

# 1 Аннотация

В работе рассматривается задача выбора структуры модели глубокого обучения. Модель — это вычислительный вероятностный граф, т.е. граф, в котором ребрами выступают нелинейные функции, а вершинами — результаты действия функцией на выборку. Каждому ребру поставлено в соответствие множество нелинейных функций, такое что линейная комбинация этих функций определяет дифференцируемую функцию заданной сигнатуры. Структурой модели назовем веса линейной комбинации этих функций.

Для нахождения оптимальной структуры предлагается ввести вероятностную интерпретацию модели, т.е. предположения о распределениях параметров и структуры модели. Проводится градиентная оптимизация параметров и гиперпараметров модели на основе байесовского вариационного вывода. Решается двухуровневая задача оптимизации: на первом уровне проводится оптимизация нижней оценки правдоподобия модели по вариационным параметрам модели. На втором уровне проводится оптимизация гиперпараметров модели. В качестве оптимизируемой функции для гиперпараметров модели предлагается обобщенная функция правдоподобия. Показано, что данная функция позволяет проводить оптимизацию несколькими алгоритмами: последовательным добавлением и удалением параметров, полным перебором, а также максимизацией нижней оценки правдоподобия модели.

Проводится сравнение с эвристическими алгоритмами выбора структуры модели. Вычислительный эксперимент проводится на синтетических данных и выборке рукописных цифр MNIST.

**Цель работы:** предложить метод выбора модели субоптимальной сложности, позволяющий проводить выбор модели в нескольких режимах (ELBO, AddDel, полный перебор, оптимизация без регуляризации и с регуляризацией).

## 2 Постановка задачи

Задана выборка

$$\mathfrak{D} = \{(\mathbf{x}_i, y_i)\}, i = 1, \dots, m, \quad (1)$$

состоящая из множества пар «объект-метка»

$$\mathbf{x}_i \in \mathbf{X} \subset \mathbb{R}^n, \quad y_i \in \mathbf{y} \subset \mathbb{Y}.$$

Метка  $y$  объекта  $\mathbf{x}$  принадлежит либо множеству:  $y \in \mathbb{Y} = \{1, \dots, Z\}$  в случае задачи классификации, где  $Z$  — число классов, либо некоторому подмножеству вещественных чисел  $y \in \mathbb{Y} \subseteq \mathbb{R}$  в случае задачи регрессии. Далее будем полагать, что объекты  $\mathbf{x}$  являются реализацией некоторой случайно величины и порождены независимо.

Определим множество архитектур моделей глубокого обучения для дальнейшего выбора оптимальной. Будем рассматривать модель как граф  $V, E$ , ребрами  $E$  которого являются функции из заданного множества функций  $\mathbf{O}$ , действующие на выборку,

а вершинами  $V$  — промежуточные представления выборки под действием данных функций. Перейдем к формальному определению модели. Пусть задан граф  $V, E$ . Пусть для каждого ребра  $\langle i, j \rangle \in E$  определено множество функций  $\mathbf{o}(i, j)$ . Граф  $V, E$  с множеством функций  $\mathbf{O}$  называется моделью, если функция, задаваемая рекурсивно как

$$f_i(\mathbf{x}) = \sum_{j \in \text{Adj}(v_i)} o(i, j)(f_j(\mathbf{x})),$$

является непрерывной дифференцируемой функцией из  $\mathbb{R}^n$  во множество  $\mathbb{Y}$  при любом  $o(i, j)$ , являющемся выпуклой комбинацией функций из множества  $\mathbf{o}(i, j)$ .

Пусть для каждого ребра  $i, j$  задан нормированный положительный вектор  $\gamma_{i,j} \in \mathbb{R}^{|\mathbf{o}(i,j)|}$ , определяющий веса функций из множества  $\mathbf{o}(i, j)$ . Будем считать, что вектор  $\gamma_{i,j}$  распределен по распределению Gumbel-Softmax:

$$\gamma_{i,j} \sim \text{GS}(c, \mathbf{m}_{i,j}).$$

где  $c$  — вектор концентрации распределения,  $\mathbf{m}_{i,j}$  — вектор средних. Обозначим за структуру модели  $\mathbf{\Gamma}$  множество всех векторов  $\gamma$ .

Для более адекватной оптимизации положим априорное распределение параметров модели  $\mathbf{w}$  зависящим от структуры модели  $\mathbf{\Gamma}$ . Пусть задано соответствие  $S$  между каждым параметром  $\mathbf{w} \in \mathbf{W}$  и параметром структуры  $\gamma_i$ . Положим распределение параметров  $\mathbf{w}$  нормальным с нулевым средним и диагональной ковариационной матрицей:

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1} \cdot S).$$

Пусть также определено правдоподобие выборки  $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma})$ .

**Определение** Правдоподобием модели  $\mathbf{f}$  назовем следующее выражение:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{A}, \mathbf{m}, c) = \int_{\mathbf{w}, \mathbf{\Gamma}} p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}) p(\mathbf{w}|\mathbf{A}) p(\mathbf{\Gamma}|\mathbf{m}, c) d\mathbf{w} d\mathbf{\Gamma}. \quad (2)$$

Пусть задано значение концентрации  $c$ . Требуется найти гиперпараметры модели  $\mathbf{A}, \mathbf{m}$  доставляющие максимум правдоподобия модели:

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

**Теорема.** При  $c \ll 0$  оптимизация (2) эквивалентна дискретной оптимизации:  $\gamma_{ij} \in \{0, 1\}, \|\gamma_i\| = 1$ .

**Доказательство.** При  $c \rightarrow 0$  плотность вероятности сконцентрирована на вершинах симплекса (из оригинальной статьи). Т.к. функция вероятности ограничена, то по теореме Лебега допустим предельный переход к распределению на вершинах.

### 3 Вариационная постановка задачи

В общем виде вычисление значения интеграла (2) является вычислительно сложной процедурой. В качестве приближенного значения интеграла будем использовать вариационную верхнюю оценку правдоподобия модели. Пусть заданы непрерывные параметрические распределения  $q_w, q_\gamma$ , аппроксимирующие апостериорные распределения  $p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)$ ,  $p(\Gamma|\mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)$ . Тогда верно следующее выражение:

$$\log p(\mathbf{y}|\mathbf{X}, \mathbf{A}, \mathbf{m}, c) \leq \mathbb{E}_{q_w, q_\gamma} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma, \mathbf{A}, \mathbf{m}, c) - D_{KL}(q_\gamma||p(\Gamma)) - D_{KL}(q_w||p(\mathbf{w})). \quad (3)$$

Разница между верхней оценкой (3) и правдоподобием модели (2) определяется дивергенцией между апостериорными распределениями  $p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)$ ,  $p(\Gamma|\mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)$  и вариационными распределениями  $q_w, q_\gamma$ .

Определим основные величины, которые характеризуют сложность модели.

**Определение** Параметрической сложностью  $C_w$  модели назовем наименьшую дивергенцию вариационных параметров при условии априорного распределения параметров:

$$C_w = \arg \min_{\mathbf{A}} D_{KL}(q|p).$$

(Примечание: кажется, здесь должны учитываться параметры и структура, т.к. в априорном распределении параметров зашита зависимость от структуры).

**Определение** Структурной сложностью  $C_\gamma$  модели назовем энтропию распределения структуры:

$$C_\gamma = \mathbb{E}_{q_\gamma} \log q_\gamma.$$

Сформулируем основные требования к оптимизационной задаче и оптимизируемым функционалам:

1. Оптимизируемые функции должны быть дифференцируемы.
2. Степень регуляризации структуры и параметров должна быть контролируемой.
3. Оптимизация должна приводить к максимуму вариационной оценки.
4. Оптимизация должна позволять калибровать параметрическую сложность модели
5. Оптимизация должна позволять калибровать структурную сложность модели.
6. Оптимизация должна позволять проводить полный перебор структуры.

Сформулируем задачу как двухуровневую задачу оптимизации. Обозначим за  $\boldsymbol{\theta}$  оптимизируемые на первом уровне величины. Обозначим за  $\mathbf{h}$  величины, оптимизируемые на втором уровне. Положим  $\boldsymbol{\theta}$  равным параметрам распределений  $q_w, q_\gamma$ . Положим  $\mathbf{h} = [\mathbf{A}, \mathbf{m}]$ .

Пусть  $L$  — приближенное значение вариационной оценки правдоподобия:

$$L = \beta \log p(\mathbf{y}|\hat{\mathbf{w}}, \hat{\Gamma}) - D_{KL}(q_\gamma||p(\Gamma)) - D_{KL}(q_w||p(\mathbf{w})),$$

где  $\hat{\mathbf{w}} \sim q_w$ ,  $\hat{\Gamma} \sim q_\gamma$ ,  $\beta$  — коэффициент, контролирующий степень регуляризации структуры и параметров.

**Теорема.** Пусть  $\beta \neq 1$ . Тогда функция  $L$  сходится по вероятности к вариационной нижней оценке правдоподобия для подвыборки  $\mathfrak{D}$  мощностью  $\beta m$ .

**Доказательство.** Рассмотрим произвольную подвыборку  $\hat{\mathfrak{D}}$  мощностью  $m_0$ . Верхняя оценка правдоподобия модели для подвыборки имеет вид:

$$\log p(\hat{\mathbf{y}}|\hat{\mathbf{X}}, \mathbf{A}, \mathbf{m}, c) \leq \mathbb{E}_{q_w, q_\gamma} \log p(\hat{\mathbf{y}}|\hat{\mathbf{X}}, \mathbf{w}, \Gamma, \mathbf{A}, \mathbf{m}, c) - D_{KL}(q_\gamma||p(\Gamma)) - D_{KL}(q_w||p(\mathbf{w})).$$

$$\log p(\hat{\mathbf{y}}|\hat{\mathbf{X}}, \mathbf{w}, \Gamma, \mathbf{A}, \mathbf{m}, c) = \sum_i \log p(\hat{\mathbf{y}}_i|\hat{\mathbf{x}}_i, \mathbf{w}, \Gamma, \mathbf{A}, \mathbf{m}, c) \xrightarrow{p}_{m \rightarrow \infty} m \mathbb{E} \log p(\mathbf{y}|\hat{\mathbf{x}}, \mathbf{w}, \Gamma, \mathbf{A}, \mathbf{m}, c).$$

Формула нижней оценки получается подстановкой.

Пусть  $Q$  — валидационная функция:

$$\begin{aligned} Q(c, c_1, c_2, c_3, \mathbf{p}) = & c_1 \log p(\mathbf{y}|\hat{\mathbf{w}}, \hat{\Gamma}) + c_2 [-D_{KL}(q_\gamma||p(\Gamma)) - D_{KL}(q_w||p(\mathbf{w}))] + \\ & + c_3 \sum_{p_k \in \mathbf{p}} D_{KL}(q_\gamma||p_k), \end{aligned}$$

где  $\mathbf{p}$  — заданные распределения на структурах,  $c_1, c_2, c_3$  — коэффициенты.

Сформулируем задачу поиска оптимальной модели как двухуровневую задачу.

$$\hat{\mathbf{h}} = \arg \max_{\mathbf{h} \in \mathbb{R}^h} Q(T^n(\boldsymbol{\theta}_0, \mathbf{h})), \quad (4)$$

где  $T$  — оператор оптимизации, решающий задачу оптимизации:

$$L(T^n(\boldsymbol{\theta}_0, \mathbf{h})) \rightarrow \max.$$

**Теорема.** Пусть  $D_{KL}(q_w||p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)) = 0$ ,  $D_{KL}(q_\gamma||p(\Gamma|\mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)) = 0$ , пусть  $c_1 = 1, c_2 = 1, c_3 = 0$ . Тогда оптимизация (4) эквивалентна оптимизации (2).

**Доказательство.** При соблюдении условий теоремы неравенство вариационной оценки превращается в равенство.

**Теорема.** Пусть  $c_1 = 1, c_3 = 0$ . При  $c_2 \rightarrow \infty$  параметрическая сложность  $C_w \rightarrow 0$ . **Доказательство (идея).** При  $C_w \rightarrow \infty$  априорное распределение  $p(\mathbf{w})$  будет стремиться к дельта-функции. Штраф за  $D_{KL}$  будет расти, параметры будут стремиться к нулю.

**Теорема.** Пусть  $c_1 = 1, c_3 = 0, c_2 > 0, c'_2 > c_2$  и логарифм распределения  $\mathbb{E}_{q_w, q_\gamma} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma, \mathbf{A}, \mathbf{m})$  является выпуклым. Пусть  $q(\mathbf{w}), q(\mathbf{w})'$  — распределения, полученные в результате оптимизаций с различными коэффициентами  $c_2, c'_2$  соответственно. Тогда  $C_w(q) \geq C_w(q')$ .

**Доказательство (не окончено).** Заметим, что  $q$  и  $q'$  можно выразить следующим образом:

$$q = \arg \max_{\hat{q}: L(\hat{q}, \mathbf{A}, \mathbf{\Gamma})} Q(\mathbf{A}, \mathbf{\Gamma}, c_2),$$

$$q' = \arg \max_{\hat{q}: L(\hat{q}, \mathbf{A}, \mathbf{\Gamma})} Q(\mathbf{A}, \mathbf{\Gamma}, c'_2).$$

Функция  $Q$  является выпуклой как сумма выпуклых функций. Отсюда справедливы следующие неравенства (по единственности точек экстремума):

$$\mathbb{E}_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) - c_2 D_{\text{KL}}(q||p) - \mathbb{E}_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) + c_2 D_{\text{KL}}(q'||p) > 0,$$

$$\mathbb{E}_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) - c'_2 D_{\text{KL}}(q'||p) - \mathbb{E}_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) + c_2 D_{\text{KL}}(q||p) > 0.$$

Вычитая неравенства получим:

$$D_{\text{KL}}(q||p) \geq D_{\text{KL}}(q'||p'),$$

$$\mathbb{E}_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) \leq \mathbb{E}_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c).$$

**Теорема** Пусть при любых выпуклых комбинациях функций  $\mathbf{O}$  модель является выпуклой. Пусть  $c_1, c_2 \geq 0, c_3 \leq 0$ . Пусть в качестве оценки правдоподобия выборки используется следующая:

$$\mathbb{E}_{q_w, q_\gamma} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) \sim \frac{1}{N} \sum_{k=1}^N \log p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}_i, \hat{\mathbf{\Gamma}}_i, \mathbf{A}, \mathbf{m}, c),$$

где  $\hat{\mathbf{w}}_i, \hat{\mathbf{\Gamma}}_i$  зафиксированные получены под действием вариационной репараметризации с зафиксированной реализацией исходных случайных величин. Тогда  $L$  и  $Q$  являются выпуклыми функциями. **Доказательство** По сумме выпуклых функций.

**Теорема (предварительно)** При  $c \rightarrow 0$   $C_\gamma \rightarrow 0$ .

**Теорема (предварительно)** При  $c \rightarrow \infty$   $C_\gamma \rightarrow \max$ .

**Утверждение (предварительно, нужно развить).** Пусть  $c_3 > 0, c < 0$  и все  $p_k \in \mathbf{p}$  отражают распределения на вершинах симплекса. Тогда оптимизация приведет к  $q_\gamma$ , сконцентрированному на одной из остальных вершин симплекса.

**Утверждение (очень предварительно).** Изменение  $c$  позволяет избежать ухода в локальный минимум.

**Утверждение (очень предварительно).** Изменение  $c_2$  позволяет избежать ухода в локальный минимум.

**Утверждение (очень предварительно).** Взаимосвязь структуры и параметров в prior позволяет получить «хорошие» модели.

**Утверждение (предварительно).** Пусть  $c_1 = c_2 = c_3 = 0$ . Пусть  $q_w \sim \mathcal{N}(\mathbf{0}, \sigma), \sigma \sim 0$ . Тогда оптимизация эквивалентна обычной оптимизации параметров с  $l_2$  - регуляризацией.

## 4 Вычислительный эксперимент

В качестве модельного эксперимента рассматривалась задача выбора модели линейной регрессии. Множество объектов  $\mathbf{X}$  было сгенерировано из трехмерного стандартного распределения:

$$\mathbf{X} \sim \mathcal{N}(0, \mathbf{I}), n = 3.$$

Множество меток было определено следующим правилом:

$$\mathbf{y} = \arg \max_{0,1} (\mathbf{X}_1 + \mathbf{X}_2),$$

третья компонента не участвовала в генерации ответа.

Рассматривались четыре возможные структуры:

1.  $f_1 = \mathbf{w}_1 \mathbf{X}_1$  (модель — регрессия только по первому признаку),
2.  $f_2 = \mathbf{w}_2 \mathbf{X}_2$  (модель — регрессия только по первому признаку),
3.  $f_3 = \mathbf{w}_3 \mathbf{X}_3$  (модель — регрессия только по шумовому признаку),
4.  $f_4 = \mathbf{w}_4 \mathbf{X}$  (модель — регрессия по всем признакам).

Ожидаемое поведение оптимизации:

1. При  $c_1 = c_2 = 1, c \sim 0$  (Evidence с низкой температурой) будет произведен выбор структуры  $f_4$ .
2. При  $c_1 = c_2 = 1, c \gg 0$  (Evidence с высокой температурой) будет произведен выбор двух структур с одинаковым весом:  $f_1, f_2$ .
3. При  $c_1 = c_2 = 0, c_3 = 1, \mathbf{p} = [[0.0, 0.0, 1.0, 0.0]], c \sim 0$  (Поощряется выбор структуры с шумовой компонентой) будет произведен выбор структуры  $f_4$ , при снижении параметра  $\beta$  выбор будет меняться в сторону  $f_3$ .

### Результаты

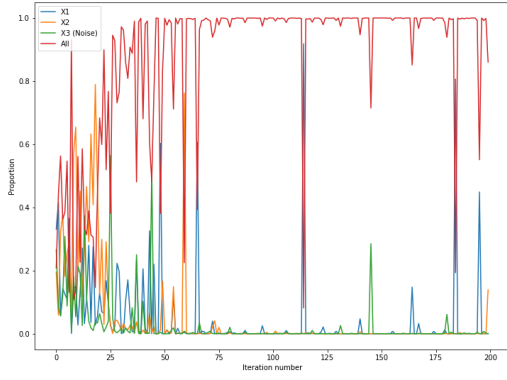


Figure 1: Evidence с низкой температурой

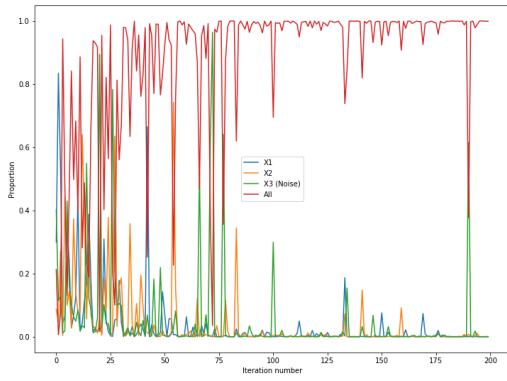


Figure 2: Evidence с высокой температурой

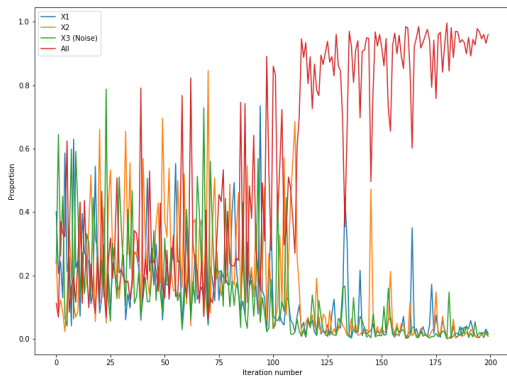


Figure 3: Evidence с высокой температурой,  $\beta = 0.01$

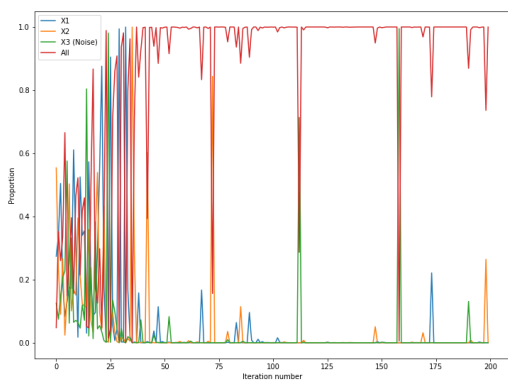


Figure 4: Поощрение выбора шумовой компоненты

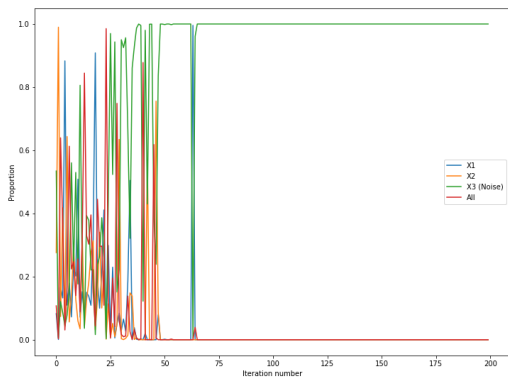


Figure 5: Поощрение выбора шумовой компоненты,  $\beta = 0.01$