

Дистилляция моделей глубокого обучения

Грабовой Андрей Валериевич

Московский физико-технический институт

МФТИ, г. Долгопрудный

Вероятностная интерпретация дистилляции моделей

Цель

Предложить вероятностную постановку задачи дистилляции моделей глубокого обучения. Развить существующие методы дистилляции и привилегированного обучения используя байесовский подход.

Задачи

1. Поставить вероятностную задачу дистилляции для задач классификации и регрессии.
2. Предложить метод байесовской дистилляции нейросетевых моделей.
3. Провести теоретический анализ полученных результатов.

Исследуемая проблема

Снижение размерности пространства параметров моделей глубокого обучения.

Список литературы

1. Грабовой А.В., Стрижов В.В. Анализ моделей привилегированного обучения и дистилляции // Автоматика и телемеханика, 2021 (на рецензировании).
2. Грабовой А.В., Бахтеев О.Ю., Стрижов В.В. Введение отношения порядка на множестве параметров аппроксимирующих моделей // Информатика и ее применения, 2020, 14(2) : 58-65.
3. Грабовой А.В., Стрижов В.В. Байесовская дистилляция моделей глубокого обучения // (текущая работа, в процессе).
4. Lopez-Paz D., Bottou L., Scholkopf B., Vapnik V. Unifying Distillation and Privileged Information // In International Conference on Learning Representations. Puerto Rico, 2016.
5. Hinton G., Vinyals O., Dean J. Distilling the Knowledge in a Neural Network // NIPS Deep Learning and Representation Learning Workshop. 2015.
6. Madala H., Ivakhnenko A. Inductive Learning Algorithms for Complex Systems Modeling. Boca Raton: CRC Press Inc., 1994.

Классическая постановка задачи обучения с учителем

Заданы:

- 1) признаки $\mathbf{x}_i \in \mathbb{R}^n$;
- 2) привилегированные признаки $\mathbf{x}_i^* \in \mathbb{R}^{n^*}$;
- 3) целевая переменная $y_i \in \mathbb{Y}$;
- 4) индексы объектов, для которых известна привилегированная информация \mathcal{I} , а для которых она не известна $\bar{\mathcal{I}}$.

Функции учителя $\mathbf{f} : \mathbb{R}^{n^*} \rightarrow \mathbb{Y}'$ и ученика $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{Y}'$ — пространство оценок. Ответ $\mathbf{s}_i = \mathbf{f}(\mathbf{x}_i^*)$ функции \mathbf{f} для объекта \mathbf{x}_i^* используется при решении оптимизационной задачи.

Требуется выбрать модель ученика \mathbf{g} из множества

$$\mathfrak{G} = \{\mathbf{g} | \mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{Y}'\}.$$

Оптимизационная задача:

$$\mathbf{g} = \arg \min_{\mathbf{g} \in \mathfrak{G}} \mathcal{L}(\mathbf{g}, \mathbf{f}, \mathbf{X}, \mathbf{X}^*, \mathbf{y}),$$

где \mathcal{L} функция ошибки.

Постановка задачи: дистилляция Хинтона

Заданы:

- 1) $\mathbf{x}_i^* = \mathbf{x}_i$ для всех $i \in \{1, 2, \dots, m\}$;
- 2) $y_i \in \mathbb{Y} = \{1, \dots, K\}$, $\mathbb{Y}' = \mathbb{R}^K$.

Параметрические семейства учителя и ученика:

$$\begin{aligned}\mathfrak{F}_{\text{cl}} &= \{\mathbf{f} | \mathbf{f} = \text{softmax}(\mathbf{v}(\mathbf{x})/T), \quad \mathbf{v} : \mathbb{R}^n \rightarrow \mathbb{R}^K\}, \\ \mathfrak{G}_{\text{cl}} &= \{\mathbf{g} | \mathbf{g} = \text{softmax}(\mathbf{z}(\mathbf{x})/T), \quad \mathbf{z} : \mathbb{R}^n \rightarrow \mathbb{R}^K\},\end{aligned}$$

где \mathbf{z}, \mathbf{v} — это дифференцируемые параметрические функции заданной структуры, T — параметр температуры.

Функция ошибки

$$\mathcal{L}_{\text{st}}(\mathbf{g}) = - \underbrace{\sum_{i=1}^m \sum_{k=1}^K y_i^k \log \mathbf{g}(\mathbf{x}_i)}_{\text{исходная функция потерь}} \Big|_{T=1} - \underbrace{\sum_{i=1}^m \sum_{k=1}^K \mathbf{f}(\mathbf{x}_i)}_{\text{слагаемое дистилляция}} \Big|_{T=T_0} \log \mathbf{g}(\mathbf{x}_i) \Big|_{T=T_0},$$

где $\cdot|_{T=t}$ параметр температуры T равняется t .

Оптимизационная задача:

$$\hat{\mathbf{g}} = \arg \min_{\mathbf{g} \in \mathfrak{G}_{\text{cl}}} \mathcal{L}_{\text{st}}(\mathbf{g}).$$

Постановка задачи: дистилляция Вапника

Заданы:

- 1) $\mathcal{I} = \{1, 2, \dots, m\}$;
- 2) $y_i \in \mathbb{Y} = \{1, \dots, K\}$, $\mathbb{Y}' = \mathbb{R}^K$.

Параметрическое семейство учителя:

$$\mathfrak{F}_{\text{cl}}^* = \left\{ \mathbf{f} \mid \mathbf{f} = \text{softmax}(\mathbf{v}^*(\mathbf{x}^*)/T), \quad \mathbf{v}^* : \mathbb{R}^{n^*} \rightarrow \mathbb{R}^K \right\},$$

где \mathbf{v}^* — это дифференцируемые параметрические функции заданной структуры, T — параметр температуры.

Функция ошибки:

$$\mathcal{L}_{\text{st}}(\mathbf{g}) = - \sum_{i=1}^m \sum_{k=1}^K y_i^k \log \mathbf{g}(\mathbf{x}_i) \Big|_{T=1} - \sum_{i=1}^m \sum_{k=1}^K \mathbf{f}(\mathbf{x}_i^*) \Big|_{T=T_0} \log \mathbf{g}(\mathbf{x}_i) \Big|_{T=T_0},$$

где $\cdot \Big|_{T=t}$ параметр температуры T равняется t .

Оптимизационная задача:

$$\hat{\mathbf{g}} = \arg \min_{\mathbf{g} \in \mathfrak{G}_{\text{cl}}} \mathcal{L}_{\text{st}}(\mathbf{g}).$$

Вероятностная постановка

Гипотеза порождения данных:

- 1) задано распределение целевой переменной $p(y_i | \mathbf{x}_i, \mathbf{g})$;
- 2) задано совместное распределение $p(y_i, \mathbf{s}_i | \mathbf{x}_i, \mathbf{g})$;
- 3) для всех $i \in \mathcal{I}$ элементы y_i и \mathbf{s}_i являются зависимыми величинами;
- 4) если $|\mathcal{I}| = 0$ то решение равно решению максимума правдоподобия.

Совместное правдоподобие истинных меток и меток учителя:

$$p(\mathbf{y}, \mathbf{S} | \mathbf{X}, \mathbf{g}, \mathcal{I}) = \prod_{i \notin \mathcal{I}} p(y_i | \mathbf{x}_i, \mathbf{g}) \prod_{i \in \mathcal{I}} p(y_i, \mathbf{s}_i | \mathbf{x}_i, \mathbf{g}).$$

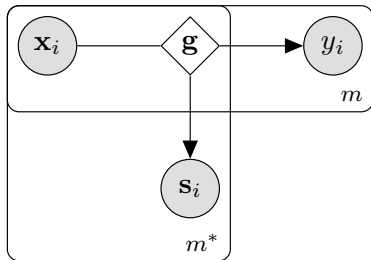
Задача оптимизации:

$$\mathbf{g} = \arg \max_{\mathbf{g} \in \mathcal{G}} p(\mathbf{y}, \mathbf{S} | \mathbf{X}, \mathbf{g}, \mathcal{I}),$$

имеет вид:

$$\sum_{i \notin \mathcal{I}} \log p(y_i | \mathbf{x}_i, \mathbf{g}) + (1 - \lambda) \sum_{i \in \mathcal{I}} \log p(y_i | \mathbf{x}_i, \mathbf{g}) \\ + \lambda \sum_{i \in \mathcal{I}} \log p(\mathbf{s}_i | \mathbf{x}_i, \mathbf{g}),$$

где $\lambda \in [0, 1]$ — метапараметр.



Частный случай: задача классификации

Заданы:

- 1) учитель $\mathbf{f} \in \mathfrak{F}_{\text{cl}}^*$ и ученик $\mathbf{g} \in \mathfrak{G}_{\text{cl}}$;
- 2) распределение истинных меток $p(y|\mathbf{x}, \mathbf{g}) = \text{Cat}(\mathbf{g}(\mathbf{x}))$;
- 3) распределение ответов учителя $p(\mathbf{s}|\mathbf{x}, \mathbf{g}) = C \prod_{k=1}^K g_k(\mathbf{x})^{s^k}$, $C < \infty$.

$$\begin{aligned} \hat{\mathbf{g}} = \arg \max_{\mathbf{g} \in \mathfrak{G}} & \sum_{i \notin \mathcal{I}} \sum_{k=1}^K y_i^k \log g_k(\mathbf{x}_i) \Big|_{T=1} + (1 - \lambda) \sum_{i \in \mathcal{I}} \sum_{k=1}^K y_i^k \log g_k(\mathbf{x}_i) \Big|_{T=1} \\ & + \lambda \sum_{i \in \mathcal{I}} \sum_{k=1}^K s_{i,k} \log g_k(\mathbf{x}_i) \Big|_{T=T_0} + \lambda \sum_{i \in \mathcal{I}} \sum_{k=1}^K \left(\log g_k(\mathbf{x}_i) \Big|_{T=T_0} + \log \log \frac{1}{g_k(\mathbf{x}_i)} \Big|_{T=T_0} \right) \end{aligned}$$

Теорема (Грабовой, 2020)

Пусть всех k выполняется $1 > 1 - \varepsilon > g_k(\mathbf{x}) > \varepsilon > 0$, тогда при

$$C = (-1)^K \frac{K^{K/2}}{2^{K(K-1)/2}} \prod_{k=1}^K g_k(\mathbf{x}) \log g_k(\mathbf{x})$$

функция $p(\mathbf{s}|\mathbf{x}, \mathbf{g}) = C \prod_{k=1}^K g_k(\mathbf{x})^{s^k}$ является плотностью распределения.

Частный случай: задача регрессии

- 1) учитель $f \in \mathfrak{F}_{rg}^* = \{f | f = \mathbf{v}^*(\mathbf{x}^*), \quad \mathbf{v}^* : \mathbb{R}^{n^*} \rightarrow \mathbb{R}\};$
- 2) ученик $g \in \mathfrak{G}_{rg} = \{g | g = \mathbf{z}(\mathbf{x}), \quad \mathbf{z} : \mathbb{R}^n \rightarrow \mathbb{R}\};$
- 3) распределение истинных меток $p(y|\mathbf{x}, g) = \mathcal{N}(y|g(\mathbf{x}), \sigma);$
- 4) распределения меток учителя $p(s|\mathbf{x}, g) = \mathcal{N}(s|g(\mathbf{x}), \sigma_s).$

Оптимизационная задача:

$$\begin{aligned} \hat{g} = \arg \min_{g \in \mathfrak{G}} \sum_{i \notin \mathcal{I}} \sigma^2 (y_i - g(\mathbf{x}_i))^2 \\ + (1 - \lambda) \sum_{i \in \mathcal{I}} \sigma^2 (y_i - g(\mathbf{x}_i))^2 + \lambda \sum_{i \in \mathcal{I}} \sigma_s^2 (s_i - g(\mathbf{x}_i))^2. \end{aligned}$$

Теорема (Грабовой, 2020)

Пусть \mathfrak{G}_{rg} — класс линейных функций $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$. Тогда решение оптимизационной задачи эквивалентно решению задачи линейной регрессии $\mathbf{y}'' = \mathbf{X}\mathbf{w} + \varepsilon$, $\varepsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$, где $\Sigma^{-1} = \text{diag}(\sigma')$ и \mathbf{y}'' имеют следующий вид:

$$\sigma'_i = \begin{cases} \sigma^2, & \text{если } i \notin \mathcal{I} \\ (1 - \lambda) \sigma^2 + \lambda \sigma_s^2, & \text{иначе,} \end{cases} \quad \mathbf{y}'' = \Sigma \mathbf{y}', \quad y'_i = \begin{cases} \sigma^2 y_i, & \text{если } i \notin \mathcal{I} \\ (1 - \lambda) \sigma^2 y_i + \lambda \sigma_s^2 s_i, & \text{иначе.} \end{cases}$$

Сводная таблица вычислительного эксперимента

Dataset	Model	CrossEntropyLoss	Accuracy	StudentSize
FashionMnist	without teacher	$0,461 \pm 0,005$	$0,841 \pm 0,002$	7850
	with teacher	$0,453 \pm 0,003$	$0,842 \pm 0,002$	7850
Synthetic	without teacher	$0,225 \pm 0,002$	$0,831 \pm 0,002$	33
	with teacher	$0,452 \pm 0,001$	$0,828 \pm 0,001$	33
Twitter	without teacher	$0,501 \pm 0,006$	$0,747 \pm 0,005$	1538
	with teacher	$0,489 \pm 0,003$	$0,764 \pm 0,004$	1538

В таблице показаны результаты вычислительного эксперимента для разных выборок. Точность аппроксимации выборки учеником улучшается при использовании модели учителя при обучении.

Байесовская дистилляция полносвязной сети

Задана модель учителя с фиксированными параметрами в виде суперпозиции отображений:

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

где \mathbf{U} матрицы линейны отображений, σ нелинейность. Вектор параметров модели учителя:

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

На основе выборки $\{\mathbf{x}_i, y_i\}_{i=1}^m$ и учителя f требуется выбрать модель ученика:

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

где \mathbf{W} , σ , \mathbf{w} вводятся аналогично учителю. Задача выбора модели g эквивалента задаче оптимизации вектора параметров \mathbf{w} .

Оптимизация параметров \mathbf{w} на основе вариационного вывода:

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

где $q(\mathbf{w})$ вариационное распределение параметров \mathbf{w} , а $p(\mathbf{w}|\mathbf{A})$ является априорным распределением вектора параметров модели ученика. Предлагается рассмотреть параметры априорного распределения $p(\mathbf{w}|\mathbf{A})$ как функцию от апостериорного распределения $p(\mathbf{u}|\mathbf{X}, \mathbf{y})$.

Построение априорного распределения

Проблема: в общем случае пространства параметров учителя \mathbf{u} и ученика \mathbf{w} разное.

Рассмотрим апостериорное распределение параметров модели учителя:

$$p(\mathbf{u}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{u}_0, \Sigma_0).$$

Пространства параметров совпадают

- ▶ число слоев совпадает $L = T$;
- ▶ размеры соответствующих слоев совпадают,

тогда $p(\mathbf{w}|\mathbf{A}) = p(\mathbf{w}|\mathcal{D})$.

Пространства параметров не совпадают В данном случае дистилляция проходит в два этапа:

- ▶ выполняется сопоставление моделей учителя и ученика;
- ▶ апостериорное распределение учителя назначается априорным распределением ученика.

Определение

Сопоставление параметрических моделей — изменение структуры модели (одной или нескольких моделей) в результате которого вектора параметров различных моделей принадлежит одному пространству параметров.

Отличие в размере скрытого слоя

- ▶ число слоев совпадает $L = T$;
- ▶ размеры соответствующих слоев не совпадают: $n_l \leq n_t$, где n_t для всех $l \in \{1, \dots, L\}, t \in \{1, \dots, T\}$.

Преобразования t -го слоя учителя:

$$\phi(t, \mathbf{u}) : \mathbb{R}^{\mathbf{p}_{\text{tr}}} \rightarrow \mathbb{R}^{\mathbf{p}_{\text{tr}} - 2n_t}$$

описывает удаление одного нейрона из t -го слоя. Новый вектор параметров обозначим $\mathbf{u}' = \phi(t, \mathbf{u})$, а выброшенные параметры обозначим \mathbf{u}''

Теорема (Грабовой, 2021)

Пусть выполняются следующие условия:

- ▶ $p(\mathbf{u}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{u}_0, \Sigma_0)$;
- ▶ число слоев совпадает $L = T$;
- ▶ для всех t, l таких, что $t = l$ выполняется $n_s \leq n_l$.

Тогда $p(\mathbf{u}'|\mathbf{X}, \mathbf{y})$ является нормальным распределением:

$$p(\mathbf{u}'|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{u}'_0 + \Sigma_0'{}'' \Sigma_0''^{-1} (\mathbf{0} - \mathbf{u}''_0), \Sigma_0' - \Sigma_0'{}'' \Sigma_0''^{-1} \Sigma_0'{}'').$$

Отличие в числе слоев нейросети

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

Преобразования t -го слоя учителя:

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

описывает удаление t -го слоя. Новый вектор параметров обозначим $\mathbf{u}' = \phi(t, \mathbf{u})$, а выброшенные параметры обозначим \mathbf{u}''

Теорема (Грабовой, 2021)

Пусть выполняются следующие условия:

- ▶ $p(\mathbf{u}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{u}_0, \Sigma_0)$;
- ▶ соответствующие размеры слоев совпадают, $n_t = n_{t-1}$;
- ▶ функция активации удовлетворяет свойству $\sigma \circ \sigma = \sigma$.

Тогда $p(\mathbf{u}'|\mathbf{X}, \mathbf{y})$ является нормальным распределением:

$$p(\mathbf{u}'|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{u}'_0 + \Sigma_0'^{''} \Sigma_0''^{-1} (\mathbf{i} - \mathbf{u}_0''), \Sigma_0' - \Sigma_0'^{''} \Sigma_0''^{-1} \Sigma_0'^{''}),$$

где $\mathbf{i} = \text{vec}(\mathbf{I})$

Порядок на множестве параметров

Задание порядка на множестве параметров:

- ▶ случайным образом (используется в рамках вычислительного эксперимента)
- ▶ на основе метода оптимального прореживания нейросети:

$$\xi = \arg \min_j h_{jj} \frac{u_j^2}{2},$$

где h_{jj} коэффициент при квадратичном члене в разложении Тейлора функции ошибки по параметрам модели.

- ▶ на основе отношения плотности апостериорного распределения параметра к плотности апостериорного распределения параметра к нулю:

$$\xi = \arg \max_j \frac{p(0|\mathbf{X}, \mathbf{t})}{p(u_j|\mathbf{X}, \mathbf{t})},$$

- ▶ выбор на основе анализа мультиколлинеарности параметров методов Белсли:

$$\xi = \arg \max_j \frac{\lambda_{\max}}{\lambda_j},$$

где λ являются сингулярными числами ковариационной матрицы параметров.

Вычислительный эксперимент

Синтетическая выборка:

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

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

Модель учителя:

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

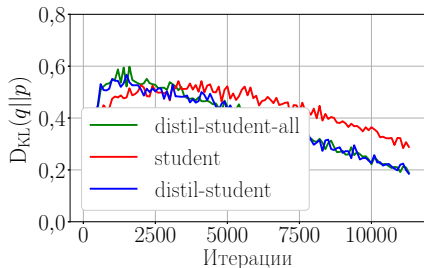
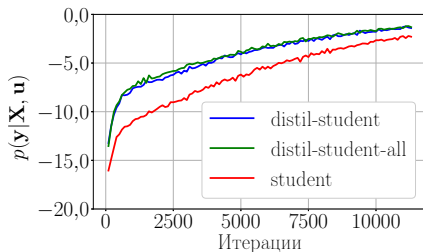
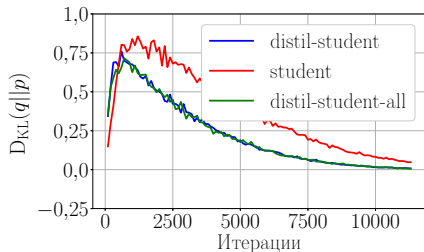
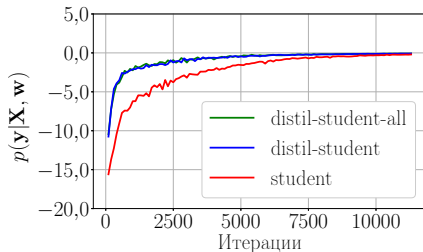
Первая конфигурация модели ученика:

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

Вторая конфигурация модели ученика:

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

Результаты экспериментов



Сводная таблица вычислительного эксперимента

	teacher	student	distil-student	distil-student-all
Эксперимент на синтетической выборке (удаление нейрона)				
Архитектура	[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)				
Архитектура	[10, 100, 50, 1]	[10, 50, 1]	[10, 50, 1]	[10, 50, 1]
Число параметро				
Разность площадей	-	0	23310	25506
Эксперимент на выборке FashionMnist				
Архитектура	[784, 800, 50, 10]	[784, 50, 10]	[784, 50, 10]	[784, 50, 10]
Число параметро				
Разность площадей	-	0	1165	1145

Выносятся на защиту

1. Проведен вероятностный анализ задачи дистилляции.
2. Выполнено обобщение классического подхода введя вероятностные предположения о природе данных.
3. Теоретические анализы сформулированы в виде теорем для задачи классификации и регрессии.
4. Поставлена задача байесовской дистилляции моделей глубокого обучения.
5. Предложен метод задания априорного распределения параметров модели ученика на основе апостериорного распределения параметров учителя.
6. Доказаны теоремы, которые позволяют проводить сведение структуры модели учителя к структуре модели ученика.
7. Проведен ряд вычислительных экспериментов, которые показывают применимость предложенных методов.

Публикации ВАК по теме

1. Грабовой А.В., Бахтеев О.Ю., Стрижов В.В. Определение релевантности параметров нейросети // Информатика и ее применения, 2019, 13(2).
2. Грабовой А.В., Бахтеев О. Ю., Стрижов В.В. Введение отношения порядка на множестве параметров аппроксимирующих моделей // Информатика и ее применения, 2020, 14(2).
3. A. Grabovoy, V. Strijov. Quasi-periodic time series clustering for human. Lobachevskii Journal of Mathematics, 2020, 41(3).
4. Грабовой А.В., Стрижов В.В. Анализ выбора априорного распределения для смеси экспертов // Журнал Вычислительной математики и математической физики, 2021. 61(5).
5. Грабовой А.В., Стрижов В.В. Анализ моделей привилегированного обучения и дистилляции // Автоматика и телемеханика, 2021 (на рецензировании).
6. T. Gadaev, A. Grabovoy, A. Motrenko, V. Strijov Numerical methods of minimum sufficient sample size estimation for linear models // in progress.
7. Базарова А.И., Грабовой А.В., Стрижов В.В. Анализ свойств вероятностных моделей в задачах обучения с экспертом // подано.
8. Грабовой А.В., Стрижов В.В. Байесовская дистилляция моделей глубокого обучения // (текущая работа).