

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 другие. 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

описывает глубокие модели машинного обучения. 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 параметров требует больших вычислительных ресурсов. Из-за этого данные модели не могут быть использованы в

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

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 — снижение 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.

Определение 2. Структура модели — множество структурных параметров модели, которые задают вид суперпозиции.

Определение 3. Сопоставление параметрических моделей — изменение структуры модели (одной или нескольких моделей) в результате которого векторы параметров различных моделей лежат в одном пространстве.

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 данной работе в качестве параметрических моделей рассматривается полносвязная нейронная сеть. В качестве структурных параметров модели выбраны число слоев, а также размер каждого скрытого слоя.

В рамках предложенного метода сопоставления параметрических моделей, не оговорен выбор порядка на множестве параметров модели учителя. Для этого предлагается упорядочивать параметры модели учителя на основе их значимости. Первый нейрон является наиболее значимым, а последний нейрон наименее значимым. Поряд-

док задается 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.

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 обозначается 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:

$$f = \sigma \circ \mathbf{U}_T \sigma \circ \mathbf{U}_{T-1} \circ \dots \circ \mathbf{U}_2 \sigma \circ \mathbf{U}_1,$$

где T — число слоев модели учителя, σ — функция активации, а \mathbf{U}_t обозначает матрицу линейного преобразования. Матрицы \mathbf{U} соединяются в вектор параметров \mathbf{u} модели учителя f :

$$(1) \quad \mathbf{u} = \text{vec}([\mathbf{U}_T, \mathbf{U}_{T-1}, \dots, \mathbf{U}_1]),$$

где vec операция векторизация последовательности матриц в вектор параметров. 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 ν учителя f :

$$(2) \quad \nu = \sum_{t=1}^T n_t n_{t-1}.$$

Для построения вектора параметров \mathbf{u} задается полный порядок на элементов матриц \mathbf{U}_t . Для полносвязной нейронной сети вводится естественный порядок, индуцированный номером слоя t , номером нейрона, и номером элемента вектора параметров нейрона: выбирается матрица \mathbf{U}_t , строка этой матрицы и элемент строки.

Например, для модели учителя в задаче регрессии:

$$(3) \quad f(\mathbf{x}) = \sigma \circ \mathbf{U}_3 \sigma \circ \mathbf{U}_2 \sigma \circ \mathbf{U}_1 \mathbf{x},$$

вектор параметров \mathbf{u} принимает 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 известно апостериорное распределение параметров $p(\mathbf{u}|\mathcal{D})$. На основе выборки \mathcal{D} и апостериорного распределения параметров учителя f требуется выбрать модель ученика из параметрического семейства функций:

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

где L число слоев модели ученика. Число параметров v модели ученика g вычисляется аналогично выражению (2). Вектор параметров модели ученика \mathbf{w} строится аналогичным образом (1). Модель g задается своим вектором параметров \mathbf{w} . Следовательно, задача выбора модели g эквивалентна задаче оптимизации вектора параметров $\mathbf{w} \in \mathbb{R}^v$.

Параметры $\hat{\mathbf{w}} \in \mathbb{R}^v$ оптимизируются при помощи вариационного вывода на основе совместного правдоподобия модели и данных:

$$(4) \quad \mathcal{L}(\mathcal{D}, \mathbf{A}) = \log p(\mathcal{D}|\mathbf{A}) = \log \int_{\mathbf{w} \in \mathbb{R}^v} p(\mathcal{D}|\mathbf{w}) p(\mathbf{w}|\mathbf{A}) d\mathbf{w},$$

где $p(\mathbf{w}|\mathbf{A})$ — априорное распределение вектора параметров модели ученика. Так как взятие интеграла (4) является вычислительно сложной задачей, используем вариационный подход [24, 25]. Для этого зададим вариационное распределение параметров модели ученика $q(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, которое аппроксимирует неизвестное апостериорное распределение $p(\mathbf{w}|\mathcal{D})$:

$$q(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \approx p(\mathbf{w}|\mathcal{D}).$$

Далее распределение $q(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ будем обозначать просто $q(\mathbf{w})$. Оптимизация параметров \mathbf{w} сводится к решению задачи:

$$(5) \quad \hat{\mathbf{w}} = \arg \min_{\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{w}} D_{\text{KL}}(q(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) || p(\mathbf{w}|\mathbf{A})) - \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}).$$

Выражение (5) не учитывает параметры учителя f . Для использования информации о распределении параметров учителя предлагается рассмотреть параметры априорного распределения $p(\mathbf{w}|\mathbf{A})$ как функцию от апостериорного распределения учителя $p(\mathbf{u}|\mathcal{D})$.

3. Построение априорного распределения ученика

Зададим апостериорное распределение параметров модели учителя как нормальное распределение:

$$(6) \quad p(\mathbf{u}|\mathcal{D}) = \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma}),$$

где \mathbf{m} и $\boldsymbol{\Sigma}$ параметры этого распределения. На основе параметров \mathbf{m} и $\boldsymbol{\Sigma}$ требуется задать параметры \mathbf{A} априорного распределения $p(\mathbf{w}|\mathbf{A})$. В случае, когда структура моделей учителя и ученика задаются числом слоев и размером этих слоев, то возможны следующие варианты: число слоев и размер каждого слоя совпадает; число слоев совпадает, а размеры различаются; не совпадает число слоев.

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 этих условий, апериорное распределение параметров модели ученика приравнивается к апостериорному распределению параметров учителя, то есть $p(\mathbf{w}|\mathbf{A}) = p(\mathbf{u}|\mathcal{D})$.

3.2. Удаление нейрона в слое учителя

Проведем согласование модели учителя и модели ученика, согласно определению 3 при помощи последовательных преобразований параметров \mathbf{u} . Рассмотрим преобразование:

$$\phi(t, \mathbf{u}) : \mathbb{R}^\nu \rightarrow \mathbb{R}^{\nu-2n_t}$$

вектора \mathbf{u} , которое описывает удаление одного нейрона из t -го слоя учителя. Обозначим новый вектор параметров $\mathbf{v} = \phi(t, \mathbf{u})$, а элементы вектора, которые были удалены как $\bar{\mathbf{v}}$. Заметим, что векторы \mathbf{v} и $\bar{\mathbf{v}}$ являются случайными величинами.

Теорема 1. Пусть выполняются следующие условия:

- 1) апостериорное распределение $p(\mathbf{u}|\mathcal{D})$ параметров модели учителя является нормальным распределением (6);
- 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 $n_t \leq n_l$.

Тогда апостериорное распределение параметров модели учителя $p(\mathbf{v}|\mathcal{D})$ также является нормальным распределением.

Доказательство. Не уменьшая общности, пусть $\phi(t, \mathbf{u})$ удаляет j -й нейрон в t -м слое, что является удалением j -й строки матрицы \mathbf{U}_t . Заметим, что удаление j -й строки матрицы \mathbf{U}_t влечет удаление j -й компоненты вектора \mathbf{z}_{t+1} , где:

$$\mathbf{z}_t = \boldsymbol{\sigma} \circ \mathbf{U}_{t-1} \boldsymbol{\sigma} \circ \dots \circ \mathbf{U}_2 \boldsymbol{\sigma} \circ \mathbf{U}_1 \mathbf{x}.$$

Удаление j -й компоненты вектора \mathbf{z}_{t+1} эквивалентно занулению j -го столбца матрицы \mathbf{U}_{t+1} . Заметим, что тогда предсказание модели не зависит от параметров j -й строки матрицы \mathbf{U}_t , а следовательно донными параметрами также можно пренебречь.

Найдем распределение вектора \mathbf{v} . Для поиска распределения вектора параметров после зануления j -го столбца матрицы \mathbf{U}_{t+1} воспользуемся формулой условной вероятности, а для удаления j -й строки матрицы \mathbf{U}_t воспользуемся маргинализацией распределения. Обозначим зануляемые параметры модели как $\boldsymbol{\nu}_1$, а удаляемые параметры как $\boldsymbol{\nu}_2$. Также обозначим все параметры, которые не были занулены как $\bar{\boldsymbol{\nu}}_1 = [\bar{\boldsymbol{\nu}}_1^T, \boldsymbol{\nu}_2^T]^T$. Итоговое распределение параметров принимает следующий вид:

$$p(\mathbf{v}|\mathcal{D}) = \int_{\boldsymbol{\nu}_2} p(\bar{\boldsymbol{\nu}}_1|\mathcal{D}, \boldsymbol{\nu}_1 = \mathbf{0}) d\boldsymbol{\nu}_2.$$

Из свойств нормального распределения следует, что распределение $p(\bar{\nu}_1|\mathcal{D}, \nu_1 = \mathbf{0})$ также является нормальным распределением с параметрами μ, Ξ :

$$\begin{aligned}\mu &= \mathbf{m}_{\bar{\nu}_1} + \Sigma_{\bar{\nu}_1, \nu_1} \Sigma_{\nu_1, \nu_1}^{-1} (\mathbf{0} - \mathbf{m}_{\nu_1}), \\ \Xi &= \Sigma_{\bar{\nu}_1, \bar{\nu}_1} - \Sigma_{\bar{\nu}_1, \nu_1} \Sigma_{\nu_1, \nu_1}^{-1} \Sigma_{\nu_1, \bar{\nu}_1},\end{aligned}$$

где введены обозначения $\mathbf{m}_{\bar{\nu}_1}, \mathbf{m}_{\nu_1}$ соответствует подвектору вектора \mathbf{m} , который относится к параметрам $\bar{\nu}_1$ и ν_1 соответственно. Ковариационная матрица $\Sigma_{\bar{\nu}_1, \nu_1}$ обозначает подматрицу матрицы Σ , которая соответствует ковариационной матрицей между параметрами $\bar{\nu}_1$ и ν_1 .

Распределение $p(\mathbf{v}|\mathcal{D})$ найдем при помощи маргинализации распределения $p(\bar{\nu}_1|\mathcal{D}, \nu_1 = \mathbf{0})$ по параметрам ν_2 . Используя свойства нормального распределения получаем распределения:

$$(7) \quad p(\mathbf{v}|\mathcal{D}) = \mathcal{N}(\mu_v, \Xi_{v,v}),$$

где μ_v обозначает подвектор вектора μ , который относится к параметру \mathbf{v} , а матрица $\Xi_{v,v}$ является подматрицей матрицы Ξ , которая относится к вектору параметров \mathbf{v} .

Теорема 1 задает апостериорное распределение параметров (7) после зануления нейронов в модели нейросети — учителя. Заметим, что аналогичным образом можно удалить сразу подмножество нейронов в рамках одного слоя. В случае, если число нейронов отличается в нескольких слоях модели нейросети учителя, то выполняется последовательно применения отображения $\phi(t, \mathbf{u})$ для каждого t -го слоя.

3.3. Удаление слоя учителя

Проведем согласование модели учителя и модели ученика, согласно определения 3 при помощи последовательных преобразований вектора параметров \mathbf{u} . Рассмотрим преобразование:

$$\psi(t, \mathbf{u}) : \mathbb{R}^\nu \rightarrow \mathbb{R}^{\nu - n_t n_{t-1}}$$

вектора \mathbf{u} которое описывает удаление одного t -го слоя. Обозначим новый вектор параметров $\mathbf{v} = \psi(t, \mathbf{u})$, а элементы вектора, которые были удалены как $\bar{\mathbf{v}}$.

Теорема 2. Пусть выполняются следующие условия:

1) апостериорное распределение параметров $p(\mathbf{u}|\mathcal{D})$ модели учителя является нормальным распределением (6);

2) соответствующие размеры слоев совпадают, $n_t = n_{t-1}$, то есть матрица \mathbf{U}_t является квадратной;

3) функция активации удовлетворяет свойству идемпотентности $\sigma \circ \sigma = \sigma$. Тогда апостериорное распределения также описывается нормальным распределением со следующей плотностью распределения:

$$(8) \quad p(\mathbf{v}|\mathcal{D}) = \mathcal{N}(\mathbf{m}_v + \Sigma_{v, \bar{\mathbf{v}}} \Sigma_{\bar{\mathbf{v}}, \bar{\mathbf{v}}}^{-1} (\mathbf{i} - \bar{\mathbf{v}}), \Sigma_{v,v} - \Sigma_{v, \bar{\mathbf{v}}} \Sigma_{\bar{\mathbf{v}}, \bar{\mathbf{v}}}^{-1} \Sigma_{\bar{\mathbf{v}}, v}),$$

где вектор \mathbf{i} задается следующим образом:

$$\mathbf{i} = [\underbrace{1, 0, \dots, 0}_{n_t}, \underbrace{0, 0, 1, \dots, 0}_{n_t}, \underbrace{0, 0, 0, 1, \dots, 0}_{n_t}, \underbrace{0, 0, \dots, 1}_{n_t}]^\top.$$

Доказательство. Рассмотрим структуру нейронной сети с T слоями и $T + 1$ слоем. Не уменьшая общности для удаления рассматривается t -й слой, для которого выполняются условия теоремы. Заметим, что если заметить t -й слой нейронной сети с $T + 1$ слоем, то он будет эквивалентным архитектуре с T слоями:

$$\begin{aligned} f &= \sigma \circ \mathbf{U}_{T+1} \sigma \circ \mathbf{U}_T \cdots \sigma \circ \mathbf{U}_t \sigma \circ \cdots \mathbf{U}_2 \sigma \circ \mathbf{U}_1 = \\ &= \sigma \circ \mathbf{U}_{T+1} \sigma \circ \mathbf{U}_T \cdots \sigma \circ \mathbf{I} \sigma \circ \cdots \mathbf{U}_2 \sigma \circ \mathbf{U}_1 = \\ &= \sigma \circ \mathbf{U}_{T+1} \sigma \circ \mathbf{U}_T \cdots \sigma \circ \sigma \circ \cdots \mathbf{U}_2 \sigma \circ \mathbf{U}_1 =^1 \\ &=^1 \sigma \circ \mathbf{U}_{T+1} \sigma \circ \mathbf{U}_T \cdots \sigma \circ \cdots \mathbf{U}_2 \sigma \circ \mathbf{U}_1 \end{aligned}$$

Получаем, что удаление t -го слоя нейросети эквивалентно приравнивая t -го слоя к единичной матрице. Распределение параметров после приравнивая к единичной матрице вычисляется при помощи условного распределения. В силу общих свойств нормального распределения условное распределения также является нормальным распределением с параметрами μ, Ξ :

$$\begin{aligned} \mu &= \mathbf{m}_v + \Sigma_{v, \bar{v}} \Sigma_{\bar{v}, \bar{v}}^{-1} (\mathbf{i} - \bar{v}), \\ \Xi &= \Sigma_{v, v} - \Sigma_{v, \bar{v}} \Sigma_{\bar{v}, \bar{v}}^{-1} \Sigma_{\bar{v}, v}, \end{aligned}$$

где вектор \mathbf{m}_v является подвектором вектора \mathbf{m} соответствующий параметрам v , а матрица $\Sigma_{v, \bar{v}}$ является подматрицей ковариационной матрицы Σ соответствующий векторам параметров v и \bar{v} .

Теорема 2 задает апостериорное распределение параметров (8) после удаления слоя нейросети. Полученное распределение $p(v|\mathfrak{D})$ является оценкой апостериорного распределения модели без одного слоя.

3.4. Выполнение последовательных преобразований

Преобразования ϕ, ψ согласовывают пространства параметров учителя f и ученика g . После сопоставления параметрических моделей получаем, что параметры модели учителя и модели ученика принадлежат одному семейству 3.1.

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 (3). Матрицы линейных преобразований имеют размер:

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

В качестве функции активации была выбрана функция активации ReLu. Модель учителя предварительно обучена на основе вариационного вывода (5), 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:

$$(9) \quad g = \sigma \circ \mathbf{W}_3 \sigma \circ \mathbf{W}_2 \sigma \circ \mathbf{W}_1,$$

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

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

В качестве функции активации была выбрана функция активации ReLu.

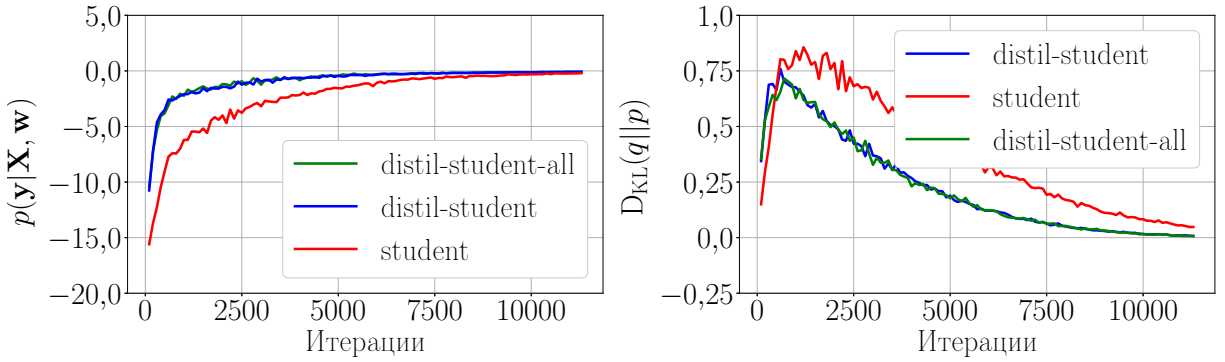


Рис. 1. Структура (9) модели ученика g . Слева: правдоподобие выборки в зависимости от номера итерации при обучении. Справа: KL-дивергенция между вариационным и априорным распределениями параметров модели.

Рис. 1 сравнивает модели ученика, со структурой (9). 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 (7), а ковариационная матрица была приравнена к единичной матрицы (на легенде обозначается distil-student); модель с полной дистилляцией согласно выражения (7) (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:

$$(10) \quad g = \sigma \circ \mathbf{W}_2 \sigma \circ \mathbf{W}_1,$$

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

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

В качестве функции активации была выбрана функция активации ReLu.

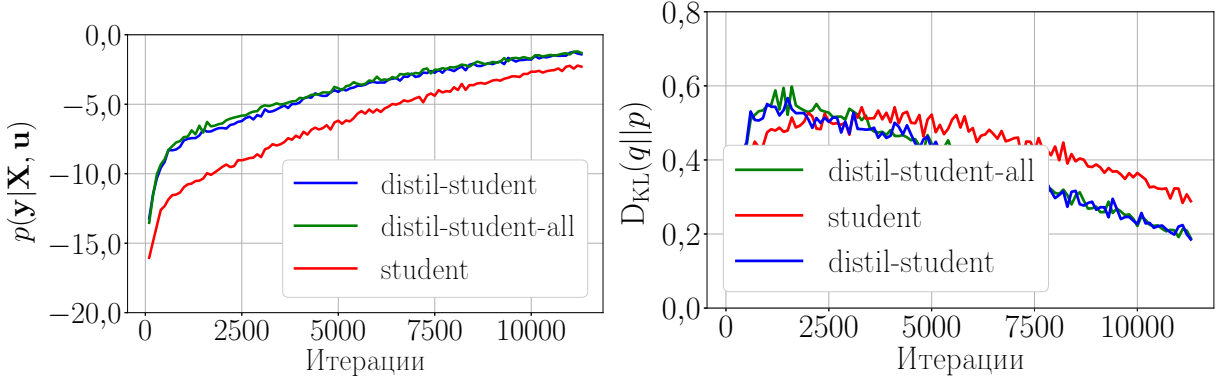


Рис. 2. Структура (10) модели ученика g . Слева: правдоподобие выборки в зависимости от номера итерации при обучении. Справа: KL-дивергенция между вариационным и априорным распределениями параметров модели.

Рис. 2 сравнивает модели ученика со структурой (10). Аналогично рис. 1 на рис. 2 представлено сравнение модели без дистилляции (student), модели с дистилляцией параметра среднего значения (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 (3). Матрицы линейных преобразований имеют размер:

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

В качестве функции активации была выбрана функция активации ReLu. Модель учителя предварительно обучена на основе вариационного вывода (5), где в качестве априорного распределения параметров выбрано стандартное нормальное распределение.

Таблица 2. Сводная таблица результатов вычислительного эксперимента.

	teacher	student	distil-student	distil-student-all
Эксперимент на синтетической выборке (удаление нейрона)				
Структура	[10, 100, 50, 1]	[10, 10, 10, 1]	[10, 10, 10, 1]	[10, 10, 10, 1]
Число параметров	6050	210	210	210
Разность площадей	-	0	16559	16864
Эксперимент на синтетической выборке (удаление слоя)				
Структура	[10, 100, 50, 1]	[10, 50, 1]	[10, 50, 1]	[10, 50, 1]
Число параметров	6050	550	550	550
Разность площадей S	-	0	23310	25506
Эксперимент на выборке FashionMnist				
Структура	[784, 800, 50, 10]	[784, 50, 10]	[784, 50, 10]	[784, 50, 10]
Число параметров	667700	39700	39700	39700
Разность площадей S	-	0	1165	1145

В качестве модели ученика были выбрана конфигурация с одним скрытым слоем (10), где матрицы линейных преобразований имеют размер:

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

В качестве функции активации была выбрана функция активации ReLu.

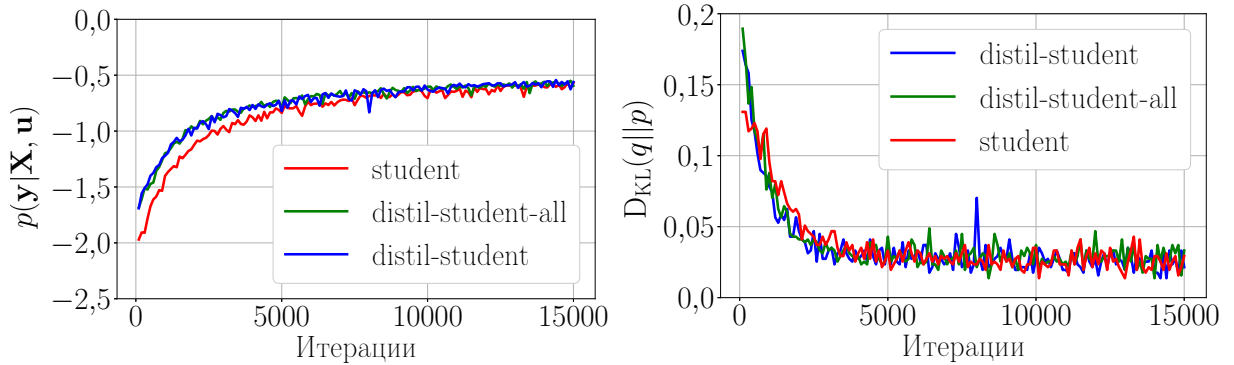


Рис. 3. Слева: правдоподобие выборки в зависимости от номера итерации при обучении. Справа: KL-дивергенция между вариационным и априорным распределениями параметров модели.

Рис. 3 сравнивает модели 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 (8). Видно, что у моделей с заданием априорного распределения на основе апостериорного распределения параметров учителя правдоподобие выборки выше, чем у модели, где в качестве априорного распределения выбрано стандартное нормальное распределение.

В табл. 2 представлен результат вычислительного эксперимента. Для численного сравнения качества моделей выбрана разность площадей графика $p(\mathbf{y}|\mathbf{X}, \mathbf{u})$ между моделью student и моделями distil-student и distil-student-all соответственно:

$$S = \sum_k^K p(\mathbf{y}|\mathbf{X}, \mathbf{u}^s) - p(\mathbf{y}|\mathbf{X}, \mathbf{u}^{ds}),$$

где $\mathbf{u}^s, \mathbf{u}^{ds}$ обозначает параметры модели студента и модели дистиллированного студента соответственно. Заметим, что площадь S имеет знак: чем большее положительное число, тем дистиллированная модель лучше, чем модель построенная без учителя. В случае, если площадь S принимает отрицательное значение, то значит модель без дистилляции является лучше чем модель с дистилляцией. В рамках вычислительного эксперимента видно, что площадь S под графиками принимает положительные значения, то есть модели ученика полученные при помощи дистилляции являются лучше чем модель ученика без дистилляции.

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. Краткое описание эксперимента представлено в табл. 2.

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.