 \mathbf{g}_{0,1}^1(\mathbf{x}) + \dots + \gamma_{0,1}^K \mathbf{g}_{0,1}^K(\mathbf{x}) = \gamma_{0,1}^1 \boldsymbol{\sigma}$$

где $\mathbf{W}_1, \dots, \mathbf{W}_K$ — матрицы параметров, $\{\mathbf{g}_{0,1}^i\}_{i=1}^K$ — базовые функции для скрытого слоя нейросети.

Структурные параметры: $\boldsymbol{\Gamma} = [\gamma_{0,1}]$.

Структура модели определяется вершиной булевого K -мерного куба.

Задача выбора структуры модели: два скрытых слоя

Двухслойная нейросеть:

$$\mathbf{f}(\mathbf{x}) = \text{softmax} \left(\mathbf{W}^T \mathbf{f}_2(\mathbf{x}) \right), \quad \mathbf{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow 2^{|\mathbb{Y}|}, \quad \mathbf{x} \in \mathbb{R}^n.$$

$$\mathbf{f}_2(\mathbf{x}) = \gamma_{1,2}^1 \mathbf{g}_{1,2}^1(\mathbf{f}_1(\mathbf{x})) + \cdots + \gamma_{1,2}^K \mathbf{g}_{1,2}^K(\mathbf{f}_1(\mathbf{x})) = \gamma_{1,2}^1 \sigma(\mathbf{W}_{K+1}^T \mathbf{f}_1(\mathbf{x})) + \cdots + \gamma_{1,2}^K \sigma(\mathbf{W}_{2K}^T \mathbf{f}_1(\mathbf{x})),$$

$$\mathbf{f}_1(\mathbf{x}) = \gamma_{0,1}^1 \mathbf{g}_{0,1}^1(\mathbf{x}) + \cdots + \gamma_{0,1}^K \mathbf{g}_{0,1}^K(\mathbf{x}) = \gamma_{0,1}^1 \sigma(\mathbf{W}_1^T \mathbf{x}) + \cdots + \gamma_{0,1}^K \sigma(\mathbf{W}_K^T \mathbf{x}),$$

где $\mathbf{W}_1, \dots, \mathbf{W}_{2K}$ — матрицы параметров, $\{\mathbf{g}_{0,1}^i, \mathbf{g}_{1,2}^i\}_{i=1}^K$ — базовые функции для скрытых слоев нейросети.

Структурные параметры: $\Gamma = [\gamma_{0,1}, \gamma_{1,2}]$.

Структура модели определяется вершинами **двух** булевых K -мерных кубов.

Картинка: связь параметров со структурой.

Семейство моделей

Задан граф V, E .

Для каждого ребра $(j, k) \in E$ определен вектор **базовых функций** $\mathbf{g}_{j,k}$ мощностью $K_{j,k}$.
Граф V, E со множеством функций $\mathcal{G} = \{\mathbf{g}_{j,k}\}_{(j,k) \in E}$ называется **семейством моделей**,
если функция, задаваемая рекурсивно как

$$\mathbf{f}_j(\mathbf{x}) = \sum_{k \in \text{Adj}(v_j)} \langle \gamma_{j,k}, \mathbf{g}_{j,k} \rangle (f_k(\mathbf{x})), \quad \mathbf{f}_0(\mathbf{x}) = \mathbf{x},$$

является непрерывной дифференцируемой функцией из \mathbb{R}^n во множество \mathbb{Y} при любых значениях векторов γ .

Модель определяется параметрами подмоделей f_j и структурными параметрами γ .

Обозначим за вектор **параметров модели** \mathbf{W} конкатенацию параметров всех подмоделей $\{f_j\}_{j=1}^{|V|}$.

Обозначим за вектор **параметров модели** $\mathbf{\Gamma}$ конкатенацию структурных параметров γ .

Критерии качества модели

Точность S модели $f(x)$ — величина ошибки на контрольной выборке.

Устойчивость модели $f(x)$ — число обусловленности матрицы \mathbf{A} :

$$\eta(\mathbf{W}) = \frac{\lambda_{\max}}{\lambda_{\min}} \quad \text{при гипотезе } \mathbf{W} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}),$$

λ_{\max} — максимальное, а λ_{\min} — минимальное собственные числа матрицы \mathbf{A} .

Статистические критерии качества модели

Статистическая сложность модели \mathbf{f} :

$$\text{MDL}(\mathbf{y}, \mathbf{f}) = -\log p(\mathbf{f}) - \log (p(\mathbf{y}|\mathbf{X}, \mathbf{f})\delta\mathfrak{D}),$$

где $\delta\mathfrak{D}$ — допустимая точность передачи информации о выборке \mathfrak{D} .

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

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

где \mathbf{m} — гиперпараметры априорного распределения структуры модели.

TODO: надо ли про bits-back?

Структурная сложность модели — энтропия апостериорного распределения структуры модели:

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

В данной работе предлагается метод оптимизации модели, учитывающий все перечисленные критерии качества модели.

Правдоподобие модели

Пусть заданы априорное распределение параметров и структуры $p(W, \Gamma)$.

Модель f оптимальна, если достигается максимум правдоподобия модели:

$$p(y|X) = \int_{W, \Gamma} p(y|X, W, \Gamma) p(W, \Gamma) dW d\Gamma.$$

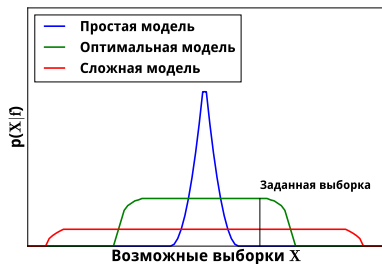
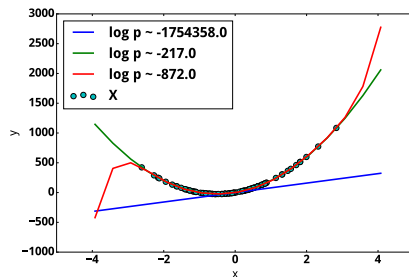


Схема выбора модели по правдоподобию



Пример: полиномы

Статистические критерии качества модели

Выбор оптимальной модели

Основные проблемы выбора оптимальной модели

- Интеграл правдоподобия невычислим аналитически.
- Задача оптимизации многоэкстремальна и невыпукла.

Требуется

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

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

Вариационная нижняя оценка правдоподобия

Интеграл правдоподобия невычислим аналитически.

Правдоподобие модели:

$$p(\mathbf{y}|\mathbf{X}) = \int_{\mathbf{W}, \mathbf{\Gamma}} p(\mathbf{y}|\mathbf{X}, \mathbf{W}, \mathbf{\Gamma}) p(\mathbf{W}, \mathbf{\Gamma}) d\mathbf{W} d\mathbf{\Gamma}.$$

Получим нижнюю оценку интеграла правдоподобия.

Пусть q — непрерывное распределение.

$$\log p(\mathbf{y}|\mathbf{X}) \geq \mathbb{E}_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{W}, \mathbf{\Gamma}) - D_{\text{KL}}(p(\mathbf{w}, \mathbf{\Gamma}) || q(\mathbf{W}, \mathbf{\Gamma})) = \log \hat{p}_q(\mathbf{y}|\mathbf{X}).$$

Полученная оценка совпадает с интегралом правдоподобия при

$$D_{\text{KL}}(q(\mathbf{W}, \mathbf{\Gamma}) || (p(\mathbf{W}, \mathbf{\Gamma}))) = 0.$$

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 $\mathbf{g}(j, k)$.

Будем считать, что вектор $\gamma_{j,k}$ распределен по распределению Gumbel-Softmax(GS):

$$p(\gamma) = (K_{j,k} - 1)! c_{\text{temp}}^{K_{j,k}-1} \left(\prod_{h=1}^{K_{j,k}-1} \alpha_h \gamma_h^{-c_{\text{temp}}-1} \right) \left(\sum_{h=1}^u \alpha_h \gamma_h^{-c_{\text{temp}}} \right),$$

где $\mathbf{m} = [\alpha_1, \dots, \alpha_{K_{j,k}}]$ — параметры сдвига распределения, c_{temp} — температура распределения.

Вариационный вывод: распределение структурных параметров

Для реализации h -й компоненты случайной величины γ справедлива следующая формула:

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

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

Для всех элементов структуры $\gamma_{j,k}$ положим:

$$p(\gamma_{j,k}) = GS([\alpha_1, \dots, \alpha_{K_{j,k}}, c_{\text{temp}}]), \quad q(\gamma_{j,k}) = GS([\hat{\alpha}_1, \dots, \hat{\alpha}_{K_{j,k}}, \hat{c}_{\text{temp}}]),$$

где $\hat{\alpha}_1, \dots, \hat{\alpha}_{K_{j,k}}, \hat{c}_{\text{temp}}$ — параметры вариационного распределения.

Вариационный вывод: распределение параметров

Пусть $\mathbf{W} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1})$ и структура модели Γ определена однозначно.

Пусть $q_{\mathbf{W}} = \mathcal{N}(\boldsymbol{\mu}_q, \mathbf{A}_q^{-1})$, $\boldsymbol{\theta} = [\boldsymbol{\mu}_q, \mathbf{A}_q^{-1}]$.

Тогда вариационная оценка имеет вид:

$$\int_{\mathbf{W}} q(\mathbf{W}) \log p(\mathcal{D}, \mathbf{W}, \mathbf{A}^{-1}) d\mathbf{W} - D_{\text{KL}}(q_{\mathbf{W}}(\mathbf{W}) || p(\mathbf{W} | \mathbf{A}^{-1})) \simeq$$
$$\sum_{i=1}^m \log p(\mathbf{x}_i | \mathbf{W}_i) - D_{\text{KL}}(q_{\mathbf{W}}(\mathbf{W}) || p(\mathbf{W} | \mathbf{A}^{-1})) = -L(\boldsymbol{\theta}, \mathbf{A}^{-1}, \mathcal{D}),$$

где $\mathbf{W}_i \sim q_{\mathbf{W}}$.

Дивергенция $D_{\text{KL}}(q_{\mathbf{W}}(\mathbf{w}) || p(\mathbf{W} | \mathbf{A}^{-1}))$ вычисляется аналитически:

$$D_{\text{KL}}(q_{\mathbf{W}}(\mathbf{W}) || p(\mathbf{W} | \mathbf{A}^{-1})) = \frac{1}{2} (\text{tr}(\mathbf{A} \mathbf{A}_q^{-1}) + \boldsymbol{\mu}_q^T \mathbf{A} \boldsymbol{\mu}_q - n + \ln |\mathbf{A}^{-1}| - \ln |\mathbf{A}_q^{-1}|).$$

Вариационная оценка на основе мултистарта

$$\log p(y|\mathbf{X}, \mathbf{A}) \geq \mathbb{E}_{q(\mathbf{w})} \log p(y, \mathbf{W}|\mathbf{X}, \mathbf{A}^{-1}) - \mathbb{E}_{q_{\mathbf{w}}}(-\log(q_{\mathbf{w}})).$$

Теорема [Бахтеев, 2016]. Пусть L — функция потерь, градиент которой — непрерывно-дифференцируемая функция с константой Липшица C . Пусть $\theta = [\mathbf{W}^1, \dots, \mathbf{W}^k]$ — начальные приближения оптимизации модели. Пусть β — шаг градиентного спуска, такой что:

- $\beta < \frac{1}{C}$,
- $\beta^{(-1)} > \max_{r \in \{1, \dots, k\}} \lambda_{\max}(\mathbf{H}(\mathbf{W}^r))$.

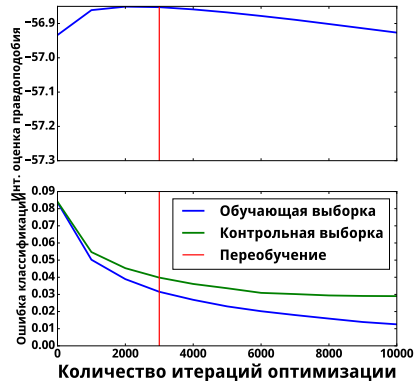
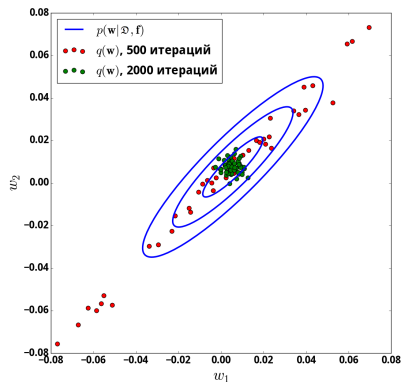
Тогда

$$\mathbb{E}_{q_{\mathbf{w}}^{\tau}}(-\log(q_{\mathbf{w}}^{\tau})) - \mathbb{E}_{q_{\mathbf{w}}^{\tau-1}}(-\log(q_{\mathbf{w}}^{\tau-1})) \sim \frac{1}{k} \sum_{r=1}^k (\beta \text{Tr}[\mathbf{H}(\mathbf{W}^r)] - \beta^2 \text{Tr}[\mathbf{H}(\mathbf{W}^r)\mathbf{H}(\mathbf{w}^r)]) + o_{\beta \rightarrow 0}(1),$$

где \mathbf{H} — гессиан функции потерь L , $q_{\mathbf{w}}^{\tau}$ — распределение $q_{\mathbf{w}}^{\tau}$ в момент оптимизации τ .

Вариационная оценка с использованием градиентного спуска

Максимизация вариационной оценки эквивалентна минимизации $D_{\text{KL}}(q(\mathbf{W})||p(\mathbf{W}|\mathcal{D}, \mathbf{A}^{-1}))$.
Градиентный спуск не минимизирует $D_{\text{KL}}(q(\mathbf{W})||p(\mathbf{W}|\mathcal{D}, \mathbf{A}^{-1}))$.



Оптимизация параметров вариационного распределения

Оптимизацию параметров вариационного распределения $q(\mathbf{W}, \mathbf{\Gamma}) = q_{\mathbf{W}}(\mathbf{W})q_{\mathbf{\Gamma}}(\mathbf{\Gamma})$ будем проводить по следующему функционалу:

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

Теорема

Пусть $c_{\text{reg}} > 0$. Тогда функция L сходится по вероятности к вариационной нижней оценке правдоподобия для подвыборки \mathcal{D} мощностью $c_{\text{reg}} m$:

$$L \rightarrow^p c_{\text{reg}} m \log \hat{p}_q(y|\mathbf{X}).$$

Теорема

Для любых значений \mathbf{A}, \mathbf{m} и вариационных параметров $q_{\mathbf{W}}$ существует точка на вершинах симплексов структуры, определяющая распределение $q_{\mathbf{\Gamma}}$, что для любой точки внутри симплексов, определяющей распределение $q'_{\mathbf{\Gamma}}$ справедливо выражение:

$$\lim_{c_{\text{temp}} \rightarrow 0} \frac{\log \hat{p}'_q(y|\mathbf{X})}{\log \hat{p}_q(y|\mathbf{X})} = -\infty.$$

Оптимизация параметров априорного распределения

Оптимизацию параметров априорного распределения будем проводить по следующему функционалу:

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

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

Общая задача оптимизации

Общая задача оптимизации — двухуровневая:

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

где

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

Оператор оптимизации

Обозначим за \mathbf{h} гиперпараметры \mathbf{A}, \mathbf{m} .

Обозначим за θ параметры распределений $q_{\mathbf{W}}, q_{\Gamma}$.

Определение

Оператором T назовем оператор стохастического градиентного спуска, производящий η шагов оптимизации:

$$\hat{\theta} = T \circ T \circ \dots \circ T(\theta_0, \mathbf{A}^{-1}) = T^\eta(\theta_0, \mathbf{A}^{-1}), \quad (1)$$

где

$$T(\theta, \mathbf{A}^{-1}) = \theta - \beta \nabla L(\theta, \mathbf{A}^{-1})|_{\hat{\mathcal{D}}},$$

γ — длина шага градиентного спуска, θ_0 — начальное значение параметров θ , $\hat{\mathcal{D}}$ — случайная подвыборка исходной выборки \mathcal{D} .

Оптимизация гиперпараметров

Перепишем итоговую задачу оптимизации:

$$\hat{\mathbf{h}} = \arg \max_{\mathbf{h}} Q(T^{\eta}(\boldsymbol{\theta}_0, \mathbf{A}^{-1})),$$

где $\boldsymbol{\theta}_0$ — начальное значение $\boldsymbol{\theta}$.

Утверждение, Luketina et al., 2016

Пусть функции L и Q являются дважды дифференцируемыми и выпуклыми. Пусть гессиан функции L аппроксимируется единичной матрицей:

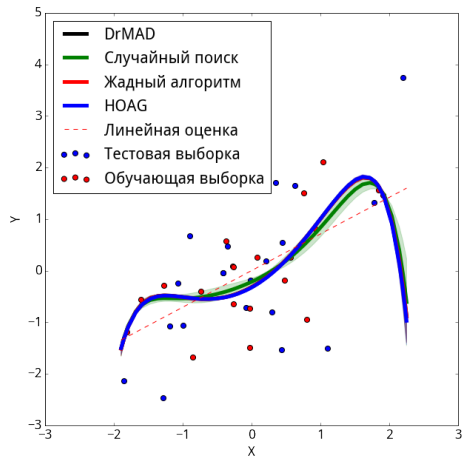
$$\mathbf{H}(L, \boldsymbol{\theta}) \approx \mathbf{I}.$$

Тогда допустима следующая оптимизация гиперпараметров:

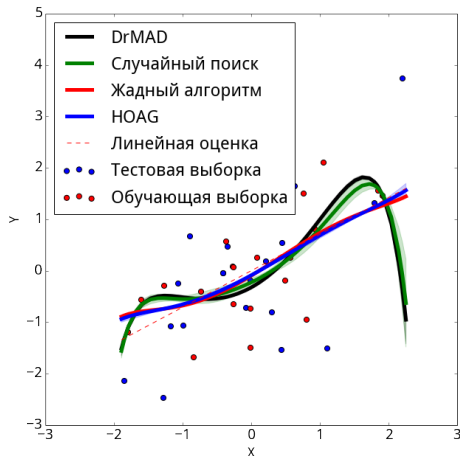
$$\mathbf{h}' = \mathbf{h} - \beta^h \nabla_{\mathbf{h}} Q(T(\boldsymbol{\theta}), \mathbf{h}),$$

где β^h — шаг оптимизации гиперпараметров.

Оптимизация гиперпараметров: пример



Кросс-Валидация



Вариационная оценка

Оптимизация правдоподобия модели

Теорема

Пусть существуют параметры распределения $q(\mathbf{W}, \mathbf{\Gamma})$, такие что $D_{KL}(q(\mathbf{W}, \mathbf{\Gamma})|p(\mathbf{W}, \mathbf{\Gamma}|\mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c_{temp})) = 0$.

Тогда двухуровневая задача оптимизация эквивалентна задаче оптимизации правдоподобия модели:

$$\arg \max_{\mathbf{A}, \mathbf{m}} p(\mathbf{y}|\mathbf{X}, \mathbf{A}, \mathbf{m}, c_{temp})$$

при $c_{reg} = c_{prior} = c_{train} > 0, c_{comb} = 0$.

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

Параметрическая сложность

TODO: тип сходимости

Теорема

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

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

Теорема

Пусть $\mathbf{f}_1 \in F(1, 1, c_{\text{prior}}, 0, \{\}, c_{\text{temp}})$, $\mathbf{f}_2 \in F(1, 1, c_{\text{prior}}, 0, \{\}, c'_{\text{temp}})$, $c_{\text{prior}} < c'_{\text{prior}}$.

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 большую чем у \mathbf{f}_2 .

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

Структурная сложность

Теорема

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

$$p(\gamma_{i,j}) = GS(\alpha_1, \dots, \alpha_{K_{i,j}}, c_{\text{temp}}).$$

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

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

Теорема

Пусть $\mathbf{f}_1 \in F(c_{\text{reg}}, c_{\text{train}}, c_{\text{prior}}, 0, \{\}, c_{\text{temp}})$, $\mathbf{f}_2 \in \lim_{c_{\text{temp}}' \rightarrow \infty} F(c_{\text{reg}}, c_{\text{train}}, c_{\text{prior}}, 0, \{\}, c_{\text{temp}}')$. 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:

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

Полный перебор

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

$$p(\gamma_{i,j}) = GS(\alpha_1, \dots, \alpha_{K_{i,j}}, c_{\text{temp}}).$$

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

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

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

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

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

Теорема

Вариационные распределения структур q_Γ последовательности вырождаются в распределения вида $\delta(\hat{\Gamma})$, где $\hat{\Gamma}$ — точка на декартовом произведении вершин симплексов структуры модели. Вариационные распределения последовательности проходят все возможные комбинации структур модели.

Заключение

- Предложен алгоритм оптимизации параметров, гиперпараметров и структурных параметров моделей глубокого обучения.
- Предложен метод выбора модели наиболее правдоподобной структуры.
- Проведено исследование свойства оптимизационных алгоритмов выбора модели.