

## Анализ моделей привилегированного обучения и дистилляции<sup>1</sup>

Данная работа посвящена методам понижения сложности аппроксимирующих моделей. Предлагается вероятностное обоснование методов дистилляции и привилегированного обучения. В работе приведены общие выводы для произвольной параметрической функции с наперед заданной структурой. Показано теоретическое обоснование для частных случаев: линейной и логистической регрессии. Теоретические результаты анализируются в вычислительном эксперименте на синтетических выборках и реальных данных. В качестве реальных данных рассматривается выборки FashionMNIST и Twitter Sentiment Analysis.

*Ключевые слова:* выбор модели; байесовский вывод; дистилляция модели; привилегированное обучение.

### 1. Введение

Повышение точности аппроксимации в задачах машинного обучения влечет за собой повышение сложности моделей и как следствие снижает их интерпретируемость. Примером такого усложнения являются следующие модели: трансформеры [1], BERT [2], ResNet [3] а также ансамбли этих моделей.

При построении модели машинного обучения используется два свойства: сложность модели и точность аппроксимации модели. Сложность влияет на время, которое модель требуется для принятия решения, а также на интерпретируемость модели, следовательно модель которая имеет меньшую сложность является более предпочтительной [4]. С другой стороны точность аппроксимации модели нужно максимизировать. В данной работе рассматривается метод *дистилляции* модели. Данный метод позволяет строить новые модели на основе ранее обученных моделей.

*Определение 1. Дистилляция модели — уменьшение сложности модели путем выбора модели в множестве более простых моделей с использованием ответов более сложной модели.*

В работе [5] Дж. Е. Хинтоном рассматривается метод дистилляции моделей машинного обучения для задачи классификации. В работе проведен ряд экспериментов, в которых проводилась дистилляция моделей для разных задач машинного обучения. Эксперимент на выборке MNIST [6], в котором избыточно сложно нейросеть была дистиллирована в нейросеть меньшей сложности. Эксперимент по Speech Recognition, в котором ансамбль моделей был *дистиллирован* в одну модель. Также в работе [5] был проведен эксперимент по обучению экспертных моделей на основе одной большой модели.

---

<sup>1</sup>Работа выполнена при поддержке ... (грант №...).

*Определение 2. Привилегированная информация — множество признаков, которые доступны только в момент выбора модели, но не в момент тестирования.*

В работе [7] В. Н. Вапником введено понятия *привилегированной информации*. В работе [8] метод дистилляции [5] используется вместе с привилегированным обучением [7]. В предложенном методе на первом этапе обучается модель *учителя* в пространстве привилегированной информации, после чего обучается модель *ученика* в исходном признаковом пространстве используя *дистилляцию* [5]. Для обучения строится функция ошибки специального вида, анализируемая в данной работе. Эта функция состоит из нескольких слагаемых, включая ошибки учителя, ученика и регуляризирующие элементы. Первые варианты подобной функции ошибки были предложены А. Г. Ивахненко [9].

*Определение 3. Учитель — фиксируемая модель, ответы которой используются при выборе модели ученика.*

*Определение 4. Ученик — модель, которая выбирается согласно какого-либо критерия.*

В данной работе предлагается рассмотреть вероятностный подход к решению задачи дистилляции модели и задачи привилегированного обучения. Подход обобщается на случай, когда привилегированная информация доступна не для всех объектов из обучающей выборки. В рамках вероятностного подхода предлагается анализ и обобщение функции ошибки [5, 8]. Рассматриваются частные задачи классификации и регрессии [9].

В рамках вычислительного эксперимента анализируются методы использующие и не использующие модель учителя при обучении модели ученика. Для анализа используются реальные выборки для задачи классификации изображений FashionMNIST [10] и для задачи классификации текстов Twitter Sentiment Analysis [11]. Выборка FashionMNIST использовалась вместо общепринятой выборки MNIST, так как последняя имеет приемлемое качество аппроксимации даже для линейного классификатора. Вычислительный эксперимент использует модели разной сложности: линейная модель, полносвязная нейронная сеть, сверточная нейронная сеть [12], модель Bi-LSTM [13] и модель BERT [2].

## 2. Постановка задачи обучения с учителем

Задано множество объектов  $\Omega$  и множество целевых переменных  $\mathbb{Y}$ . Множество  $\mathbb{Y} = \{1, \dots, K\}$  для задачи классификации, где  $K$  число классов, множество  $\mathbb{Y} = \mathbb{R}$  для задачи регрессии. Для каждого объекта из  $\omega_i \in \Omega$  задана целевая переменная  $\mathbf{y}_i = \mathbf{y}(\omega_i)$ . Множество целевых переменных для всех объектов обозначим  $\mathbf{Y}$ . Для множества  $\Omega$  задано отображение в некоторое признаковое пространство  $\mathbb{R}^n$ :

$$\varphi : \Omega \rightarrow \mathbb{R}^n, \quad |\Omega| = m,$$

где  $n$  размерность признакового пространства, а  $m$  количество объектов в множестве  $\Omega$ . Отображение  $\varphi$  отображает объект  $\omega_i \in \Omega$  в соответствующий ему вектор признаков  $\mathbf{x}_i = \varphi(\omega_i)$ . Пусть для объектов  $\Omega^* \subset \Omega$  задана привилегированная информация:

$$\varphi^* : \Omega^* \rightarrow \mathbb{R}^{n^*}, \quad |\Omega^*| = m^*,$$

где  $m^* \leq m$  — число объектов с привилегированной информацией,  $n^*$  — число признаков в пространстве привилегированной информации. Отображение  $\varphi^*$  отображает объект  $\omega_i \in \Omega^*$  в соответствующий ему вектор признаков  $\mathbf{x}_i^* = \varphi^*(\omega_i)$ .

Множество индексов объектов, для которых известна привилегированная информация, обозначим  $\mathcal{I}$ :

$$\mathcal{I} = \{1 \leq i \leq m \mid \text{для } i\text{-го объекта задана привилегированная информация}\},$$

а множество индексов объектов, для которых не известна привилегированная информация, обозначим  $\{1, \dots, m\} \setminus \mathcal{I} = \bar{\mathcal{I}}$ .

Пусть на множестве привилегированных признаков задана функция учителя  $\mathbf{f}(\mathbf{x}^*)$ :

$$\mathbf{f} : \mathbb{R}^{n^*} \rightarrow \mathbb{Y}^*,$$

где  $\mathbb{Y}^* = \mathbb{Y}$  для задачи регрессии и  $\mathbb{Y}^*$  является единичным симплексом  $\mathcal{S}_K$  в пространстве размерности  $K$  для задачи классификации. Модель учителя  $\mathbf{f}$  ставит объекты  $\mathbf{X}^*$  в соответствие объектам  $\mathbf{S}$ , то есть  $\mathbf{f}(\mathbf{x}_i^*) = \mathbf{s}_i$ .

Требуется выбрать модель ученика  $\mathbf{g}(\mathbf{x})$  из множества:

$$(1) \quad \mathfrak{G} = \{\mathbf{g} \mid \mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{Y}^*\},$$

например для задачи классификации множество  $\mathfrak{G}$  может быть параметрическим семейством функций линейных моделей:

$$\mathfrak{G}_{\text{lin,cl}} = \{\mathbf{g}(\mathbf{W}, \mathbf{x}) \mid \mathbf{g}(\mathbf{W}, \mathbf{x}) = \text{softmax}(\mathbf{W}\mathbf{x}), \quad \mathbf{W} \in \mathbb{R}^{n \times K}\}.$$

### 3. Постановка задачи: Хинтона и Вапника

Рассмотрим описание метода предложенного в работах [5, 8]. В рамках данных работ предполагается, что для всех данных доступна привилегированная информация  $\mathcal{I} = \{1, 2, \dots, m\}$ . В работе [5] решается задача классификации вида:

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

где  $y_i$  — это класс объекта, также будем обозначать  $\mathbf{y}_i$  вектором вероятности для класса  $y_i$ .

В постановке Хинтона рассматривается параметрическое семейство функций:

$$(2) \quad \mathfrak{G}_{\text{cl}} = \{\mathbf{g} \mid \mathbf{g} = \text{softmax}(\mathbf{z}(\mathbf{x})/T), \quad \mathbf{z} : \mathbb{R}^n \rightarrow \mathbb{R}^K\},$$

где  $\mathbf{z}$  — это дифференцируемая параметрическая функция заданной структуры,  $T$  — параметр температуры. В качестве модели учителя  $\mathbf{f}$  рассматривается функция из множества  $\mathfrak{F}_{\text{cl}}$ :

$$(3) \quad \mathfrak{F}_{\text{cl}} = \{\mathbf{f} \mid \mathbf{f} = \text{softmax}(\mathbf{v}(\mathbf{x})/T), \quad \mathbf{v} : \mathbb{R}^n \rightarrow \mathbb{R}^K\},$$

где  $\mathbf{v}$  — это дифференцируемая параметрическая функция заданной структуры,  $T$  — параметр температуры. Параметр температуры  $T$  имеет следующие свойства:

1) при  $T \rightarrow 0$  получаем вектор, в котором один из классов имеет единичную вероятность;

2) при  $T \rightarrow \infty$  получаем равновероятные классы.

Функция потерь  $\mathcal{L}$  в которой учитывается перенос информации от модели учителя  $\mathbf{f}$  к модели ученика  $\mathbf{g}$  имеет следующий вид:

$$(4) \quad \mathcal{L}_{st}(\mathbf{g}) = - \underbrace{\sum_{i=1}^m \sum_{k=1}^K y_i^k \log \mathbf{g}(\mathbf{x}_i)}_{\text{исходная функция потерь}} \Big|_{T=1} - \underbrace{\sum_{i=1}^m \sum_{k=1}^K \mathbf{f}(\mathbf{x}_i) \Big|_{T=T_0} \log \mathbf{g}(\mathbf{x}_i) \Big|_{T=T_0}}_{\text{слагаемое дистилляция}},$$

где  $\cdot|_{T=t}$  обозначает, что параметр температуры  $T$  в предыдущей функции равняется  $t$ .

Получаем оптимизационную задачу:

$$(5) \quad \hat{\mathbf{g}} = \arg \min_{\mathbf{g} \in \mathcal{G}_{cl}} \mathcal{L}_{st}(\mathbf{g}).$$

Работа [8] обобщает метод предложенный в работе [5]. Решение задачи оптимизации (5) зависит только от вектора ответов модели учителя  $\mathbf{f}$ . Следовательно признаковые пространства учителя и ученика могут различаться. В этом случае получаем следующую постановку задачи:

$$\mathcal{D} = \{(\mathbf{x}_i, \mathbf{x}_i^*, y_i)\}_{i=1}^m, \quad \mathbf{x}_i \in \mathbb{R}^n, \quad \mathbf{x}_i^* \in \mathbb{R}^{n^*}, \quad y_i \in \{1, \dots, K\},$$

где  $\mathbf{x}_i$  это информация доступна на этапах обучения и контроля, а  $\mathbf{x}_i^*$  это информация доступна только на этапе обучения. Модель учителя принадлежит множеству моделей  $\mathfrak{F}_{cl}^*$ :

$$(6) \quad \mathfrak{F}_{cl}^* = \{\mathbf{f} | \mathbf{f} = \text{softmax}(\mathbf{v}^*(\mathbf{x}^*)/T), \quad \mathbf{v}^* : \mathbb{R}^{n^*} \rightarrow \mathbb{R}^K\},$$

где  $\mathbf{v}^*$  — это дифференцируемая параметрическая функция заданной структуры,  $T$  — параметр температуры. Множество моделей  $\mathfrak{F}_{cl}^*$  отличается от множества моделей  $\mathfrak{F}_{cl}$  из выражения (3). В множестве  $\mathfrak{F}_{cl}$  модели используют пространство исходных признаков, а в множестве  $\mathfrak{F}_{cl}^*$  модели используют пространство привилегированных признаков. Функция потерь (4) в случае модели учителя  $\mathbf{f} \in \mathfrak{F}_{cl}^*$  переписывается в следующем виде:

$$(7) \quad \mathcal{L}_{st}(\mathbf{g}) = - \sum_{i=1}^m \sum_{k=1}^K y_i^k \log \mathbf{g}(\mathbf{x}_i) \Big|_{T=1} - \sum_{i=1}^m \sum_{k=1}^K \mathbf{f}(\mathbf{x}_i^*) \Big|_{T=T_0} \log \mathbf{g}(\mathbf{x}_i) \Big|_{T=T_0},$$

где  $\cdot|_{T=t}$  обозначает, что параметр температуры  $T$  в предыдущей функции равняется  $t$ .

Требуется построить модель, которая использует привилегированную информацию  $\mathbf{x}_i^*$  при обучении. Для этого рассмотрим двухэтапную модель обучения предложенную в работе [8]:

- 1) выбираем оптимальную модель учителя  $\mathbf{f} \in \mathfrak{F}_{\text{cl}}^*$ ;
- 2) выбираем оптимальную модель ученика  $\mathbf{g} \in \mathfrak{G}_{\text{cl}}$  используя дистилляцию [5].

Модель ученика — это функция, которая минимизирует (7). Модель учителя — это функция, которая минимизирует кросс-энтропийную функции ошибки:

$$\mathcal{L}_{th}(\mathbf{f}) = - \sum_{i=1}^m \sum_{k=1}^K y_i^k \log \mathbf{f}(\mathbf{x}_i^*).$$

#### 4. Постановка задачи: вероятностный подход

##### 4.1. Метод максимального правдоподобия

Задано распределения целевой переменной  $p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{g})$ . Для поиска  $\hat{\mathbf{g}}$  воспользуемся методом максимального правдоподобия. В качестве  $\hat{\mathbf{g}}$  выбирается функция, которая максимизирует правдоподобие модели:

$$(8) \quad \hat{\mathbf{g}} = \arg \max_{\mathbf{g} \in \mathfrak{G}} \prod_{i=1}^N p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{g}),$$

где множество  $\mathfrak{G}$  задается (1).

##### 4.2. Подход дистилляции модели учителя в модель ученика

Рассмотрим вероятностную постановку, в которой должны быть выполнены ограничения:

- 1) задано распределение целевой переменной  $p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{g})$ ;
- 2) задано совместное распределение целевой переменной и ответов модели учителя  $p(\mathbf{y}_i, \mathbf{s}_i|\mathbf{x}_i, \mathbf{g})$ ;
- 3) для всех  $\omega \in \Omega^*$  элементы  $\mathbf{y}(\omega)$  и  $\mathbf{s}(\omega)$  являются зависимыми величинами, так как ответы учителя должны коррелировать с истинными ответами;
- 4) если  $|\Omega^*| = 0$  то решение должно соответствовать решению (8).

Рассмотрим совместное правдоподобие истинных меток и меток учителя:

$$(9) \quad p(\mathbf{Y}, \mathbf{S}|\mathbf{X}, \mathbf{g}, \mathcal{I}) = \prod_{i \notin \mathcal{I}} p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{g}) \prod_{i \in \mathcal{I}} p(\mathbf{y}_i, \mathbf{s}_i|\mathbf{x}_i, \mathbf{g}).$$

Распишем  $p(\mathbf{y}_i, \mathbf{s}_i|\mathbf{x}_i, \mathbf{g})$  по формуле условной вероятности:

$$(10) \quad p(\mathbf{y}_i, \mathbf{s}_i|\mathbf{x}_i, \mathbf{g}) = p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{g}) p(\mathbf{s}_i|\mathbf{y}_i, \mathbf{x}_i, \mathbf{g})$$

Подставляя выражения (10) в (9) получаем.

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

Заметим, что  $\mathbf{y}_i$  и  $\mathbf{s}_i$  зависимы только через переменную  $\mathbf{x}_i$ , тогда  $p(\mathbf{s}_i|\mathbf{y}_i, \mathbf{x}_i, \mathbf{g}) = p(\mathbf{s}_i|\mathbf{x}_i, \mathbf{g})$ . Получаем совместное правдоподобие:

$$(11) \quad p(\mathbf{Y}, \mathbf{S}|\mathbf{X}, \mathbf{g}, \mathcal{I}) = \prod_{i \notin \mathcal{I}} p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{g}) \prod_{i \in \mathcal{I}} p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{g}) \prod_{i \in \mathcal{I}} p(\mathbf{s}_i|\mathbf{x}_i, \mathbf{g}).$$

Используя (11) получаем следующую оптимизационную задачу для поиска  $\hat{\mathbf{g}}$

$$(12) \quad \hat{\mathbf{g}} = \arg \max_{\mathbf{g} \in \mathcal{G}} \prod_{i \notin \mathcal{I}} p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{g}) \prod_{i \in \mathcal{I}} p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{g}) \prod_{i \in \mathcal{I}} p(\mathbf{s}_i | \mathbf{x}_i, \mathbf{g}).$$

Для удобства, будем минимизировать логарифм, тогда из (12) получаем:

$$(13) \quad \hat{\mathbf{g}} = \arg \max_{\mathbf{g} \in \mathcal{G}} \sum_{i \notin \mathcal{I}} \log p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{g}) + (1 - \lambda) \sum_{i \in \mathcal{I}} \log p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{g}) + \lambda \sum_{i \in \mathcal{I}} \log p(\mathbf{s}_i | \mathbf{x}_i, \mathbf{g}),$$

где параметр  $\lambda \in [0, 1]$  введен для взвешивания ошибок на истинных ответах и ошибок относительно ответов учителя.

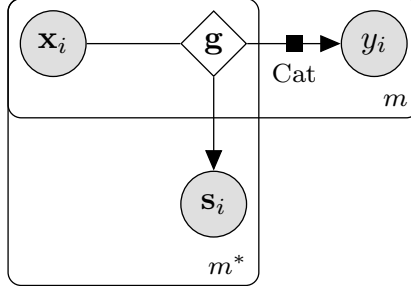


Рис. 1. Вероятностная модель в формате плоских нотаций.

На рис. 1 показан вид вероятностной модели в графовой нотации, для произвольной функции  $\mathbf{g}$ . Для каждой реализации  $\mathbf{g}$  соответствующий блок требует уточнение. На рис 3 показана более подробная реализация в случае, когда модель  $\mathbf{g}$  это линейная модель.

## 5. Обучение с учителем для задачи классификации и регрессии

### 5.1. Случай классификации

Для задачи многоклассовой классификации рассматриваются следующие вероятностные предположения:

- 1) рассматривается функция учителя  $\mathbf{f} \in \mathfrak{F}_{\text{cl}}^*$  (6);
- 2) рассматривается функция ученика следующего вида  $\mathbf{g} \in \mathfrak{G}_{\text{cl}}$  (2);
- 3) для истинных меток рассматривается категориальное распределение  $p(y | \mathbf{x}, \mathbf{g}) = \text{Cat}(\mathbf{g}(\mathbf{x}))$ , где  $\mathbf{g}(\mathbf{x})$  задает вероятность каждого класса;
- 4) для меток учителя введем плотность распределения

$$(14) \quad p(\mathbf{s} | \mathbf{x}, \mathbf{g}) = C \prod_{k=1}^K g_k(\mathbf{x})^{s^k},$$

где  $g^k$  обозначает вероятность класса  $k$ , которую предсказывает модель ученика, а  $s^k$  — вероятность класса  $k$ , которую предсказывает модель учителя.

**Теорема 1.** Пусть вероятность каждого класса отделима от нуля и единицы, то есть для 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  $p(\mathbf{s}|\mathbf{x}, \mathbf{g})$  определенная в (14) является плотностью распределения.

**Доказательство.** Во-первых покажем, что для произвольного вектора ответов  $\mathbf{s} \in \mathcal{S}_K$  выполняется  $p(\mathbf{s}|\mathbf{x}, \mathbf{g}) \geq 0$ . Заметим, что для всех  $k$  выполняется, что  $\log g_k(\mathbf{x}) < 0$ , тогда

$$C = \underbrace{\frac{K^{K/2}}{2^{K(K-1)/2}}}_{>0} \prod_{k=1}^K \underbrace{g_k(\mathbf{x})}_{>\varepsilon} \underbrace{(-\log g_k(\mathbf{x}))}_{>0} > 0,$$

тогда с учетом того, что  $g_k(\mathbf{x}) > 0$  и  $C > 0$  получаем, что  $p(\mathbf{s}|\mathbf{x}, \mathbf{g}) \geq 0$ . Во-вторых покажем, что интеграл по всему пространству ответов  $\mathcal{S}_K$  является конечным:

(15)

$$\begin{aligned} \int_{\mathcal{S}_K} p(\mathbf{s}|\mathbf{x}, \mathbf{g}) ds &= \int_{\mathcal{S}_K} \prod_{k=1}^K g_k(\mathbf{x})^{s_k} ds = \prod_{k=1}^K \int_{\mathcal{S}_K} g_k(\mathbf{x})^{s_k} ds \\ &= \prod_{k=1}^K \int_0^1 \frac{r^{K-1} \sqrt{K}}{(K-1)! \sqrt{2^{K-1}}} g_k(\mathbf{x})^r dr = \prod_{k=1}^K \underbrace{\frac{\sqrt{K}}{(K-1)! \sqrt{2^{K-1}}}}_D \int_0^1 r^{K-1} g_k(\mathbf{x})^r dr \\ &= D^K \prod_{k=1}^K \int_0^1 r^{K-1} \exp(r \log g_k(\mathbf{x})) dr \\ &= (-D)^K \prod_{k=1}^K \log g_k(\mathbf{x}) (\Gamma(K) - \Gamma(K, -\log g_k(\mathbf{x}))) \\ &= (-D)^K (K-1)!^K \prod_{k=1}^K \log g_k(\mathbf{x}) (1 - g_k(\mathbf{x}) \exp_{K-1}(-\log g_k(\mathbf{x})) + g_k(\mathbf{x})) \\ &= \frac{(-\sqrt{K})^K}{2^{K(K-1)/2}} \prod_{k=1}^K \log g_k(\mathbf{x}) (1 - g_k(\mathbf{x}) \exp_{K-1}(-\log g_k(\mathbf{x})) + g_k(\mathbf{x})) < \infty, \end{aligned}$$

где  $\Gamma(K)$  является гамма функцией,  $\Gamma(K, -\log g_k(\mathbf{x}))$  является неполной гамма функцией,  $\exp_n(x)$  является суммой Тейлора из первых  $n$  слагаемых. В рамках приближенных расчетов будем считать, что  $\exp_n(x) \approx \exp(x)$ , тогда с учетом (15) получаем:

$$(16) \quad C(\mathbf{g}, \mathbf{x}) = \int_{\mathcal{S}_K} p(\mathbf{s}|\mathbf{x}, \mathbf{g}) ds \approx (-1)^K \frac{K^{K/2}}{2^{K(K-1)/2}} \prod_{k=1}^K g_k(\mathbf{x}) \log g_k(\mathbf{x})$$

Полученное выражение (16) заканчивает доказательство теоремы.

Из теоремы 1 следует, что плотность введенная для меток учителя является плотностью распределения, следовательно можно воспользоваться выражением (13). Используя предположения 1)–4) и подставляя в (13) получаем следующую оптимиза-

ционную задачу:

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

где выражение  $\cdot \Big|_{T=t}$  обозначает, что в предыдущую функцию softmax требуется подставить значение температуры  $T$  равное некоторому значению  $t$ .

Проанализировав выражение (17) получаем, что первые три слагаемые совпадают со слагаемыми в выражении (4) при  $\mathcal{I} = \{1, \dots, m\}$ , и  $\lambda = \frac{1}{2}$ , а третье слагаемое является некоторым регуляризатором, который получен из вида распределения.

Анализируя первые три слагаемых в выражении (17) получаем, что при  $T_0 = 1$  получаем сумму кросс энтропий между двумя распределениями для каждого объекта:

1) первое распределение это выпуклая комбинация с весом  $1 - \lambda$  и  $\lambda$ : распределения задаваемое метками объектов  $\text{Cat}(\mathbf{y})$  и распределения задаваемого моделью учителя  $\text{Cat}(\mathbf{s})$

2) второе распределение это распределение задаваемое моделью ученика  $\text{Cat}(\mathbf{g}(\mathbf{x}))$ .

Получаем, что модель ученика восстанавливает плотность не исходных меток, а новую плотность, которая является выпуклой комбинаций плотности исходных меток и меток учителя.

## 5.2. Случай регрессии

Для задачи регрессии рассматриваются следующие вероятностные предположения:

1) рассматривается функция учителя  $\mathbf{f} \in \mathfrak{F}_{\text{rg}}^*$ :

$$\mathfrak{F}_{\text{rg}}^* = \{ \mathbf{f} | \mathbf{f} = \mathbf{v}^*(\mathbf{x}^*), \quad \mathbf{v}^* : \mathbb{R}^{n^*} \rightarrow \mathbb{R} \},$$

где  $\mathbf{v}^*$  — это дифференцируемая параметрическая функция;

2) рассматривается функция ученика  $\mathbf{g} \in \mathfrak{G}_{\text{rg}}$ :

$$\mathfrak{G}_{\text{rg}} = \{ \mathbf{g} | \mathbf{g} = \mathbf{z}(\mathbf{x}), \quad \mathbf{z} : \mathbb{R}^n \rightarrow \mathbb{R}^K \},$$

где  $\mathbf{z}$  — это дифференцируемая параметрическая функция;

3) истинные метки имеют нормальное распределение

$$p(y|\mathbf{x}, \mathbf{g}) = \mathcal{N}(y|\mathbf{g}(\mathbf{x}), \sigma);$$

4) метки учителя распределены

$$p(s|\mathbf{x}, \mathbf{g}) = \mathcal{N}(s|\mathbf{g}(\mathbf{x}), \sigma_s);$$



Используя предположения 1)–4) и подставляя в (13) получаем следующую оптимизационную задачу:

$$(18) \quad \begin{aligned} \hat{g} = \arg \min_{g \in \mathcal{G}} \sum_{i \notin \mathcal{I}} \sigma^2 (y_i - \mathbf{g}(\mathbf{x}_i))^2 \\ + (1 - \lambda) \sum_{i \in \mathcal{I}} \sigma^2 (y_i - \mathbf{g}(\mathbf{x}_i))^2 + \lambda \sum_{i \in \mathcal{I}} \sigma_s^2 (s_i - \mathbf{g}(\mathbf{x}_i))^2. \end{aligned}$$

Выражение (18) записано с точностью до аддитивной константы относительно  $\mathbf{g}$ .

*Теорема 2.* Пусть множество  $\mathcal{G}$  описывает класс линейных функций вида  $\mathbf{g}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ . Тогда решение оптимизационной задачи (18) эквивалентно решению следующей задачи линейной регрессии:

$$(19) \quad \mathbf{y}'' = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}),$$

где  $\boldsymbol{\Sigma}^{-1} = \text{diag}(\boldsymbol{\sigma}')$  и  $\mathbf{y}''$  имеют следующий вид:

$$(20) \quad \begin{aligned} \sigma'_i &= \begin{cases} \sigma^2, & \text{если } i \notin \mathcal{I} \\ (1 - \lambda) \sigma^2 + \lambda \sigma_s^2, & \text{иначе} \end{cases}, \\ \mathbf{y}'' &= \boldsymbol{\Sigma} \mathbf{y}', \\ y'_i &= \begin{cases} \sigma^2 y_i, & \text{если } i \notin \mathcal{I} \\ (1 - \lambda) \sigma^2 y_i + \lambda \sigma_s^2 s_i, & \text{иначе} \end{cases}. \end{aligned}$$

*Доказательство.* Обозначим  $\mathbf{a}_{\mathcal{J}} = [a_i | i \in \mathcal{J}]^\top$ , где  $\mathbf{a}$  произвольный вектор, а  $\mathcal{J}$  произвольное не пустое индексное множество. Подвектор вектора ответов  $\mathbf{y}$ , для элементов которого доступна привилегированная информация обозначим  $\mathbf{y}_{\mathcal{I}} = [y_i | i \in \mathcal{I}]^\top$ . Аналогично обозначим матрицу  $\mathbf{X}_{\mathcal{I}} = [\mathbf{x}_i | i \in \mathcal{I}]^\top$ .

В случае линейной модели  $\mathbf{g}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  выражение (18) принимает вид:

$$\begin{aligned} \hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathcal{W}} \sigma^2 (\mathbf{y}_{\bar{\mathcal{I}}} - \mathbf{X}_{\bar{\mathcal{I}}} \mathbf{w})^\top (\mathbf{y}_{\bar{\mathcal{I}}} - \mathbf{X}_{\bar{\mathcal{I}}} \mathbf{w}) \\ + \sigma^2 (1 - \lambda) (\mathbf{y}_{\mathcal{I}} - \mathbf{X}_{\mathcal{I}} \mathbf{w})^\top (\mathbf{y}_{\mathcal{I}} - \mathbf{X}_{\mathcal{I}} \mathbf{w}) + \sigma_s^2 \lambda (\mathbf{s}_{\mathcal{I}} - \mathbf{X}_{\mathcal{I}} \mathbf{w})^\top (\mathbf{s}_{\mathcal{I}} - \mathbf{X}_{\mathcal{I}} \mathbf{w}). \end{aligned}$$

Раскроем скобки и сгруппируем:

$$\begin{aligned} \hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathcal{W}} \sigma^2 (\mathbf{w}^\top \mathbf{X}_{\bar{\mathcal{I}}}^\top \mathbf{X}_{\bar{\mathcal{I}}} \mathbf{w} - 2 \mathbf{y}_{\bar{\mathcal{I}}}^\top \mathbf{X}_{\bar{\mathcal{I}}} \mathbf{w}) \\ + (1 - \lambda) \sigma^2 (\mathbf{w}^\top \mathbf{X}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}} \mathbf{w} - 2 \mathbf{y}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}} \mathbf{w}) + \lambda \sigma_s^2 (\mathbf{w}^\top \mathbf{X}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}} \mathbf{w} - 2 \mathbf{s}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}} \mathbf{w}) \end{aligned}$$

Продифференцируем выражение, приравняем к нулю и сгруппируем элементы:

$$(21) \quad \begin{aligned} (\sigma^2 \mathbf{X}_{\bar{\mathcal{I}}}^\top \mathbf{X}_{\bar{\mathcal{I}}} + (1 - \lambda) \sigma^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}} + \lambda \sigma_s^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}}) \mathbf{w} = 2 \sigma^2 \mathbf{X}_{\bar{\mathcal{I}}}^\top \mathbf{y}_{\bar{\mathcal{I}}} \\ + 2 (1 - \lambda) \sigma^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{y}_{\mathcal{I}} + 2 \lambda \sigma_s^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{s}_{\mathcal{I}}. \end{aligned}$$

Воспользуемся следующими равенствами:

$$(22) \quad \begin{aligned} \sigma^2 \mathbf{X}_{\bar{\mathcal{I}}}^\top \mathbf{X}_{\bar{\mathcal{I}}} + (1 - \lambda) \sigma^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}} + \lambda \sigma_s^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{X}_{\mathcal{I}} &= \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}, \\ 2 \sigma^2 \mathbf{X}_{\bar{\mathcal{I}}}^\top \mathbf{y}_{\bar{\mathcal{I}}} + 2 (1 - \lambda) \sigma^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{y}_{\mathcal{I}} + 2 \lambda \sigma_s^2 \mathbf{X}_{\mathcal{I}}^\top \mathbf{s}_{\mathcal{I}} &= 2 \mathbf{X} \mathbf{y}', \end{aligned}$$

где  $\Sigma$  и  $\mathbf{y}'$  из условия задачи (20).

Подставляя (22) в (21) получаем:

$$\mathbf{w} = 2 (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X} \Sigma^{-1} \mathbf{y}'',$$

что соответствует решению задачи (19).

Теорема 2 показывает, что обучения с учителем для задачи регрессии можно свести к классической задаче оптимизации для задачи линейной регрессии.

## 6. Вычислительный эксперимент

Проводится вычислительный эксперимент для анализа качества моделей, которые получены путем дистилляции модели учителя в модель ученика. Как показано в теореме 2 задачу регрессии с учителем можно свести к задаче регрессии без учителя, поэтому в эксперименте более подробно рассматривается случай классификации. Во всех частях вычислительного эксперимента для поиска оптимальных параметров нейросетей использовался градиентный метод оптимизации Адам [14].

### 6.1. Выборка FashionMNIST

В данной части проводится эксперимент для задачи классификации для выборки FashionMNIST [10]. В качестве модели учителя  $\mathbf{f}$  рассматривается модель нейросети с двумя сверточными слоями и с тремя полносвязными слоями, в качестве функции активации рассматривается ReLu. Модель учителя содержит 30 тысяч обучаемых параметров. В качестве модели ученика рассматривается модель логистической регрессии для многоклассовой классификации. Модель ученика содержит 7850 обучаемых параметров.

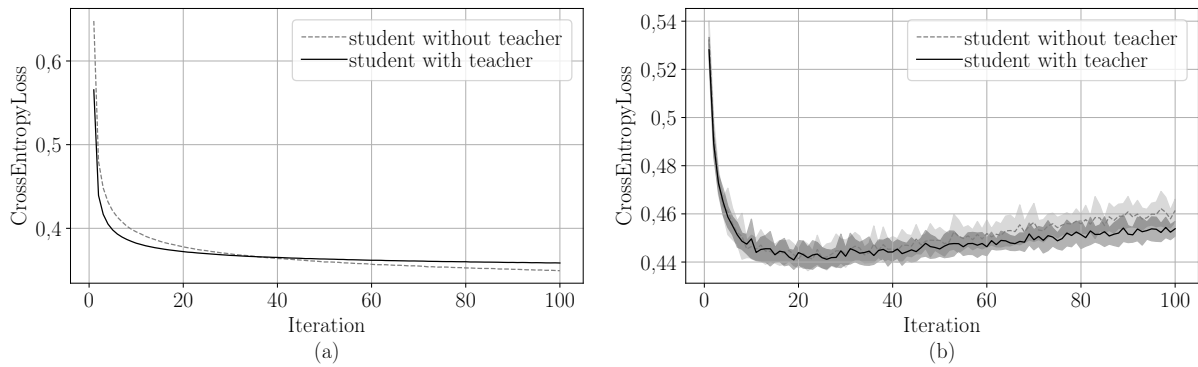


Рис. 2. Зависимость кросс-энтропии между истинными метками и предсказанными учеников вероятностями классов: а) на обучающей выборке; б) на тестовой выборке.

На рис. 2 показан график зависимости кросс-энтропии между истинными метками объектов и вероятностями, которые предсказывает модель ученика. На графике сравнивается модель, которая обучалась без учителя (в задаче оптимизации (17) присутствует только первое слагаемое) с моделью, которая была получена путем дистилляции модели нейросети в линейную модель. На графике видно, что обе модели начинают переобучаться после 30-й итерации, но модель, которая получена путем дистилляции переобучается не так быстро, что следует из того, что ошибка на тестовой выборке растет медленней, а на обучающей выборке падает также медленней.

## 6.2. Синтетический эксперимент

Проанализируем модель на синтетической выборке. Выборка построенная следующим образом:

$$\begin{aligned}\mathbf{W} &= [\mathcal{N}(w_{jk}|0, 1)]_{n \times K}, & \mathbf{X} &= [\mathcal{N}(x_{ij}|0, 1)]_{m \times n}, \\ \mathbf{S} &= \text{softmax}(\mathbf{XW}), & \mathbf{y} &= [\text{Cat}(y_i|\mathbf{s}_i)],\end{aligned}$$

где функция softmax берется построчно. Строки матрицы  $\mathbf{S}$  будем рассматривать как предсказание учителя, то есть учитель знает истинные вероятности каждого класса. На рис. 3 показана вероятностная модель в графовой нотации. В эксперименте число признаков  $n = 10$ , число классов  $K = 3$ , для обучения было сгенерировано  $m_{\text{train}} = 1000$  и  $m_{\text{test}} = 100$  объектов.

На рис. 4 показано распределение по классам для каждого объекта обучающей выборки. Видно, что все классы являются равновероятными.

Построим в качестве ученика простую линейную модель, которая минимизирует крос-энтропийную (первое слагаемое в формуле (17)). Представление данной модели в виде графовой модели показано на рис. 3.

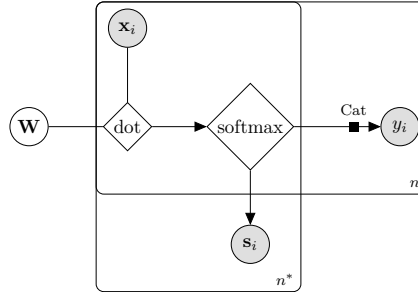


Рис. 3. Вероятностная модель используемая в синтетическом эксперименте.

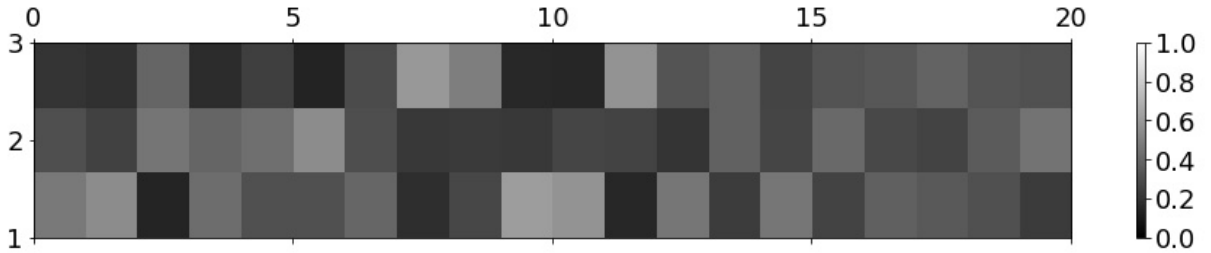


Рис. 4. Истинное распределение объектов по классам.

На рис. 5 показано распределение вероятностей классов, которое предсказала модель. Видно, что данное распределение является не соответствующим истинному, так как модель сосредотачивает всю вероятность в одном классе.

Рассмотрим модель, которая учитывает информацию о истинных распределениях на классах для каждого объекта. Для этого будем минимизировать первые три слагаемых в формуле (17), при  $T_0 = 1$  и  $\lambda = 0,75$ . В качестве меток учителя  $s_{i,k}$  использовались истинные вероятности для каждого класса для данного объекта. На

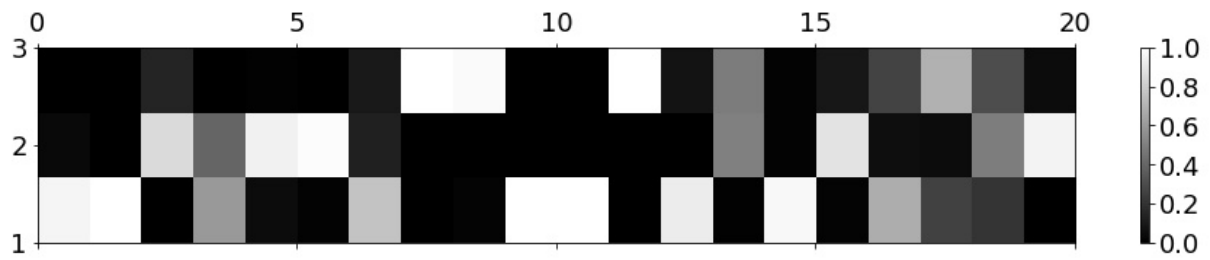


Рис. 5. Распределение предсказанное моделью без использования информации об истинном распределении на классах.

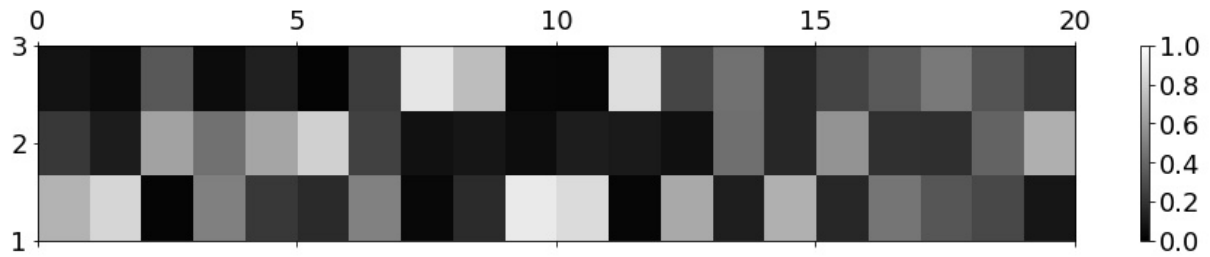


Рис. 6. Распределение предсказанное моделью с использованием информации об истинном распределении на классах.

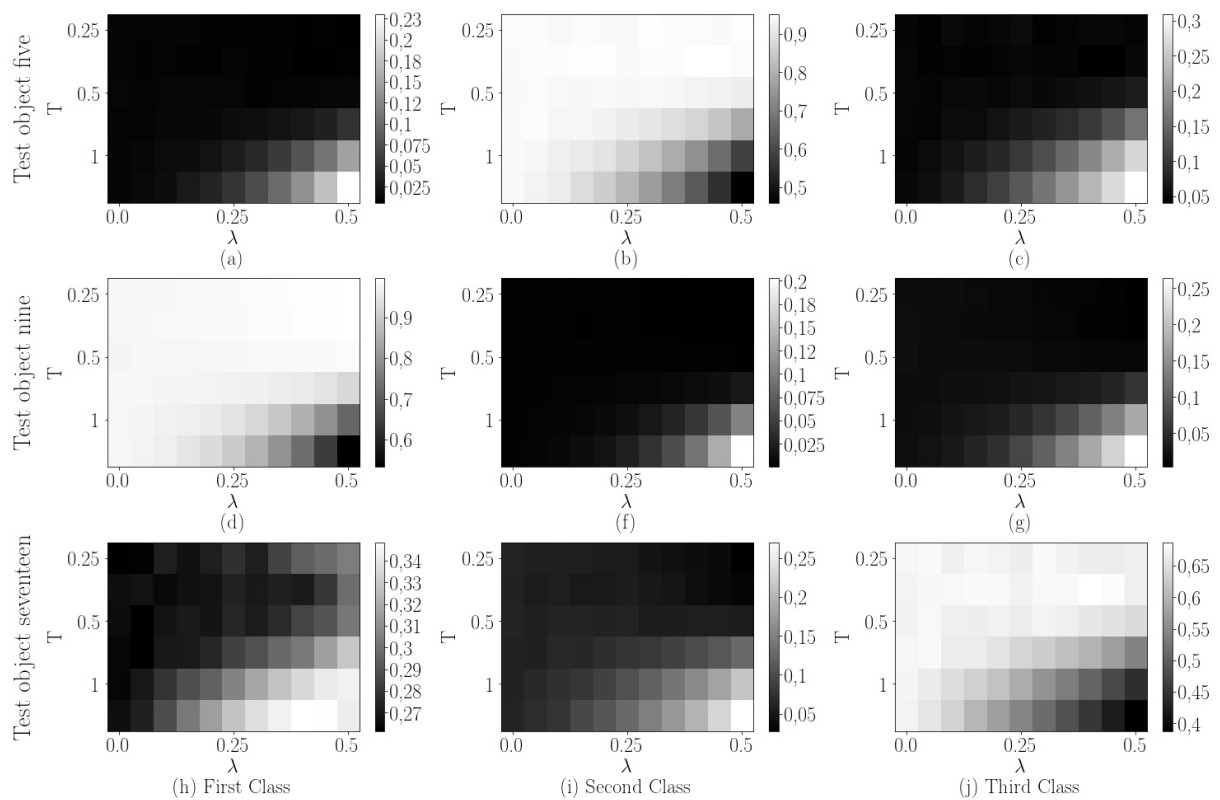


Рис. 7. Вероятности классов для разных объектов.

рис 6 показано распределение, которое дала модель в данном случае, видно, что распределения являются сглаженными и концентрации всей вероятности в одном классе

Таблица 1. Сводная таблица результатов 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

не наблюдается.

Заметим, что в данном примере предполагается, что модель учителя учитывает не только метки классов, а и распределение на метках классов, в то время как в выборке  $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ , имеются только точечные оценки в виде меткок.

В данном примере используются истинные распределения в качестве предсказаний учителя, но их можно заменить предсказаниями модели учителя, которая предсказывает не только сами меток, а и их распределение для каждого объекта.

На рис. 7 показана зависимость вероятности верного класса от температуры  $T$  и параметра доверия  $\lambda$  для одного из объекта из тестовой выборке. Видно, что при увеличении температуры распределение на классас становится более равномерным.

### 6.3. Выборка Twitter Sentiment Analysis

В данной части проводится эксперимент на выборке Twitter Sentiment Analysis. Данная выборка содержит короткие сообщения, для которых нужно предсказать эмоциональный окрас: содержит твит позитивный окрас или негативный. Выборка разделена на 1,18 миллиона твитов для обучения и 0,35 миллиона твитов для тестирования. В твитах была выполнена следующая предобработка:

- все твиты были переведены в нижний регистр;
- все никнеймы вида “@andrey” были заменены на токен “name”;
- все цифры были заменены на токен “number”.

Результаты данной части эксперимента показаны в табл. 1. В качестве модели учителя использовалась модель Bi-LSTM с 170 тысячами параметров для обучения. В качестве эмбедингов обучалась матрица из 30 миллионов параметров в единой процедуре с моделью BI-LSTM. Обученная модель предсказывает с точностью 0,835. В качестве модели ученика рассматривается модель с 1538 параметрами, но в качестве эмбедингов рассматривается переобученная модель BERT.

Программное обеспечение для проведения экспериментов и проверки результатов находится в [15].

## 7. Заключение

В данной работе проанализирована задача обучения модели ученика с помощью модели учителя. Исследован метод дистилляции и привилигированного обучения.

Предложено вероятностное обоснование дистилляции. Введены вероятностные предположения описывающие дистилляцию моделей. В рамках данных вероятностных предположений проанализированы модели для задачи классификации и регрессии. Результат анализа сформулирован в виде теоремы 1 и теоремы 2.

Теорема 2 показала, что обучения линейной регрессии с учителем эквивалентно замене обучающей выборке и вероятностных предположений о распределении истинных ответов. Для задачи классификации ответы учителя дают дополнительную информацию в виде распределения классов для каждого объекта из обучающей выборки. Данная информация не может быть переписана в виде классической задачи классификации. Для использования данной информации требуется использовать распределение, которое представлено в теореме 1.

В вычислительном эксперименте сравнивается модель ученика, которая обучена без использования учителя и с использованием модели учителя. В таблице 1 показаны результаты вычислительного эксперимента для разных выборок. Из таблицы видно, что точность аппроксимации выборки учеником улучшается при использовании модели учителя. Задачи регрессии не приведена в вычислительном эксперименте, так как в теореме 2 была показана эквивалентность классическому решению задачи линейной регрессии. Для задачи классификации проведен вычислительный эксперимент.

В дальнейшем предполагается обобщить метод максимального правдоподобия для дистилляции моделей используя Байесовский подход выбора моделей машинного обучения. Также в рамках байесовского подхода планируется улучшить методы для получения улучшения качества не только для задачи классификации, но и для задачи регрессии.

## СПИСОК ЛИТЕРАТУРЫ

1. *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.
2. *Devlin J., Chang M., Lee K., Toutanova K.* BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding // arXiv preprint arXiv:1810.04805. 2018.
3. *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.
4. *Бахтеев О. Ю., Стрижов В. В.* Выбор моделей глубокого обучения субоптимальной сложности // АиТ. 2018. № 8. С. 129–147.
5. *Hinton G., Vinyals O., Dean J.* Distilling the Knowledge in a Neural Network // NIPS Deep Learning and Representation Learning Workshop. 2015.
6. *LeCun Y., Cortes C., Burges C.* The MNIST dataset of handwritten digits, 1998. <http://yann.lecun.com/exdb/mnist/index.html>.

7. *Vapnik V., Izmailov R.* Learning Using Privileged Information: Similarity Control and Knowledge Transfer // Journal of Machine Learning Research. 2015. No 16. P. 2023–2049.
8. *Lopez-Paz D., Bottou L., Scholkopf B., Vapnik V.* Unifying Distillation and Privileged Information // In International Conference on Learning Representations. Puerto Rico, 2016.
9. *Madala H., Ivakhnenko A.* Inductive Learning Algorithms for Complex Systems Modeling. Boca Raton: CRC Press Inc., 1994.
10. *Xiao H., Rasul K., Vollgraf R.* Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms // arXiv preprint arXiv:1708.07747. 2017.
11. *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.
12. *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.
13. *Hochreiter S., Schmidhuber J.* Long short-term memory // Neural Computation. 1997. V. 9. No 8. P. 1735–1780.
14. *Kingma D, Ba J.* Adam: A Method for Stochastic Optimization // arXiv preprint arXiv:1412.6980. 2014.
15. Код вычислительного эксперимента. URL: <https://github.com/andriygav/PrivilegeLearning>

Грабовой А.В., Московский физико-технический институт, студент,  
Москва, [grabovoy.av@phystech.edu](mailto:grabovoy.av@phystech.edu)

Стрижов В.В., Московский физико-технический институт, профессор,  
Москва, [strijov@phystech.edu](mailto:strijov@phystech.edu)