

# Оптимизация критерия заданного нейросетевой моделью в задаче детоксификации текста

А. А. Пилькевич

Московский физико-технический институт

*Научный руководитель:* д.ф.-м.н. К. В. Воронцов, д.ф.-м.н. В. В. Стрижов

*Консультант:* А. С. Попов

2022

# Детоксификация предложений

## Задача

Стилизация токсичных (частично обценных) входных предложений к нейтральному варианту. Для оценки качества используется нейросетевая модель степени токсичности.

## Проблема

Модель невозможно использовать в качестве функции потерь из-за различия входных данных, так как теряется её дифференцируемость по параметрам.

## Решение

Предлагается «адаптировать» оценочную модель, чтобы она принимала распределения вероятностей токенов, выданных детоксификатором.

# Данные и оценка качества детоксификации

Дано множество пар  $(s_t, s_d)$ :

$s_t$ — токсичное предложение	$s_d$ — нейтральная версия
<i>«Её муженька козла на кол надо посадить.»</i> <i>«Это твари а не люди.»</i>	<i>«Её мужа нужно наказать.»</i> <i>«Это плохие люди.»</i>

**Требуется:** по токсичному предложению построить нейтральное.

**Качество стилизация текста оценивается:**

- ▶ точностью соответствия заданному стилю,
- ▶ качеством сохранения смысла после перефразы.

# Автоматическая оценка качества детоксификации

**Style Transfer Accuracy (STA)** — Conversational RuBERT<sup>1</sup>. Модель классификации, предсказывающая вероятность токсичности предложения. Отвечает за стилизацию текста.

**chrF** — F-score на основе символьных  $n$ -грамм<sup>2</sup>. Отвечает за сохранение смысла:

$$\text{chrF}_\beta = (1 + \beta^2) \frac{\text{chrP} \cdot \text{chrR}}{\beta^2 \text{chrP} + \text{chrR}},$$

- ▶ chrP — доля символьных  $n$ -грамм из предлагаемого предложения, которые имеются в оригинальном.
- ▶ chrR — доля символьных  $n$ -грамм из оригинального предложения, которые представлены в предлагаемом:

---

<sup>1</sup><https://huggingface.co/DeepPavlov/rubert-base-cased-conversational>

<sup>2</sup>Popović M. *chrF: character  $n$ -gram F-score for automatic MT evaluation* – 2015.

# Детоксификация как машинный перевод

Поставим задачу детоксификации как задачу машинного перевода (seq-to-seq) для пар предложений  $(t, d)$ :

- ▶  $t = \tau_f(s_t)$ ,  $d = \tau_f(s_d)$ ,
- ▶  $\tau_f$  — токенизатор, переводящий текст в последовательность BPE-токенов из словаря  $V_f$ .

Архитектура задаётся моделью кодировщик-декодировщик  $f_\theta$  на основе предобученного ruT5-base<sup>3</sup>.

$$\mathcal{L}_{\text{CE}} = - \sum_{i=1}^n \log f_\theta(d_i | d_{<i}, t) \longrightarrow \min_{\theta},$$

$f_\theta(* | d_{<i}, t) \in [0, 1]^{|V_f|}$  — распределение вероятностей.

---

<sup>3</sup>Raffel C. et al. *Exploring the limits of transfer learning with a unified text-to-text transformer* – 2019.

# Детоксификация как задача стилизации

Предлагается использовать STA-модель  $g_{\text{toxic}}$  в качестве функции потерь:

$$\mathcal{L}_{\text{TP}} = g_{\text{toxic}}(\tau_g(s_{\text{detox}})) \rightarrow \min_{\theta}.$$

- ▶  $\tau_g$  — токенизатор, переводящий текст в последовательность BPE-токенов из словаря  $V_g$ ,
- ▶  $s_{\text{detox}}$  — результат работы детоксификатора:

$$s_{\text{detox}} = \{\arg \max_{d_i} f_{\theta}(d_i | d_{<i}, t)\}_{i=1}^n.$$

$\mathcal{L}_{\text{TP}}$  не дифференцируема по параметрам детоксификатора  $\theta$ , в силу недифференцируемости функции  $\tau_g$  и  $\arg \max$ !

$$\text{ТЫ СЛИШКОМ ТОКСИЧНЫЙ} \longrightarrow \begin{cases} [\text{ТЫ}, \text{СЛИШКОМ}, \text{ТОКС}, \text{ИЧ}, \text{НЫЙ}] & \text{от } \tau_f \\ [\text{ТЫ}, \text{\#UNK}, \text{ТОКСИЧН}, \text{ЫЙ}] & \text{от } \tau_g. \end{cases}$$

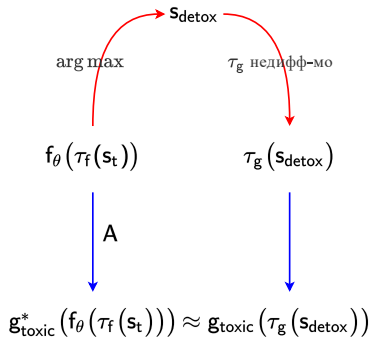
# Адаптер — аппроксимация векторного представления

Хотим при фиксированных параметрах STA-модели  $g_{\text{toxic}}$  выполнения:

$$g_{\text{toxic}}(\tau_g(s_d)) \approx g_{\text{toxic}}^*(f_\theta(\tau_f(s_t))).$$

$g_{\text{toxic}}^*$  — STA-модель, в которой заменили входной эмбединг слой *адаптером*  $A \in \mathbb{R}^{|V_f| \times e}$ , где  $e$  — размерность векторного представления токенов.

Причём  $g_{\text{toxic}}^*$  принимает  $f_\theta(*) \in [0, 1]^{n \times |V_f|}$  — «зашумлённые one-hot вектора».



# Обучение адаптера

**Выход STA-модели:**

$P_{\text{toxic}} = g_{\text{toxic}}(\tau_g(s_d))$  — вероятность токсичности предложения  $s_d$ ,

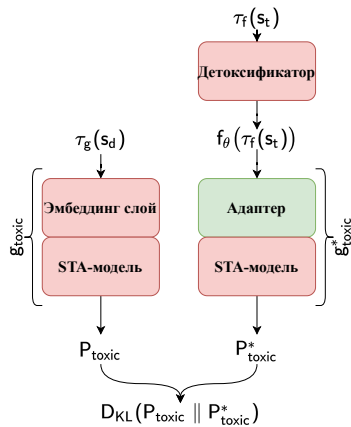
$P_{\text{toxic}}^* = g_{\text{toxic}}^*(f_\theta(\tau_f(s_t)))$  — вероятность токсичности с использованием адаптера для  $s_t$ .

**Функция потерь адаптера:**

$$D_{\text{KL}}(P_{\text{toxic}} \parallel P_{\text{toxic}}^*) \longrightarrow \min_A.$$

## Алгоритм обучения A

1. Обучается детоксификатор  $f_\theta$  на задаче seq-to-seq.
2. Обучается адаптер A при фиксированных параметрах  $f_\theta$  и STA-модели.





# Дообучение детоксификатора

При дообучении детоксификатора  $f_\theta$  используется схожая идея с обучением порождающих моделей<sup>4</sup>.

Пусть CE, TP — значение кросс-энтропии  $\mathcal{L}_{\text{CE}}$  и выход STA-модели  $\mathcal{L}_{\text{TP}}$ .  
 $F(*, *)$  — произвольная функция для агрегации функционалов,  $F : \mathbb{R}^2 \rightarrow \mathbb{R}$ .

## Алгоритм дообучения $f_\theta$

1.  $N$  батчей обучается  $f_\theta$  на  $F(\text{CE}, \text{TP})$ , при фиксированных параметрах  $A$ .
2.  $M$  батчей обучается  $A$  на  $D_{\text{KL}}$ , при фиксированных параметрах  $f_\theta$ .

---

<sup>4</sup>Goodfellow I. et al. *Generative adversarial nets* – 2014.

# Эксперименты по детоксификации

Обучающая выборка: 11136 пар  $(s_t, s_d)$ , 10% использовались для валидации во время обучения. Тестовой выборка: 800 предложений.

Результаты экспериментов в различных конфигурациях обучения:

Подход	F(CE, TP)	STA	chrF1	STA*chrF1
seq-to-seq	CE	0.739	0.578	0.427
GAN style	CE · TP	0.754	0.574	0.439
GAN style	CE + $w_2$ TP	0.813	0.569	0.462
seq-to-seq	TP	0.998	0.119	0.119

Наилучшее качество показывает линейное взвешивание функций потерь, что позволяет задать приоритет для оптимизируемого функционала.

# Проверенные подходы обучения

Проверка статистической значимости, когда на одном батче сперва оптимизировался детоксификатор, затем адаптер:

Подход	F(CE, TP)	STA	chrF1	STA*chrF1
seq-to-seq	CE	$0.744 \pm 0.010$	$0.575 \pm 0.002$	$0.430 \pm 0.042$
Same batch	$CE \cdot TP$	$0.776 \pm 0.015$	$0.569 \pm 0.004$	$0.442 \pm 0.006$
Same batch	$CE + w_2 TP$	$0.774 \pm 0.011$	$0.561 \pm 0.012$	$0.435 \pm 0.013$

Эксперимент с одновременным обучением детоксификатора и адаптера на одинаковую функцию потерь:

Подход	STA	chrF1	STA*chrF1
seq-to-seq CE	0.742	0.577	0.428
Adapter on embs	0.639	0.544	0.348
Adapter on logits	0.708	0.569	0.403

# Выносятся на защиту

1. Предложен алгоритм использования нейросетевой модели в качестве функции потерь в условиях отсутствия дифференцируемость, так как нет отображения между токенами различных токенайзеров.
2. Продемонстрирована работоспособность и эффективность предложенного метода.
3. Подобраны оптимальные параметры совместного обучения детоксификатора и адаптера.
4. Реализован и опубликован код для воспроизведения экспериментов из представленной работы<sup>5</sup>.

---

<sup>5</sup><https://github.com/Intelligent-Systems-Phystech/Pilkevich-BS-Thesis>