

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 семейство моделей глубокого обучения для дальнейшего выбора оптимальной модели. Будем рассматривать семейство моделей как граф V, E .

Каждому ребру $(i, j) \in E$ сопоставим множество функций $\mathbf{g}^{i,j}$ мощности $K^{i,j}$. Вершины V — промежуточные представления выборки под действием данных функций.

Перейдем к формальному определению модели. Пусть задан граф V, E . Пусть для каждого ребра $(i, j) \in E$ определено множество функций $\mathbf{g}^{i,j}$. Граф V, E называется семейством моделей, если функция, задаваемая рекурсивно как

$$f_j(\mathbf{x}) = \sum_{k \in \text{Adj}(v_j)} \langle \gamma^{j,k}, \mathbf{g}^{j,k}(\mathbf{f}_k(\mathbf{x})) \rangle, \quad \mathbf{f}_0(\mathbf{x}) = \mathbf{x}$$

является дифференцируемой по параметрам функцией из \mathbb{R}^n во множество \mathbb{Y} при любых значениях векторов $\gamma^{j,k}$.

Параметрами модели \mathbf{W} будем называть конкатенацию всех параметров подмоделей \mathbf{f}_j . Структурой модели $\mathbf{\Gamma}$ будем называть конкатенацию всех структурных параметров $\gamma^{j,k}$. Моделью будем называть совокупность параметров \mathbf{W} и гиперпараметров \mathbf{W} .

Пусть все векторы $\gamma^{i,j}$ являются нормированными и положительными. Пусть для каждого структурного параметра $\gamma^{j,k} \in \mathbf{\Gamma}$ определено априорное Gumbel-softmax распределение $p(\gamma^{j,k} | \mathbf{m}^{j,k}, c_{\text{temp}})$ с параметром средних \mathbf{m} и температурой c_{temp} .

Пусть для структуры модели определено априорное распределение $p(\mathbf{\Gamma} | \mathbf{m})$, где \mathbf{m} — некоторое распределение. (нужно ли определять его явно здесь?) Пусть для каждого параметра $w \in \mathbf{W}$ определено множество $\mathcal{S}(w)$ структурных параметров γ , соответствующих базовым функциям \mathbf{g} , для которых определен этот параметр:

$$w \sim \mathcal{N}(\mathbf{0}, a^{-1} \cdot (\sum_{\gamma \in \mathcal{S}(w)} \gamma)),$$

где a — гиперпараметр, входящий в диагональную матрицу \mathbf{A}^{-1} . Пусть также определено правдоподобие выборки $p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \mathbf{\Gamma})$.

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

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

Требуется найти гиперпараметры модели \mathbf{A}, \mathbf{m} доставляющие максимум правдоподобия модели:

$$\arg \max_{\mathbf{A}, \mathbf{m}} p(\mathbf{y} | \mathbf{X}, \mathbf{A}, \mathbf{m}, c_{\text{temp}}), \quad (3)$$

а также соответствующие параметры и структур модели (см. вывод Байесва, первый уровень).

Докажем теорему о дискретности решения задачи нахождения оптимальных параметров модели.

Теорема Пусть $\mathbf{\Gamma}_1$ и $\mathbf{\Gamma}_2$ — реализации $\mathbf{\Gamma}$, такие что:

- $\mathbf{\Gamma}_1$ не содержит в себе точки внутри симплексов γ .

- Γ_2 содержит в себе точки внутри симплексов γ .

Тогда для любых положительно определенных матриц \mathbf{A}_1 и \mathbf{A}_2 и векторов $\mathbf{m}_1, \mathbf{m}_2$ справедлива следующая формула:

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

Доказательство. По теореме из оригинальной статьи

$$p\left(\lim_{c_{\text{temp}} \rightarrow 0} \mathbf{m} \text{ лежит на вершинах произведения симплексов}\right) = 1.$$

Тогда апостериорная вероятность Γ :

$$p(\Gamma_1 | \mathbf{y}, \mathbf{W}, \mathbf{X}, \mathbf{A}_1, \mathbf{m}_1, c_{\text{temp}}) \propto p(\Gamma) p(\mathbf{y} | \Gamma, \mathbf{W}, \mathbf{X}, \mathbf{A}_1, \mathbf{m})$$

будет стремиться к нулю при наличии точки внутри симплексов. Что и требовалось доказать.

TODO: еще бы хотелось расписать, что гамма должна в дискретном случае концентрироваться на одной вершине, но пока непонятно как сформулировать.

3 Вариационный вывод

В общем виде вычисление значения интеграла (2) является вычислительно сложной процедурой. В качестве приближенного значения интеграла будем использовать вариационную верхнюю оценку правдоподобия модели. Пусть задано непрерывное параметрическое распределение q , аппроксимирующие апостериорное распределение $p(\mathbf{W}, \Gamma | \mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c_{\text{temp}})$.

Тогда верно следующее выражение:

$$\log p(\mathbf{y} | \mathbf{X}, \mathbf{A}, \mathbf{m}, c) \geq \mathbb{E}_q \log p(\mathbf{y} | \mathbf{X}, \mathbf{W}, \Gamma, \mathbf{A}, \mathbf{m}, c_{\text{temp}}) - D_{KL}(q || p(\mathbf{w}, \Gamma | \mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c_{\text{temp}})). \quad (4)$$

Разница между верхней оценкой (4) и правдоподобием модели (2) определяется дивергенцией между вариационным распределением q и апостериорным распределением $p(\mathbf{W}, \Gamma | \mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c_{\text{temp}})$.

В дальнейшем будем использовать следующую форму вариационного распределения:

$$q = q_{\mathbf{W}} q_{\Gamma} : \\ q_{\mathbf{W}} \sim \mathcal{N}(\boldsymbol{\mu}_q, \mathbf{A}_q^{-1}), \quad q_{\Gamma} = \prod_{(j,k) \in E} q_{\gamma}^{j,k}, \quad q_{\gamma} \sim \mathcal{GS}(\mathbf{m}^{j,k}, c_q).$$

В дальнейшем будем обозначать за \mathbf{m}_q конкатенацию всех векторов средних $\mathbf{m}^{j,k}$.

Докажем теорему о дискретности задачи

Пусть задано значение концентрации c_{temp} .

3.1 О параметрической сложности

3.2 О структурной сложности

3.3 О переборе вариантов

3.4 Общая теорема

4 Вариационная постановка задачи

Теорема. При устремлении температуры c_{temp} задача оптимизации (2) становится дискретной. 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.

Определим основные величины, которые характеризуют сложность модели.

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

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

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

Определение Структурной сложностью C_γ 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.
7. Оптимизация должна позволять баланс регуляризации.

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_γ . Положим $\mathbf{h} = [\mathbf{A}, \mathbf{m}]$.

Пусть L — приближенное значение вариационной оценки правдоподобия:

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

где $\hat{\mathbf{w}} \sim q_w$, $\hat{\boldsymbol{\Gamma}} \sim q_\gamma$, β — 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 βm .

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

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

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

Формула нижней оценки получается подстановкой.

Пусть Q — валидационная функция:

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

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

Сформулируем задачу поиска оптимальной модели как двухуровневую задачу.

$$\hat{\mathbf{h}} = \arg \max_{\mathbf{h} \in \mathbb{R}^h} Q(T^n(\boldsymbol{\theta}_0, \mathbf{h})), \quad (5)$$

где T — оператор оптимизации, решающий задачу оптимизации:

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

Теорема. Пусть $D_{KL}(q_w || p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)) = 0$, $D_{KL}(q_\gamma || p(\boldsymbol{\Gamma} | \mathbf{y}, \mathbf{X}, \mathbf{A}, \mathbf{m}, c)) = 0$, пусть $c_1 = 1, c_2 = 1, c_3 = 0$. Тогда оптимизация (5) 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 $E_{q_w, q_\gamma} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c)$ 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. Заметим, что q и q' можно выразить следующим образом:

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

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

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

$$E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) - c_2 D_{KL}(q||p) - E_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) + c_2 D_{KL}(q'||p') \geq 0,$$

$$E_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) - c'_2 D_{KL}(q'||p') - E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) + c'_2 D_{KL}(q||p) \geq 0.$$

Вычитая неравенства получим:

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

$$E_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) \leq E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c).$$

С учетом полученных неравенств распишем доказываемое утверждение:

$$\begin{aligned} & \max_p -D_{KL}(q||p) - \max_{p'} -D_{KL}(q'||p') \propto \\ & \propto \max_p -c'_2 D_{KL}(q||p) + E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) - E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) - \\ & - \max_{p'} -c'_2 D_{KL}(q'||p') + E_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) + E_{q'} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \mathbf{\Gamma}, \mathbf{A}, \mathbf{m}, c) \leq 0, \end{aligned}$$

что и т.д.

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:

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

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

Доказательство По сумме выпуклых функций.

Теорема При $c \rightarrow 0$ $C_\gamma \rightarrow 0$.

Доказательство (идея) При $c \rightarrow 0$ $p(\Gamma)$ вероятностная масса концентрируется на вершинах симплекса. $D_{\text{KL}}(p||q)$ будет равняться $-\infty$ при q_γ , распределенном не на вершинах симплекса. Тогда q_Γ становится бинарным.

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

Утверждение (предварительно, нужно развить). Пусть $c_3 > 0, c \ll 0$ и все $p_k \in \mathbf{p}$ отражают распределения на вершинах симплекса. Тогда оптимизация приведет к q_γ , 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.

5 Вычислительный эксперимент

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

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

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 $c_1 = c_2 = 1c \sim 0$ (Evidence с низкой температурой) будет произведен выбор структуры f_4 .

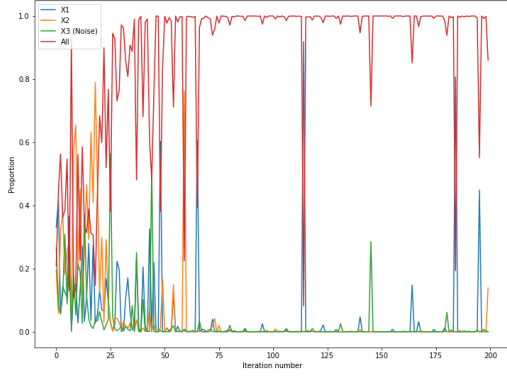


Figure 1: Evidence с низкой температурой

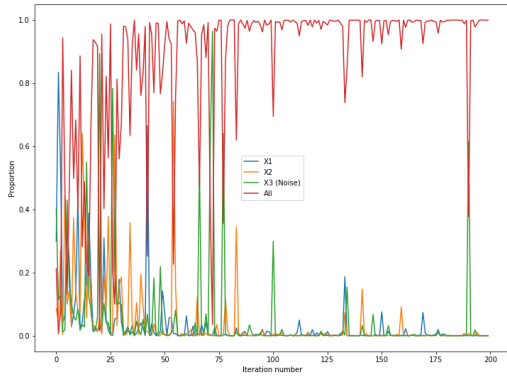


Figure 2: Evidence с высокой температурой

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 β выбор будет меняться в сторону f_3 .

Результаты

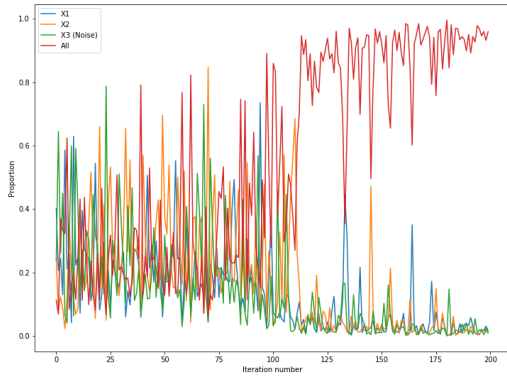


Figure 3: Evidence с высокой температурой, $\beta = 0.01$

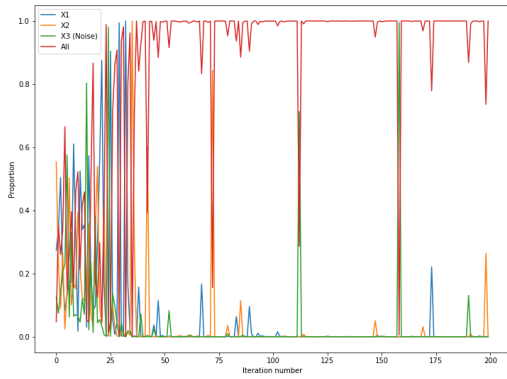


Figure 4: Поощрение выбора шумовой компоненты

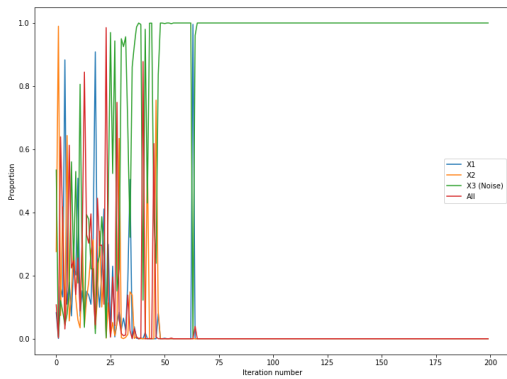


Figure 5: Поощрение выбора шумовой компоненты, $\beta = 0.01$