

Байесовская дистилляция моделей глубокого обучения¹

Исследуется проблема понижения сложности аппроксимирующих моделей. Рассматриваются методы основанные на дистилляции моделей глубокого обучения. Вводится понятия учителя и ученика. Предполагается, что модель ученика имеет меньшее число параметров, чем модель учителя. Предлагается байесовский подход к выбору модели ученика. Авторами предложен метод назначения априорного распределения параметров ученика на основе апостериорного распределения параметров модели учителя. Так как пространства параметров учителя и ученика не совпадают, предлагается механизм сопоставления этих пространств путем изменения структуры учителя. Проводится теоретический анализ предложенного механизма сопоставления. Вычислительный эксперимент проводится на синтетических и реальных данных. В качестве реальных данных рассматривается выборка FashionMNIST.

Ключевые слова: выбор модели; байесовский вывод; дистилляция модели; локальные преобразования; преобразования вероятностных пространств.

1. Введение

Исследуется проблема снижения числа обучаемых параметров моделей машинного обучения. Примерами таких моделей, с избыточным число параметров, являются: AlexNet [5], VGGNet [6], ResNet [12], BERT [8, 7], mT5 [10], GPT3[9] и т.д. Табл. 1

Таблица 1. Число параметров в моделях машинного обучения.

Название	AlexNet	VGGNet	ResNet	BERT	mT5	GPT3
Год	2012	2014	2015	2018	2020	2020
Тип данных	изображение	изображение	изображение	текст	текст	текст
Число параметров, млрд	0,06	0,13	0,06	0,34	13	175

приводит описание глубоких моделей машинного обучения. Видно, что число параметров моделей машинного обучения с годами растет. Это влечет снижение интерпретируемости моделей. Данная проблема рассматривается в специальном классе задач по состязательным атакам (англ. adversarial attack) [4]. Большое число

¹Работа выполнена при поддержке ... (грант № ...).

параметров требует больших вычислительных ресурсов. Из-за этого данные модели не могут быть использованы в мобильных устройствах. Для снижения числа параметров предложен метод дистилляции модели [14, 26, 27]. Дистиллируемая модель с большим числом параметров называется *учитель*, а модель получаемая путем дистилляции называется *ученик*. При оптимизации параметров модели ученика используется модель учителя с фиксированными параметрами.

Определение 1. Дистилляция модели — уменьшение сложности модели путем выбора модели в множестве более простых моделей на основе параметров и ответов более сложной фиксированной модели.

Идея дистилляции предложена в работах Дж. Е. Хинтона и В. Н. Вапником [14, 26, 27]. В этих работах предлагается использовать ответы учителя в качестве целевой переменной для обучения модели ученика. Поставлен ряд экспериментов, в которых проводилась дистилляция моделей для задачи классификации машинного обучения. Базовый эксперимент на выборке MNIST [15] показал применимость метода для дистилляции избыточно сложной нейросетевой модели в нейросетевую модель меньшей сложности. Эксперимент по распознаванию речи, в котором модель строилась путем дистилляции ансамбля моделей. Также в работе [14] был проведен эксперимент по обучению экспертных моделей на основе одной модели с большим числом параметров при помощи предложенного метода дистилляции на ответах учителя.

В работе [3] предложен метод передачи селективности нейронов (англ. neuron selectivity transfer) основанный на минимизации специальной функции потерь основанной на максимальном среднем отклонении (англ. maximum mean discrepancy) между выходами всех слоев модели учителя и ученика. Вычислительный эксперимент показал эффективность данного метода для задачи классификации изображений на примере выборок CIFAR [1] и ImageNet [2].

В данной работе предлагается метод основанный на байесовском выводе. В качестве априорного распределения параметров модели ученика предлагается использовать апостериорное распределение параметров модели учителя. Решается задача сопоставления пространства параметров модели учителя и модели ученика. Авторы предлагают подход, основанный на последовательном сопоставлении пространств параметров модели ученика и учителя. В результате сопоставления, параметры модели учителя и модели ученика лежат в одном пространстве. Как следствие в качестве априорного распределения параметров модели ученика выбирается апостериорное распределение параметров модели учителя.

В рамках вычислительного эксперимента проводится теоретический анализ. Предложенный метод дистилляции анализируется на примере синтетической выборки, а также реальной на выборке FashionMnist [19].

2. Постановка задачи дистилляции

Задана выборка

$$\mathfrak{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m, \quad \mathbf{x}_i \in \mathbb{R}^n, \quad y_i \in \mathbb{Y},$$

где \mathbf{x}_i, y_i — признаковое описание и целевая переменная i -го объекта, число объектов в обучающей выборке обозначается m . Размер признакового описания объектов

обозначим n . Множество $\mathbb{Y} = \{1, \dots, K\}$ для задачи классификации, где K число классов, множество $\mathbb{Y} = \mathbb{R}$ для задачи регрессии.

Задана модель учителя в виде суперпозиций линейных и нелинейных преобразований:

$$f = \sigma \circ \mathbf{U}_T \circ \sigma \circ \mathbf{U}_{T-1} \circ \dots \circ \sigma \circ \mathbf{U}_1,$$

где T — число слоев модели учителя, σ — функция активации, а \mathbf{U}_t обозначает матрицу линейного преобразования. Матрицы \mathbf{U} являются параметрами модели учителя f . Каждая матрица \mathbf{U}_t имеет размер $n_t \times n_{t-1}$, где $n_0 = n$, а $n_T = 1$ для задачи регрессии и $n_T = K$ для задачи классификации на K классов. Число параметров p_{tr} нейросетевой модели f вычисляется следующим образом:

$$p_{\text{tr}} = \sum_{t=1}^T n_t n_{t-1}.$$

Зададим полный порядок на множестве параметров модели учителя f . Вектор параметров модели учителя f обозначим \mathbf{u} . В данной работе для полносвязной нейронной сети рассматривается естественный порядок индуцированный, номером слоя t , номером нейрона и номером элемента вектора параметров нейрона (выбирается матрица \mathbf{U}_t , строка матрицы и элемент строки).

Например рассмотрим модель учителя для решения задачи регрессии:

$$(1) \quad f(\mathbf{x}) = \sigma \circ \mathbf{U}_3 \circ \sigma \circ \mathbf{U}_2 \circ \sigma \circ \mathbf{U}_1 \mathbf{x},$$

где для рассмотренной модели вектор параметров \mathbf{u} имеет следующий вид:

$$\mathbf{u} = [u_1^{1,1}, \dots, u_1^{1,n}, \dots, u_1^{n_1,1}, \dots, u_1^{n_1,n}, u_2^{1,1}, \dots, u_2^{1,n_1}, \dots, u_2^{n_2,1}, \dots, u_2^{n_2,n_1}, u_3^{1,1}, \dots, u_3^{1,n_2}].$$

Пусть для вектора параметров учителя f задано апостериорное распределение параметров $p(\mathbf{u}|\mathfrak{D})$.

На основе выборки \mathfrak{D} и учителя f требуется выбрать модель ученика из параметрического семейства функций:

$$g = \sigma \circ \mathbf{W}_L \circ \dots \circ \sigma \circ \mathbf{W}_1, \quad \mathbf{W}_l \in \mathbb{R}^{n_s \times n_{s-1}},$$

где L число слоев модели ученика. Аналогичным образом определяются n_0, n_L и p_{st} . Аналогичным образом введем вектор параметров модели ученика \mathbf{w} . Заметим, что каждая модель g из семейства задается своим вектором параметров \mathbf{w} . Следовательно задача выборки модели g эквивалента задаче оптимизации вектора параметров $\mathbf{w} \in \mathbb{R}^{p_{\text{st}}}$.

Оптимизация параметров $\hat{\mathbf{w}} \in \mathbb{R}^{p_{\text{st}}}$ проводится при помощи вариационного вывода на основе совместного правдоподобия модели и данных:

$$(2) \quad \mathcal{L}(\mathfrak{D}, \mathbf{A}) = \log p(\mathfrak{D}|\mathbf{A}) = \log \int_{\mathbf{w} \in \mathbb{R}^{p_g}} p(\mathfrak{D}|\mathbf{w}) p(\mathbf{w}|\mathbf{A}) d\mathbf{w},$$

где $p(\mathbf{w}|\mathbf{A})$ — априорное распределение вектора параметров модели ученика. Так как вычисление интеграла (2) является вычислительно сложной задачей, используем вариационный подход [24, 25]. Пусть задано вариационное распределение параметров модели ученика $q(\mathbf{w})$, которое аппроксимирует неизвестное апостериорное распределение $p(\mathbf{w}|\mathcal{D})$. Выбор параметров \mathbf{w} сводится к решению оптимизационной задачи:

$$(3) \quad \hat{\mathbf{w}} = \arg \min_{q, \mathbf{w}} D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w}|\mathbf{A})) - \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}),$$

заметим, что выражение (3) не учитывает параметры учителя f . Для использования информации о параметрах учителя предлагается рассмотреть параметры априорного распределения $p(\mathbf{w}|\mathbf{A})$ как функцию от апостериорного распределения $p(\mathbf{u}|\mathcal{D})$.

3. Построение априорного распределения ученика

Пусть апостериорное распределение параметров модели учителя является нормальным распределением:

$$(4) \quad p(\mathbf{u}|\mathcal{D}) = \mathcal{N}(\mathbf{u}_0, \Sigma_0),$$

где \mathbf{u}_0 и Σ_0 параметры апостериорного распределения.

3.1. Учитель и ученик принадлежат одному семейству

Рассмотрим следующие условия:

- 1) число слоев модели учителя равняется числу слоев модели ученика $L = T$;
- 2) размеры соответствующих слоев совпадают, другими словами, для всех t, l таких, что $t = l$ выполняется $n_l = n_t$, где n_t обозначает размер t -го слоя учителя, а n_l размер l -го слоя ученика.

В случае выполнения условий 1)–2) априорное распределение модели ученика равняется апостериорному распределению параметров учителя, то есть $p(\mathbf{w}|\mathbf{A}) = p(\mathbf{w}|\mathcal{D})$.

3.2. Удаление нейрона в слое учителя

Рассмотрим следующие условия:

- 1) число слоев модели учителя равняется числу слоев модели ученика $L = T$;
- 2) размеры соответствующих слоев не совпадают, другими словами, для всех t, l таких, что $t = l$ выполняется $n_l \leq n_t$, где n_t обозначает размер t -го слоя учителя, а n_l размер l -го слоя ученика.

Проведем согласование модели учителя и модели ученика при помощи последовательных преобразований параметров \mathbf{u} . Рассмотрим элементарное преобразование t -го слоя учителя:

$$\phi(t) : \mathbb{R}^{\text{ptr}} \rightarrow \mathbb{R}^{\text{ptr} - 2n_t}$$

вектора \mathbf{u} которое описывает удаление одного нейрона из t -го слоя. Другими словами, преобразование $\phi(t)$ зануляет одну из строк матрицы \mathbf{U}_t . Заметим, что данное зануление неизбежно производит зануление соответственного столбца матрицы \mathbf{U}_{t+1} .

Обозначим новый вектор параметров $\mathbf{u}' = \phi(t, \mathbf{u})$, а подмножество элементов, которые были удалены как \mathbf{u}'' . Аналогичным образом введем обозначения параметров нормального распределения $p(\mathbf{u}|\mathcal{D})$, а именно $\mathbf{u}'_0, \mathbf{u}''_0, \Sigma'_0, \Sigma''_0, \Sigma'''_0, \Sigma'''_0$.

Теорема 1. Пусть выполняются следующие условия:

- 1) апостериорное распределение параметров модели учителя является нормальным распределением (4);
- 2) число слоев модели учителя равняется числу слоев модели ученика $L = T$;
- 3) размеры соответствующих слоев не совпадают, другими словами, для всех t, l таких, что $t = l$ выполняется $n_s \leq n_l$.

Тогда апостериорное распределения параметров модели учителя после удаления одного нейрона имеет следующий вид:

$$(5) \quad p_{f'}(\mathbf{u}'|\mathfrak{D}) = \mathcal{N}(\mathbf{u}'_0 + \Sigma_0'^{\prime\prime} \Sigma_0''^{-1} (\mathbf{0} - \mathbf{u}_0''), \Sigma_0' - \Sigma_0'^{\prime\prime} \Sigma_0''^{-1} \Sigma_0''),$$

где f' обозначает модель f без одного нейрона.

Теорема 1 задает апостериорное распределение параметров (5) после зануления нейронов в модели нейросети. Заметим, что аналогичным образом можно удалить сразу подмножество нейронов в рамках одного слоя. В случае, если число нейронов отличается в нескольких слоях модели нейросети, то выполняется последовательно применения отображения $\phi(t, \mathbf{u})$ для каждого t -го слоя.

Заметим, что в рамках данного механизма не оговорен выбор нейронов для удаления в рамках одного слоя. Для этого предлагается выполнить упорядочивания нейронов в рамках каждого слоя. Где первый нейрон является наиболее значимым, а последний нейрон наименее значимым. Например порядок можно задать на основе отношения плотности распределения параметра к плотности распределения данного параметра в нуле [24] или на основе метода Белсли [25] и т.д. В рамках данной работы порядок на параметрах в рамках одного слоя задается случайным образом.

3.3. Удаление слоя учителя

Пусть в модели учителя требуется убрать t -й слой. Рассмотрим следующие условия:

- 1) соответствующие размеры слоев совпадают, $n_t = n_{t-1}$;
- 2) функция активации удовлетворяет следующему свойству $\sigma \circ \sigma = \sigma$.

Проведем согласование модели учителя и модели ученика при помощи последовательных преобразований вектора параметров \mathbf{u} . Рассмотрим элементарное преобразования:

$$\psi(t) : \mathbb{R}^{\text{Ptr}} \rightarrow \mathbb{R}^{\text{Ptr} - n_t n_{t-1}}$$

вектора \mathbf{u} которое описывает удаление одного t -го слоя. Другими словами, преобразования $\psi(t)$ превращает матрицу \mathbf{U}_t в единичную матрицу \mathbf{I} .

Теорема 2. Пусть выполняются следующие условия:

- 1) апостериорное распределение параметров модели учителя является нормальным распределением (4);
- 2) соответствующие размеры слоев совпадают, $n_t = n_{t-1}$;
- 3) функция активации удовлетворяет следующему свойству $\sigma \circ \sigma = \sigma$.

Тогда апостериорное распределения параметров модели учителя после удаления одного слоя имеет следующий вид:

$$(6) \quad p_{f'}(\mathbf{u}'|\mathfrak{D}) = \mathcal{N}(\mathbf{u}'_0 + \sigma_0'^{\prime\prime} \sigma_0''^{-1} (\mathbf{i} - \mathbf{u}_0''), \sigma_0' - \sigma_0'^{\prime\prime} \sigma_0''^{-1} \sigma_0''),$$

где f' обозначает модель f без одного слоя, $\mathbf{i} = [1, 0, \dots, 0, 0, 1, \dots, 0, 0, 1, \dots, 0, 0, \dots, 1]^T$.

Теорема 2 задает апостериорное распределение параметров (6) после удаления слоя нейросети. Полученное распределение $p(\mathbf{u}'|\mathfrak{D})$ является оценкой апостериорного распределения модели без одного слоя.

Заметим, что в рамках данного механизма не оговорен выбор слоя для удаления. Предполагается, что выбор слоя для удаления задается экспертно.

3.4. Выполнение последовательных преобразований

Локальные преобразования ϕ, ψ позволяют согласовать пространства параметров учителя f и ученика g . После сопоставления параметров в качестве априорного распределения параметров ученика $p(\mathbf{w}|\mathfrak{D})$ рассматривается полученное апостериорное распределение параметров учителя $p(\mathbf{u}'|\mathfrak{D})$. Получаем, следующее приближение априорного распределения:

$$p_g(\mathbf{w}|\mathfrak{D}) = p_{f'}(\mathbf{w}|\mathfrak{D}),$$

где данное априорное распределение используется для поиска оптимальных параметров модели ученика $\hat{\mathbf{w}}$ с помощью (3).

4. Вычислительный эксперимент

Проводится вычислительный эксперимент для анализа предложенного метода дистилляции на основе апостериорного распределения параметров модели учителя.

4.1. Синтетические данные

Проанализируем модель на синтетической выборке. Выборка построенная следующим образом:

$$\begin{aligned} \mathbf{w} &= [w_j : w_j \sim \mathcal{N}(0, 1)]_{n \times 1}, \quad \mathbf{X} = [x_{ij} : x_{ij} \sim \mathcal{N}(0, 1)]_{m \times n}, \\ \mathbf{y} &= [y_i : y_i \sim \mathcal{N}(\mathbf{x}_i^T \mathbf{w}, \beta)]_{m \times 1}, \end{aligned}$$

где $\beta = 0,1$ — уровень шума в данных. В эксперименте число признаков $n = 10$, для обучения и тестирования было сгенерировано $m_{\text{train}} = 900$ и $m_{\text{test}} = 124$ объекта.

В качестве модели учителя рассматривалась модель многослойный перцептрон с двумя скрытыми слоями (1). Матрицы линейных преобразований имеют размер:

$$\mathbf{U}_1 \in \mathbb{R}^{100 \times 10}, \mathbf{U}_2 \in \mathbb{R}^{50 \times 100}, \mathbf{U}_3 \in \mathbb{R}^{1 \times 50}.$$

В качестве функции активации была выбрана функция активации ReLu. Модель учителя предварительно обучена на основе вариационного вывода (3), где в качестве априорного распределения параметров выбрано стандартное нормальное распределение.

В качестве модели ученика были выбраны две конфигурации. Первая конфигурация получается путем удаления нейронов в модели учителя:

$$(7) \quad g = \sigma \circ \mathbf{W}_3 \circ \sigma \circ \mathbf{W}_2 \circ \sigma \circ \mathbf{W}_1,$$

где σ является нелинейной функцией активации, а матрицы линейных преобразований имеют размер:

$$\mathbf{W}_1 \in \mathbb{R}^{10 \times 10}, \mathbf{W}_2 \in \mathbb{R}^{10 \times 10}, \mathbf{W}_3 \in \mathbb{R}^{1 \times 10}.$$

В качестве функции активации была выбрана функция активации ReLu.

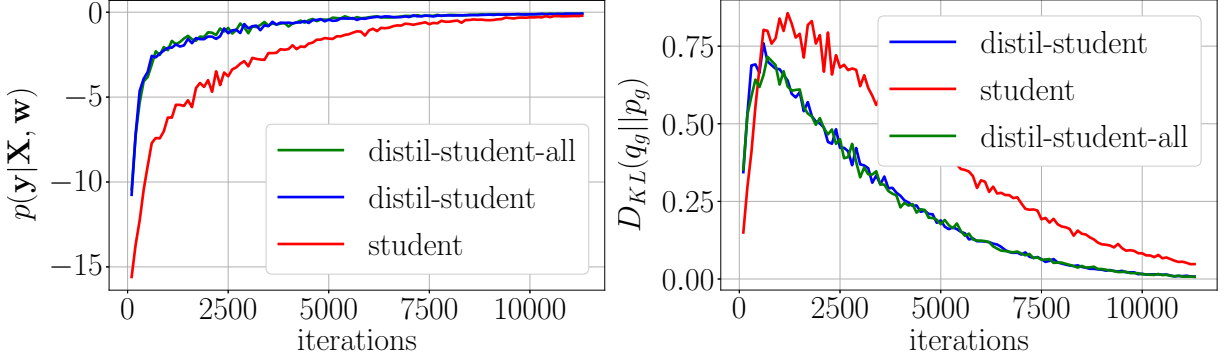


Рис. 1. Структура (7) модели ученика g . Слева: правдоподобие выборки в зависимости от номера итерации при обучении. Справа: KL-дивергенция между вариационным и априорным распределениями параметров модели.

На рис. 1 представлено сравнение моделей ученика основанных на структуре (7). Представлено сравнение разных моделей: модель без дистилляции, где в качестве априорного распределения выбирается стандартное нормальное распределение (на легенде обозначается student); модель с частичной дистилляцией, где в качестве среднего значения параметров выбираются параметры согласно выражения (5), а ковариационная матрица была приравнена к единичной матрицы (на легенде обозначается distil-student); модель с полной дистилляцией согласно выражения (5) (на легенде обозначается distil-student-all). Видно, что модели ученика, где в качестве априорного распределения выбраны распределения основанные на апостериорном распределение учителя имеют больше правдоподобие, чем модель где в качестве априорного распределения выбрано стандартное нормальное. Также заметим, что использования параметра среднего из апостериорного распределения дает основной вклад при дистилляции, так как качество моделей distil-student и distil-student-all совпадает.

Вторая конфигурация получается путем удаления слоя модели учителя:

$$(8) \quad g = \sigma \circ \mathbf{W}_2 \circ \sigma \circ \mathbf{W}_1,$$

где σ является нелинейной функцией активации, а матрицы линейных преобразований имеют размер:

$$\mathbf{W}_1 \in \mathbb{R}^{1 \times 50}, \mathbf{W}_2 \in \mathbb{R}^{50 \times 10}.$$

В качестве функции активации была выбрана функция активации ReLu.

На рис. 2 представлено сравнение моделей ученика основанных на структуре (8). Аналогично рис. 1 на рис. 2 представлено сравнение модели без дистилляции

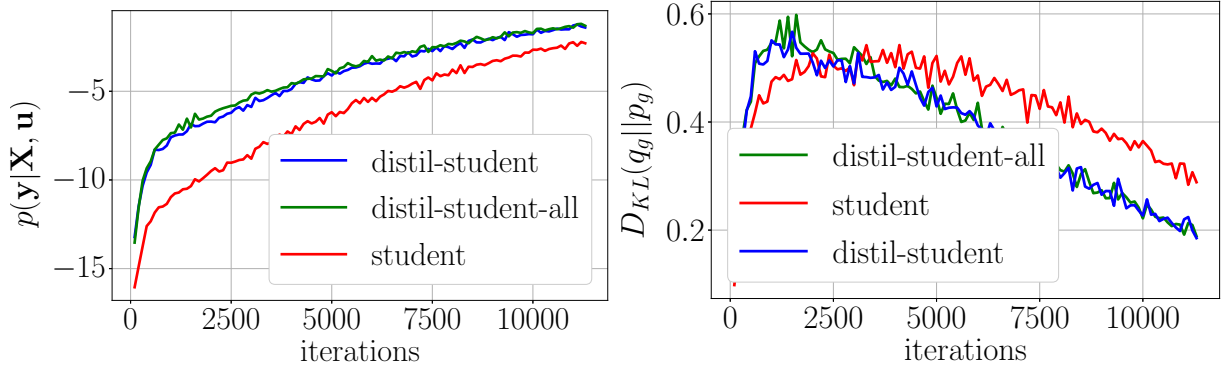


Рис. 2. Структура (8) модели ученика g . Слева: правдоподобие выборки в зависимости от номера итерации при обучении. Справа: KL-дивергенция между вариационным и априорным распределениями параметров модели.

(student), модели с дистилляцией параметра среднего значение (distil-student), модели с полной дистилляцией (distil-student-all). В рамках данного эксперимента, по дистилляции модели учителя в модель ученика с меньшим числом параметров получены результаты, которые подтверждают, что задание априорного распределения параметров ученика позволяет улучшить число итераций при выборе оптимальных параметров модели ученика.

4.2. Выборка FashionMnist

В рамках данного эксперимента проводился анализ байесовского подхода к дистилляции на реальных данных. В качестве реальных данных выбрана выборка FashionMnist[19] которая является задачей классификации изображений на 10 классов.

В качестве модели учителя рассматривалась модель многослойный перцептрон с двумя скрытыми слоями (1). Матрицы линейных преобразований имеют размер:

$$\mathbf{U}_1 \in \mathbb{R}^{800 \times 784}, \mathbf{U}_2 \in \mathbb{R}^{50 \times 800}, \mathbf{U}_3 \in \mathbb{R}^{10 \times 50},$$

В качестве функции активации была выбрана функция активации ReLu. Модель учителя предварительно обучена на основе вариационного вывода (3), где в качестве априорного распределения параметров выбрано стандартное нормальное распределение.

В качестве модели ученика были выбрана конфигурация с одним скрытым слоем (8), где матрицы линейных преобразований имеют размер:

$$\mathbf{W}_1 \in \mathbb{R}^{50 \times 784}, \mathbf{W}_2 \in \mathbb{R}^{50 \times 10}.$$

В качестве функции активации была выбрана функция активации ReLu.

На рис. 3 представлено сравнение моделей ученика с разными априорными распределениями на параметры. Аналогично синтетическому эксперименту сравнивался случай стандартного нормального распределения, использования априорного распределения с заданным средним значением параметров на основе апостериорного распределения, с определением полного априорного распределения на основе формулы (6). Видно, что у моделей с заданием априорного распределения на основе

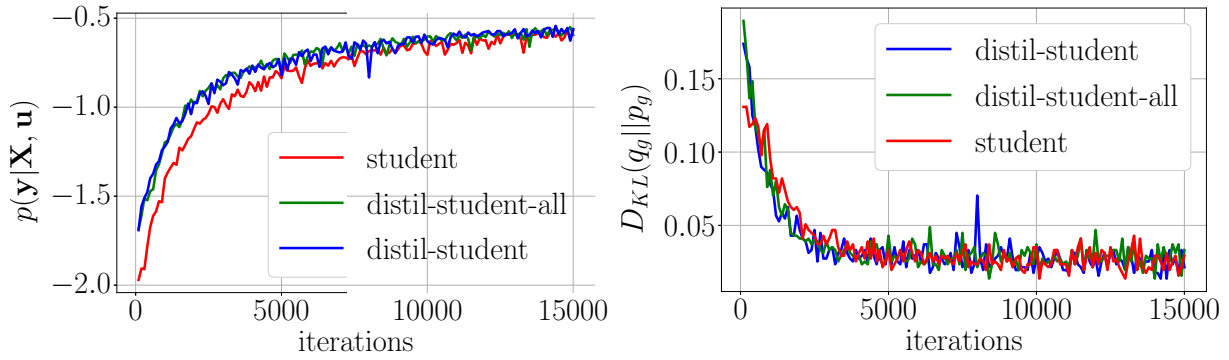


Рис. 3. Слева: правдоподобие выборки в зависимости от номера итерации при обучении. Справа: KL-дивергенция между вариационным и априорным распределениями параметров модели.

апостериорного распределения параметров учителя правдоподобие выборки выше, чем у модели, где в качестве априорного распределения выбрано стандартное нормальное распределение.

5. Заключение

В данной работе проанализирована байесовская дистилляция модели учителя в модель ученика на основе вариационного вывода. В рамках данной работы дистилляция основывается на задании априорного распределения параметров модели ученика. Априорное распределение параметров модели ученика задается на основе апостериорного распределения параметров модели учителя. Механизм преобразования структуры модели учителя в структуру модели ученика представлен в теореме 1 и теореме 2.

Теорема 1 описывает механизм сопоставления параметров модели учителя и ученика в случае, если число слоев совпадает, но размер слоев различается. Теорема 2 описывает механизм сопоставления параметров модели учителя и ученика в случае, если число слоев различается.

В вычислительном эксперименте сравнивается модель ученика, которая обучена без использования распределения параметров учителя и модель ученика, где в качестве априорного распределения параметров выбрано апостериорное распределение параметров модели учителя после сопоставления.

СПИСОК ЛИТЕРАТУРЫ

1. Alex Krizhevsky and Vinod Nair and Geoffrey Hinton CIFAR-10 (Canadian Institute for Advanced Research) // <http://www.cs.toronto.edu/~kriz/cifar.html>
2. Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L. Imagenet: A large-scale hierarchical image database // IEEE conference on computer vision and pattern recognition, 2009. P. 248–255.

3. *Huang, Zehao and Wang, Naiyan* Like What You Like: Knowledge Distill via Neuron Selectivity Transfer // arXiv e-prints, 2017.
4. *Kui Ren and Tianhang Zheng and Zhan Qin and Xue Liu* Adversarial Attacks and Defenses in Deep Learning // Engineering, 2020. P. 346–360.
5. *Alex Krizhevsky, Ilya Sutskever, Geoffrey Hinton* ImageNet Classification with Deep Convolutional Neural Networks // NIPS, 2012.
6. *Karen Simonyan and Andrew Zisserman* Very Deep Convolutional Networks for Large-Scale Image Recognition // NIPS, 2014.
7. *Vaswani A., Shazeer N., Parmar N., Uszkoreit J., Jones L., Gomez A., Kaiser L., Polosukhin I.* Attention Is All You Need // In Advances in Neural Information Processing Systems. 2017. V. 5. P. 6000–6010.
8. *Devlin J., Chang M., Lee K., Toutanova K.* BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding // arXiv preprinted, 2018.
9. *Tom B. Brown et al* GPT3: Language Models are Few-Shot Learners // arXiv preprinted, 2020.
10. *Linting Xue and Noah Constant and Adam Roberts and Mihir Kale and Rami Al-Rfou and Aditya Siddhant and Aditya Barua and Colin Raffel.* mT5: A massively multilingual pre-trained text-to-text transformer // arXiv preprinted, 2021.
11. *Yang, Ziqing and Cui, Yiming and Chen, Zhipeng and Che, Wanxiang and Liu, Ting and Wang, Shijin and Hu, Guoping* TextBrewer: An Open-Source Knowledge Distillation Toolkit for Natural Language Processing // Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations. 2020. P. 9–16.
12. *He K., Zhang X., Ren S., Sun J.* Deep Residual Learning for Image Recognition // Proc. of the IEEE Conference on Computer Vision and Pattern Recognition. Las Vegas, 2016. P. 770–778.
13. *Бахтеев О. Ю., Стрижов В. В.* Выбор моделей глубокого обучения субоптимальной сложности // АИТ. 2018. № 8. С. 129–147.
14. *Hinton G., Vinyals O., Dean J.* Distilling the Knowledge in a Neural Network // NIPS Deep Learning and Representation Learning Workshop. 2015.
15. *LeCun Y., Cortes C., Burges C.* The MNIST dataset of handwritten digits, 1998. <http://yann.lecun.com/exdb/mnist/index.html>.
16. *Vapnik V., Izmailov R.* Learning Using Privileged Information: Similarity Control and Knowledge Transfer // Journal of Machine Learning Research. 2015. No 16. P. 2023–2049.
17. *Lopez-Paz D., Bottou L., Scholkopf B., Vapnik V.* Unifying Distillation and Privileged Information // In International Conference on Learning Representations. Puerto Rico, 2016.

18. *Madala H., Ivakhnenko A.* Inductive Learning Algorithms for Complex Systems Modeling. Boca Raton: CRC Press Inc., 1994.
19. *Xiao H., Rasul K., Vollgraf R.* Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms // arXiv preprint arXiv:1708.07747. 2017.
20. *Wilson T., Kozareva Z., Nakov P., Rosenthal S., Stoyanov V., Ritter A.* SemEval-2013 Task 2: Sentiment Analysis in Twitter // Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013). Atlanta, 2013. P. 312–320.
21. *LeCun Y., Boser B., Denker J., Henderson D., Howard R., Hubbard W., Jackel L.* Backpropagation Applied to Handwritten Zip Code Recognition // Neural Computation. 1989. V. 1. No 4. P. 541–551.
22. *Hochreiter S., Schmidhuber J.* Long short-term memory // Neural Computation. 1997. V. 9. No 8. P. 1735–1780.
23. *Kingma D, Ba J.* Adam: A Method for Stochastic Optimization // arXiv preprint arXiv:1412.6980. 2014.
24. *Graves A.* Practical Variational Inference for Neural Networks // Advances in Neural Information Processing Systems, 2011. Vol. 24. P. 2348–2356.
25. *Grabovoy A.V., Bakhteev O.Y., Strijov V.V.* Estimation of relevance for neural network parameters // Informatics and Applications, 2019. Vol.13 No 2. P. 62–70.
26. *Vapnik V., Izmailov R.* Learning Using Privileged Information: Similarity Control and Knowledge Transfer // Journal of Machine Learning Research. 2015. No 16. P. 2023–2049.
27. *Lopez-Paz D., Bottou L., Scholkopf B., Vapnik V.* Unifying Distillation and Privileged Information // In International Conference on Learning Representations. Puerto Rico, 2016.