

# Методы детекции машинно-сгенерированных фрагментов в документах на базе анализа смены стиля

Анастасия Евгеньевна Вознюк  
Научный руководитель: к.ф.-м.н. А. В. Грабовой

Кафедра интеллектуальных систем ФПМИ МФТИ  
Специализация: Интеллектуальный анализ данных  
Направление: 01.03.02 Прикладные математика и информатика

27 июня 2024

# Цель детекции машинно-сгенерированных фрагментов

## Проблема

В документах все чаще и чаще встречаются фрагменты, написанные языковыми моделями, и необходимо уметь находить такие фрагменты.



Детекция  
фрагментов

## Цель

Предложить методы детекции фрагментов в документе в случае:

1. смены авторов по известным позициям;
2. единственной смены авторов на произвольной позиции.



Детекция смены  
стиля

## Решение

Использовать другие фрагменты в качестве контекста и искать смену стилистики.

## Общая постановка нашей задачи

Определим документ как конечную последовательность символов из заданного алфавита  $\mathcal{W}$ . Пространство документов:

$$\mathbb{D} = \left\{ \left[ t_j \right]_{j=1}^n \mid t_j \in \mathcal{W}, n \in \mathbb{N} \right\}.$$

Дан набор из  $N$  документов

$$\mathcal{D} = \bigcup_{i=1}^N \mathcal{D}^i, \mathcal{D}^i \in \mathbb{D}.$$

Определим множество авторов, тексты которых встречаются в наборе  $\mathcal{D}$ :

$$\mathcal{C} = \{0, \dots, k-1\}.$$

# Детекция фрагментов в документе

Для каждого документа  $d \in \mathbb{D}$  существует разбиение на фрагменты различного авторства:

$$\mathbb{T} = \left\{ \left[ t_{s_j}, t_{f_j}, C_j \right]_{j=1}^J \mid t_{s_j} = t_{f_{j-1}}, \quad s_j \in \mathbb{N}_0, \quad f_j \in \mathbb{N}, \quad C_j \in \mathcal{C} \right\},$$

где  $J$  — количество фрагментов разного авторства в документе,  $t_{s_j}$  и  $t_{f_j}$  — начало и конец  $j$ -ого фрагмента одного автора,  $C_j$  — автор  $j$ -ого фрагмента.

## Модель детектора

$$\phi : \mathbb{D} \rightarrow \mathbb{T} \quad \phi : g \circ f,$$

где  $f$  — модель выделения фрагментов, а  $g$  — классификатор авторов.



## Определение авторов параграфов

Пусть для документа  $d \in \mathbb{D}$  известны его параграфы  $\mathcal{P} = (p_1, \dots, p_n)$  и  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  — векторные представления параграфов.

### Вероятностная модель для определения авторов

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = p(\mathbf{y} | \mathbf{x}) = \frac{\exp \Phi(\mathbf{x}, \mathbf{y})}{\sum_{\mathbf{y}' \in \mathcal{Y}^n} \exp(\Phi(\mathbf{x}, \mathbf{y}'))},$$

где  $y_i \in \mathcal{C}$  — метка автора для параграфа  $p_i$ ,

$\mathcal{Y}^n$  — все возможные последовательности меток длины  $n$ ,

а функция  $\Phi(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^d$  определена как :

$$\Phi(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \left( \log \phi_{\text{EMIT}}(y_i \rightarrow x_i) + \log \phi_{\text{TRANS}}(y_{i-1} \rightarrow y_i) \right),$$

где  $\phi_{\text{EMIT}}$  определяет вероятность появления метки  $y_i$  для  $\mathbf{x}_i$ ,

а  $\phi_{\text{TRANS}}$  определяет вероятность появления метки  $y_i$  и  $y_{i-1}$  в качестве соседей,

## Задача оптимизации для определения авторов параграфов

Функция потерь для выборки документов  $\mathcal{D}$  и последовательностей меток  $\hat{\mathcal{Y}}$ :

$$\begin{aligned}\mathcal{L}(\hat{\mathcal{Y}}, \mathcal{D}) &= - \sum_{i=1}^{|\mathcal{D}|} \log(p(\mathbf{y}^i | \mathbf{x}^i)) = - \sum_{i=1}^{|\mathcal{D}|} \log \left[ \frac{\exp \Phi(\mathbf{x}^i, \mathbf{y}^i)}{\sum_{\mathbf{y}' \in \mathcal{Y}^n} \exp(\Phi(\mathbf{x}^i, \mathbf{y}'))} \right] = \\ &= \sum_{i=1}^{|\mathcal{D}|} \underbrace{\left( \log \left[ \sum_{\mathbf{y}' \in \mathcal{Y}^n} \exp(\Phi(\mathbf{x}^i, \mathbf{y}')) \right] - \Phi(\mathbf{x}^i, \mathbf{y}^i) \right)}_{\pi_n},\end{aligned}$$

$\mathcal{Y}^n$  — все возможные последовательности меток длины  $n$ ,

# Быстрый пересчёт функции потерь

$$\pi_i(\mathbf{x}) := \log \sum_{\mathbf{y} \in \mathcal{Y}^i} \exp(\Phi(\mathbf{x}, \mathbf{y})),$$

$$\pi_i^j(\mathbf{x}) := \log \sum_{\substack{\mathbf{y} \in \mathcal{Y}^i \\ y^i = j}} \exp(\Phi(\mathbf{x}, \mathbf{y})),$$

где  $j \in \mathcal{C}$ ,  $1 \leq i \leq |\mathbf{x}|$

## Утверждение 1

$\pi_i^j(\mathbf{x})$  выражается через  $\pi_{i-1}^j(\mathbf{x})$

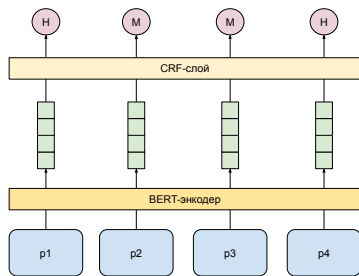
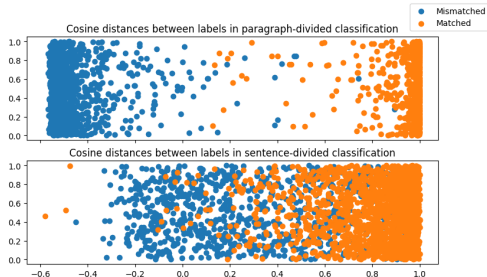


Схема модели для определения авторов параграфов

# Эксперименты с определением авторов фрагментов

Для детекции фрагментов по параграфам был сгенерирован новый датасет на основе статей с Medium.com. В текстах из 4-6 параграфов некоторые параграфы заменяли на машинно-сгенерированные.



Разбиение документа	Метод	Точность
По предложениям	RoBERTa	0.74
	RoBERTa-CRF	0.75
По параграфам	RoBERTa	0.89
	RoBERTa-CRF	0.93



## Детекция смены стиля

Пусть для документа  $d \in \mathbb{D}$  известно, что

$$\exists l_d \in \mathbb{N}_0 \quad g([t_0, t_l)) = 0, \quad g([t_{l+1}, t_{|d|})) = 1$$

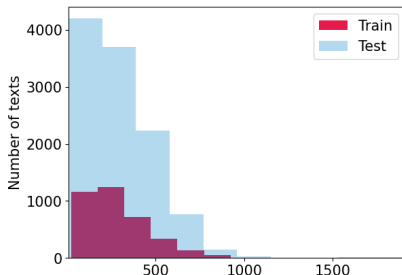
Тогда достаточно с помощью модели  $f$  найти позицию смены авторов. Назовем эту позицию сменой стиля. Параметризуем  $f$  длиной скользящего окна  $\ell$ . Функция  $f$  будет для каждого токена в документе оценивать вероятность его позиции быть позицией смены стиля:

### Модель смены стиля

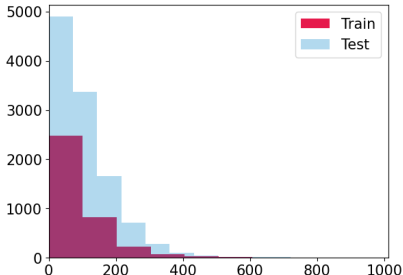
$$\begin{aligned} f_\ell &: \mathbb{D} \times \mathbb{N}_0 \rightarrow [0, 1], \\ f_\ell(d, i) &= p(t_i | t_{i-\ell}, \dots, t_{i-1}, t_{i+1}, t_{i+\ell}) \\ l_d &= \arg \max_{0 \leq i \leq |d|} f_\ell(d, i) \end{aligned}$$



# Эксперименты со сменой стиля



(a) Длины текстов



(b) Позиция смены автора

Распределения статистик в текстах после токенизации

Для детекции смены стиля данные были взяты из набора данных с соревнования **SemEval2024 Task 8 Subtask C**. Данные были дополнительно аугментированы для увеличения размера выборки и внесения разнообразия в позиции смены автора.

## Данные для экспериментов со сменой стиля

Для оценки качества предлагается использовать метрику MAE — Mean Absolute Error, среднее отклонение предсказанного индекса от истинного, считается на уровне слов в документе. По итогам участия в соревновании были получены следующие результаты:

Модель	Исходный датасет	Новый датасет
RoBERTa-base	31.56	30.71
RoBERTa-large	25.25	20.66
longformer-base	23.16	22.94
longformer-large	22.97	20.33
DeBERTaV3-base	16.12	13.98
DeBERTaV3-large	15.16	<b>13.38</b>
Top 1 соревнования	15.68	-

Метрика MAE на исходных и новых (аугментированных) данных. Дополнительно приведено лучшее решение с таблицы результатов соревнования. Longformer предлагался в качестве базового решения

## Выносятся на защиту

1. Модель детекции смены авторства на основе марковской линейной цепочки для документов, в которых смена авторов происходит только на уровне параграфов или по другим известным позициям.
2. Модель детекции смены авторства на основе трансформерных моделей для документов, в которых смена авторов происходит единожды на произвольной позиции в документе.

### Публикации

1. A. Voznyuk et al.. Leveraging Transfer Learning for Detecting Boundaries of Machine-Generated Texts // Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024) (на опубликовании).