

# Байесовская дистилляция моделей на базе трансформеров

Игорь Николаевич Игнашин

Московский физико-технический институт

*Научный руководитель:* к. ф.-м. н. А. В. Грабовой

18 мая 2024 г.

# Слайд об исследованиях

Исследуется проблема снижения размерности пространства параметров аппроксимирующих моделей.

## Цель исследования:

Адаптация методов построения выравнивающих преобразований структуры модели учителя в модель ученика для моделей трансформеров.

## Решение:

Построение последовательности выравнивающих преобразований позволяющих выровнять структуры модели учителя в модель ученика на базе архитектуры трансформера.

# Дистилляция Дж. Хинтона<sup>1</sup>

Заданы

- 1) признаки  $\mathbf{x}_i \in \mathbb{R}^n$ ,
- 2)  $y_i \in \mathbb{Y} = \{1, \dots, K\}$ ,  $\mathbb{Y}' = \mathbb{R}^K$ .

Параметрические семейства учителя и ученика:

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

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

где  $\mathbf{z}, \mathbf{v}$  — дифференцируемые по параметрам функции заданной структуры,  $T$  — параметр температуры. Оптимальная модель учителя  $\hat{\mathbf{f}} \in \mathfrak{F}_{\text{cl}}$ .

Функция ошибки

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

где  $\cdot|_{T=t}$  фиксирует температуру  $T$ .

Оптимальная модель выбирается из класса,  $\hat{\mathbf{g}} = \arg \min_{\mathbf{g} \in \mathfrak{G}_{\text{cl}}} \mathcal{L}(\mathbf{g})$ .

---

<sup>1</sup>Hinton G., et al [Distilling the knowledge in a neural network](#) // NIPS, 2015.

# Байесовская постановка задачи дистилляции

Задана обучающая выборка  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$   $x_i \in \mathbb{R}^n$ ,  $y_i \in \mathbb{Y}$ .

Модель учителя

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

где  $\mathbf{U}$  матрицы линейных отображений,  $\sigma$  монотонная вектор-функция.

Параметры учителя фиксированы

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

На основе выборки  $\{x_i, y_i\}_{i=1}^m$  и значений учителя  $\mathbf{f}(\hat{\mathbf{u}}, \mathbf{x})$  требуется выбрать модель ученика:

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

где  $\mathbf{W}$ ,  $\sigma$  вводятся как и отображения учителя. Задача выбора модели  $\mathbf{g}$  состоит в оптимизации вектора  $\mathbf{w}$ . Решается вариационным выводом

$$\hat{\mathbf{w}}, \hat{\mu}, \hat{\Sigma} = \arg \min_{\mu, \Sigma, \mathbf{w}} D_{\text{KL}}(q(\mathbf{w}|\mu, \Sigma) || p(\mathbf{w}|\mathbf{A})) - \mathbb{E}_{\mathbf{w} \sim q} \sum_{i=1}^m \log p(y_i | \mathbf{x}_i, \mathbf{w}).$$

Априорное распределение  $p(\mathbf{w}|\mathbf{A})$  задается как функция от апостериорного распределения параметров учителя  $p(\mathbf{u}|\mathbf{X}, \mathbf{y})$ . Оно задано

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

Так как размерности  $\mathbf{u}$  и  $\mathbf{w}$  не совпадают, то применяется выравнивающее преобразование – приведение параметров моделей в одно общее пространство:

$$p(\mathbf{w}|\mathbf{A}) = p(v|\mathbf{X}, \mathbf{y}), \quad v = \psi(t, \mathbf{u}), \quad \psi(t) : \mathbb{R}^{\text{Ptr}} \rightarrow \mathbb{R}^{\text{Ptr} - n_t n_{t-1}}.$$

---

<sup>1</sup>Грабовой А.В. Априорное распределение параметров в задачах выбора моделей глубокого обучения, 2022.

# Отличие архитектуры трансформера от полносвязной сети

Задана выборка  $\mathcal{D} = \{x_i, y_i\}_{i=1}^n$ , где  $x_i = [x_i^1, x_i^2, \dots, x_i^{k_i}]$ ,  $x_i^k \in \overline{0, 1, \dots, v_o}$ ,  
 $y_j = [y_j^1, y_j^2, \dots, y_j^{l_j}]$ ,  $y_j^k \in \overline{0, 1, \dots, v_d}$ .

Модель учителя

$$f(x, y) = G \circ D_T \circ \dots \circ D_2 \circ D_1 \circ E[x, y].$$

Модель ученика

$$g(x, y) = G \circ D_T \circ \dots \circ D_{k+1} \circ D_{k-1} \circ \dots \circ D_2 \circ D_1 \circ E[x, y].$$

Энкодер

$$E : \underbrace{\mathbb{R}^{L \times d_{model}}}_{\text{токены } x} \times \underbrace{\mathbb{R}^{M \times v_{eng}}}_{\text{токены } y} \rightarrow \underbrace{\mathbb{R}^{L \times d_{model}}}_{\text{память}} \times \underbrace{\mathbb{R}^{M \times d_{model}}}_{\text{токен эмбединги } y}.$$

Слой декодера

$$D_k : \underbrace{\mathbb{R}^{L \times d_{model}}}_{\text{память}} \times \underbrace{\mathbb{R}^{M \times d_{model}}}_{\text{эмбединги декодера}} \rightarrow \underbrace{\mathbb{R}^{L \times d_{model}}}_{\text{память}} \times \underbrace{\mathbb{R}^{M \times d_{model}}}_{\text{новые эмбединги декодера}}.$$

Генератор

$$G : \underbrace{\mathbb{R}^{L \times d_{model}}}_{\text{память}} \times \underbrace{\mathbb{R}^{M \times d_{model}}}_{\text{эмбединги декодера}} \rightarrow \underbrace{\mathbb{R}^{M \times d_{model} \times v_d}}_{\text{логиты}}.$$

# Байесовская дистилляция трансформера

Задана модель учителя, суперпозиция

$$\log p(\mathbf{y}_{1:t+1}) = \mathbf{f}(\mathbf{x}, \mathbf{y}_{0:t}) = \mathbf{G}(\mathbf{U}_{T+1}) \circ \mathbf{D}_T(\mathbf{U}_T) \circ \mathbf{D}_{T-1}(\mathbf{U}_{T-1}) \circ \cdots \circ \mathbf{D}_1(\mathbf{U}_1) \circ \mathbf{E}(\mathbf{U}_0)[\mathbf{x}, \mathbf{y}_{0:t}],$$

где  $\mathbf{D}$  слои декодера,  $\mathbf{E}$  энкодер,  $\mathbf{G}$  генератор, возвращающий логиты следующих токенов. Параметры учителя фиксированы

$$\mathbf{u} = \text{vec}([\mathbf{U}_{T+1}, \mathbf{U}_T, \mathbf{U}_{T-1}, \cdots, \mathbf{U}_1, \mathbf{U}_0]),$$

где  $\mathbf{U}_k$  векторизированные параметры соответствующих модулей.

На основе выборки  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^m$  и значений учителя  $\mathbf{f}(\hat{\mathbf{u}}, \mathbf{x}, \mathbf{y})$  требуется выбрать модель ученика:

$$\mathbf{g}(\mathbf{x}, \mathbf{y}_{0:t}) = \mathbf{G}(\mathbf{W}_{L+1}) \circ \mathbf{D}_L(\mathbf{W}_L) \circ \cdots \circ \mathbf{D}_1(\mathbf{W}_1) \circ \mathbf{E}(\mathbf{W}_0)[\mathbf{x}, \mathbf{y}_{0:t}], \quad L \leq T,$$

где  $\mathbf{G}, \mathbf{D}, \mathbf{E}$  вводятся как и отображения учителя. Задача выбора модели  $\mathbf{g}$  состоит в оптимизации вектора  $\mathbf{w} = \text{vec}([\mathbf{W}_{L+1}, \mathbf{W}_L, \mathbf{W}_{L-1}, \cdots, \mathbf{W}_1, \mathbf{W}_0])$ . Решается вариационным выводом

$$\hat{\mathbf{w}}, \hat{\mu}, \hat{\Sigma} = \arg \min_{\mu, \Sigma, \mathbf{w}} D_{\text{KL}}(q(\mathbf{w}|\mu, \Sigma) || p(\mathbf{w}|\mathbf{A})) - \mathbb{E}_{\mathbf{w} \sim q} \sum_{i=1}^m \log p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w}).$$

Априорное распределение  $p(\mathbf{w}|\mathbf{A})$  задается как функция от апостериорного распределения параметров учителя  $p(\mathbf{u}|\mathbf{X}, \mathbf{Y})$ . Оно задано:

$$p(\mathbf{u}|\mathbf{X}, \mathbf{Y}) = \mathcal{N}(\mathbf{m}, \Sigma).$$

Так как размерности  $\mathbf{u}$  и  $\mathbf{w}$  не совпадают, то применяется выравнивающее преобразование – приведение параметров моделей в одно общее пространство:

$$p(\mathbf{w}|\mathbf{A}) = p(v|\mathbf{X}, \mathbf{y}), \quad v = \psi(\mathbf{u}), \quad \psi: \mathbb{R}^{\text{Ptr}} \rightarrow \mathbb{R}^{\text{Pst}}.$$

# Выравнивание модели трансформера

Построим выравнивающее преобразование  $\psi$  из пространства параметров модели учителя в пространство параметров модели ученика:

$$\psi : \mathbb{R}^{\text{Ptr}} \rightarrow \mathbb{R}^{\text{Ptr} - 4 \cdot d_{\text{model}} \times d_k - 2d_v \times d_{\text{model}} - 2 \cdot N_{\text{heads}} \cdot d_v \times d_{\text{model}} - 2d_{\text{model}} \times d_{\text{ff}} - d_{\text{ff}} - d_{\text{model}}}.$$

Модель учителя:

$$\mathbf{f}(\mathbf{x}, \mathbf{y}) = \mathbf{G}(\mathbf{U}_{T+1}) \circ \mathbf{D}_T(\mathbf{U}_T) \circ \cdots \circ \mathbf{D}_2(\mathbf{U}_2) \circ \mathbf{D}_1(\mathbf{U}_1) \circ \mathbf{E}(\mathbf{U}_0)[\mathbf{x}, \mathbf{y}]$$

$$\mathbf{D}_k(\mathbf{x}, \mathbf{y}) \equiv [\mathbf{x}, \mathbf{D}_k^0(\mathbf{x}, \mathbf{y})]$$

$$\mathbf{D}_k^0(\mathbf{x}, \mathbf{y}) = \underbrace{(\mathbf{I} + \mathbf{FFN})}_{\text{feed-forward}} \circ \underbrace{(\mathbf{I} + \mathbf{MHA})[\mathbf{x}]}_{\text{multi-head attention}} \circ \underbrace{(\mathbf{I} + \mathbf{SA})}_{\text{self attention}}(\mathbf{y})$$

$$\mathbf{FFN}(\mathbf{y}) = (\mathbf{b}_2 + \sigma \circ \mathbf{U}_2^F \circ (\mathbf{b}_1 + \mathbf{y} \cdot \mathbf{U}_1^F))$$

$$\mathbf{MHA}(\mathbf{x}, \mathbf{y}) = \text{heads}[\mathbf{U}_M^h](\mathbf{x}, \mathbf{y}) \cdot \mathbf{U}^M \quad \mathbf{SA}(\mathbf{y}) = \text{heads}[\mathbf{U}_S^h](\mathbf{y}) \cdot \mathbf{U}^S,$$

где  $\sigma$  функция активации, параметры модуля  $\mathbf{U}_k$ :

$$\mathbf{U}_1^F \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}, \mathbf{U}_2^F \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}, \mathbf{b}_1 \in \mathbb{R}^{d_{\text{ff}}}, \mathbf{b}_2 \in \mathbb{R}^{d_{\text{model}}},$$

$$\mathbf{U}^M, \mathbf{U}^S \in \mathbb{R}^{N_{\text{heads}} \cdot d_v \times d_{\text{model}}}, \mathbf{U}_M^h, \mathbf{U}_S^h \in \mathbb{R}^{2 \cdot d_k \times d_{\text{model}} + d_v \times d_{\text{model}}}.$$

Модель ученика:

$$\mathbf{g}(\mathbf{x}, \mathbf{y}) = \mathbf{G}(\mathbf{W}_{T+1}) \circ \cdots \circ \mathbf{D}_{k+1}(\mathbf{W}_{k+1}) \circ \mathbf{D}_{k-1}(\mathbf{W}_{k-1}) \circ \cdots \circ \mathbf{D}_1(\mathbf{W}_1) \circ \mathbf{E}(\mathbf{W}_0)[\mathbf{x}, \mathbf{y}].$$

Эквивалентность моделей:

$$\mathbf{f} \mid_{\mathbf{U}_2^F, \mathbf{b}_2, \mathbf{U}^M, \mathbf{U}^S=0} \equiv \mathbf{g}.$$

Параметры  $\mathbf{u}$  модели  $\mathbf{f}$  делятся на **удаляемые**  $\xi_2 = \text{vec}([\mathbf{b}_1, \mathbf{U}_1^F, \mathbf{U}_M^h, \mathbf{U}_S^h])$ ,

**зануляемые**  $\xi_1 = \text{vec}([\mathbf{b}_2, \mathbf{U}_2^F, \mathbf{U}^S, \mathbf{U}^M])$ , оставшиеся  $v = \text{vec}([\mathbf{U}_{T+1}, \cdots, \mathbf{U}_{k+1}, \mathbf{U}_{k-1}, \cdots, \mathbf{U}_0])$ .

# Решение задачи выравнивания структур моделей

Параметры  $\mathbf{u}$  модели  $\mathbf{f}$ :

1. удаляемые  $\xi_2 = \text{vec}([b_1, U_1^F, U_M^h, U_S^h])$ ,
2. зануляемые  $\xi_1 = \text{vec}([b_2, U_2^F, U^S, U^M])$ ,
3. оставшиеся  $v = \text{vec}([U_{T+1}, \dots, U_{k+1}, U_{k-1}, \dots, U_0])$ .

Апостериорное распределение параметров  $v = \psi_{\text{transformer}}(\mathbf{u}, k)$  модели  $\mathbf{f}$ :

$$p(v|\mathcal{D}) = \int_{\xi_2} p(\bar{\xi}_1|\mathcal{D}, \xi_1 = \mathbf{0}) d\xi_2,$$

где  $\bar{\xi}_1 = [v, \xi_2]$ .

Из свойства распределения  $p(\bar{\xi}_1|\mathcal{D}, \xi_1 = \mathbf{0}) = \mathcal{N}(\mu, \Xi)$ , с параметрами  $\mu, \Xi$ :

$$\begin{aligned}\mu &= \mathbf{m}_{\bar{\xi}_1} + \Sigma_{\bar{\xi}_1, \xi_1} \Sigma_{\xi_1, \xi_1}^{-1} (\mathbf{0} - \mathbf{m}_{\xi_1}), \\ \Xi &= \Sigma_{\bar{\xi}_1, \bar{\xi}_1} - \Sigma_{\bar{\xi}_1, \xi_1} \Sigma_{\xi_1, \xi_1}^{-1} \Sigma_{\xi_1, \bar{\xi}_1}.\end{aligned}$$

Маргинализация нормального распределения  $p(v|\mathcal{D}) = \mathcal{N}(\mu_v, \Xi_{v,v})$ .

## Theorem (Игнашин, 2024)

Пусть апостериорное распределение параметров модели учителя  $\mathbf{u}$  имеют распределение  $\mathcal{N}(\mathbf{m}, \Sigma)$ . Модель ученика имеет схожую структуру с моделью учителя, но без одного слоя декодера  $\mathbf{D}_k(\mathbf{U}_k)$ . Тогда апостериорное распределение имеет вид:

$$p(\psi_{\text{transformer}}(\mathbf{u}, k)|\mathcal{D}) = \mathcal{N}(\mathbf{m}_v + \Sigma_{v, \xi_1} \Sigma_{\xi_1, \xi_1}^{-1} (\mathbf{0} - \mathbf{m}_{\xi_1}), \Sigma_{v,v} - \Sigma_{v, \xi_1} \Sigma_{\xi_1, \xi_1}^{-1} \Sigma_{\xi_1, v}),$$



# Случай некоррелированных параметров

Рассмотрим частный случай, когда ковариационные матрицы нулевые  $\Sigma_{\bar{\xi}_1, \xi_1} = 0$ , то есть **зануляемые параметры  $\xi_1$**  и остальные параметры  $\bar{\xi}_1$  – некоррелированы. Тогда распределение  $p(\bar{\xi}_1 | \mathcal{D}, \xi_1 = \mathbf{0}) = \mathcal{N}(\mu, \Xi)$  имеет параметры  $\mu, \Xi$ :

$$\begin{aligned}\mu &= \mathbf{m}_{\bar{\xi}_1}, \\ \Xi &= \Sigma_{\bar{\xi}_1, \bar{\xi}_1}.\end{aligned}$$

Следовательно апостериорное распределение параметров модели учителя  $v$ :

$$p(v | \mathcal{D}) = \mathcal{N}(m_v, \Sigma_{v,v}).$$

Априорное распределение модели ученика :

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

# Описание эксперимента

Обучающий датасет: 40 тысяч руско-английских переводов предложений.

В качестве аппроксимации вариационного вывода применяется обучение кроссэнтропийного лосса.

Инициализация модели ученика:

$$\mathbf{w}_0 = \mathbf{m}_v.$$

Аппроксимация  $\mathbf{m}$  апостериорного распределения параметров модели учителя

$$p(\mathbf{u}|\mathbf{X}, \mathbf{Y}) = \mathcal{N}(\mathbf{m}, \Sigma):$$

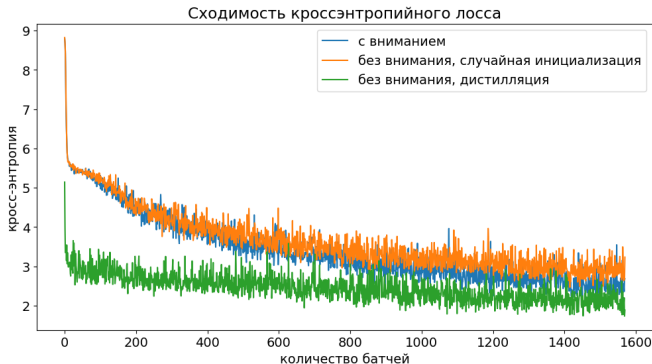
$$\mathbf{m} = \mathbf{u}_K,$$

где  $\mathbf{u}_K$  – параметры модели учителя на  $K$ -ой итерации обучения.

Выбранные рекуррентные модели:

1. Модель учителя,
2. Модель ученика  $\mathbf{w}_0 \sim \text{Uniform}[-0.08, 0.08]$ ,
3. Модель ученика  $\mathbf{w}_0 = \mathbf{m}_v$ .

# Результат вычислительного эксперимента



	модель учителя	модель ученика	модель ученика (дистиллированная)
BLEU	23.70	13.35	17.73

Дистиллированная модель сходится лучше, а также показывает лучшую метрику BLEU на отложенной выборке, чем модель той же структуры, но с произвольной инициализацией.

# Выносятся на защиту

Получены следующие результаты :

1. Предложен метод удаления слоя декодера трансформера.
2. Предложены методы выравнивания структур моделей глубокого обучения с механизмом внимания.
3. Проведен теоретический анализ предложенных методов
4. Доказана теорема о выравнивании модели трансформера.
5. Проведен вычислительный эксперимент, показывающий состоятельность метода.

Публикации и выступления на конференциях:

1. О задаче поиска равновесного распределения потоков. «66-я Всероссийская научная конференция МФТИ», 2024.
2. Игнашин И.Н., Ярмошик Д.В. Модификации алгоритма Frank –Wolfe в задаче поиска равновесного распределения транспортных потоков. // Компьютерные исследования и моделирование, 2024 Т. 16 № 1 С. 53–68.