

Выявление манипуляций и их мишеней в новостных текстах.

Лукьяненко Иван

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем

Научный руководитель: д.ф.-м.н. Воронцов Константин Вячеславович

Лето 2023 г.

Propaganda Detection

Цель

Обобщить методологию задачи Propaganda Detection. Разработать базовые модели для задачи выявления манипуляции и их мишеней на русском языке.

Актуальность

Задача Propaganda Detection является актуальной в условиях повсеместной цифровизации. Отсутствие исследования данной задачи в новостях на русском языке.

Решение

Разработать базовые модели для решения задач в области Propaganda Detection.

Propaganda Detection

1. 2017: Analytical study of language in Propagandistic News¹
2. 2019: Document-level Propaganda Detection²
3. 2019: Span Identification and Span Classification³
4. SemEval-20: Task 11: Detection of Propaganda Techniques in News Articles competition
5. SemEval-21: Task 6: Multimodal Propaganda Detection: Text and Images
6. SemEval-23: Task 3: Detecting the Category, the Framing, and the Persuasion Techniques in Online News in a Multi-lingual Setup

¹Truth of Varying Shades: Analyzing Language in Fake News and Political Fact-Checking.

²Proppy: Organizing the news based on their propagandistic content.

³Fine-Grained Analysis of Propaganda in News Articles.

Span Identification

Постановка задачи

Подзадача Span Identification - как многоклассовой классификации токенов.

- I $p(x^t, \theta)$ - параметрическое семейство моделей для описания распределения вероятности классов каждого токена, $\{x_i^t\}_{i=1}^N$ - последовательность токенов текста $t \in T$, где T - множество всех текстов
- I параметр $\theta \in \Theta$, где Θ - пространство параметров модели
- I $\{y_i^t\}_{i=1}^N$ - класс i -го токена формате one-hot в тексте t
- I $L_1(y, \hat{y}) = \sum_{c \in C} y_c \log(\hat{y}_c)$, где вектор \hat{y} - распределение вероятности классов, $c \in C$ - класс из множества допустимых классов

Задача поиска оптимальных параметров модели сформулирована как:

$$\frac{1}{|T|} \sum_{t \in T} \sum_{i \in t} L_1(y_i^t, p(x^t, \theta)_i) \rightarrow \min_{\theta \in \Theta}$$

Span Targeting

Постановка задачи

Подзадача Span Targeting - задача семантического сопоставления текстовых фрагментов.

- I $p(x_i^t, x_j^t, x^t, \theta)$ - параметрическое семейство моделей для описания семантической связи между предложениями, x_i^t, x_j^t, x^t - фрагмент, мишень, полный текст
- I параметр $\theta \in \Theta$, где Θ - пространство параметров модели
- I $y_{ij} \in [0, 1]$ - связаны ли мишень j и фрагмент манипуляции i
- I $L_2(y, \hat{y}) = \sum_{j \in J} \sum_{i \in I} y_{ij} \log(\hat{y}_{ij}) + (1 - y_{ij}) \log(1 - \hat{y}_{ij})$, где вектор \hat{y} - вероятность семантической связи фрагмента и мишени.

Задача поиска оптимальных параметров модели сформулирована как:

$$\frac{1}{|I| + |J|} \sum_{i \in I} \sum_{j \in J} L_2(y_{ij}, p(x_i^t, x_j^t, x^t, \theta)) \rightarrow \min_{\theta \in \Theta}$$

Оценка качества Span Identification

Метрики качества

Для оценки качества и сравнения моделей необходимо ввести метрики качества решения задачи.

Пусть M - множество токенов, выделенных моделью, E - множество токенов, выделенных экспертом. Введем точность и полноту.

$$\begin{aligned} C(m, e, h) &= \frac{|m \cap e|}{h}, \\ P(M, E) &= \frac{1}{|M|} \sum_{m \in M, e \in E} C(m, e, |m|) \\ R(M, E) &= \frac{1}{|E|} \sum_{m \in M, e \in E} C(m, e, |e|) \end{aligned} \tag{1}$$

XLM-RoBERTa + адаптеры

XLM-RoBERTa^a - мультязычная большая языковая модель.

Определение: *Адаптер* - оператор аппроксимации векторного представления.

Для каждой из задач необходимо дообучение своего адаптера. В обеих задачах в качестве адаптера используется обучаемый линейный оператор $L : H \rightarrow C$, где H - пространство векторных представлений, C - пространство классов.

^aUnsupervised Cross-lingual Representation Learning at Scale

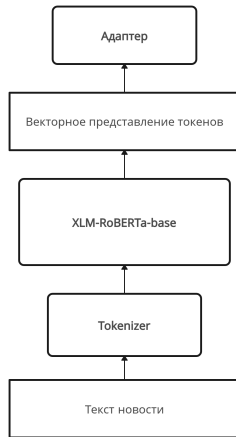


Рис.: Архитектура SI

Архитектура модели ST

Cross-Encoder - архитектура для решения задачи семантического сопоставления текстовых фрагментов.

Основная идея данного подхода заключается в том, что сопоставляемые фрагменты подаются на вход общему энкодеру В своей работе я расширяю данную архитектуру за счет использования контекста всей новости.

$$\begin{aligned}y_{span,target} &= h_1 = \text{Encoder}_1(span, target), \\ y_{full\ text} &= h_2 = \text{Encoder}_2(full\ text), \\ h_3 &= \text{concat}(y_{span,target}, y_{full\ text}), \\ \text{output} &= \text{adapter}(h_3)\end{aligned}$$

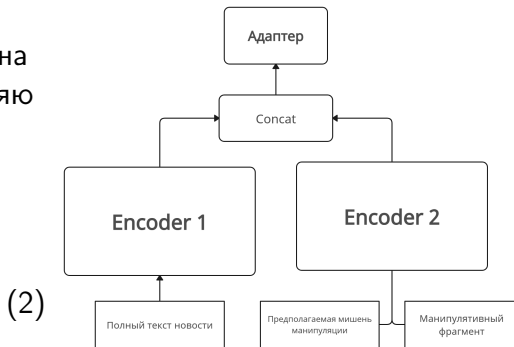


Рис.: Архитектура ST

Малоранговое приближение обновления весов

Малоранговое приближение матриц обновления весов в больших лингвистических моделях

Пусть $W_0 \in \mathbb{R}^{d \times k}$ - предобученные веса,
 $\Delta W \in \mathbb{R}^{d \times k}$ - изменение весов после дообучения
под конкретную задачу. ΔW - будем
аппроксимировать произведением двух матриц
малого ранга $B \in \mathbb{R}^{d \times r} = 0$, $A \in \mathbb{R}^{r \times k} \sim \mathcal{N}(0, \sigma^2)$,
где $r \ll \min(d, k)$

$$Wh = W_0x + \Delta Wx = W_0x + BAx \quad (3)$$

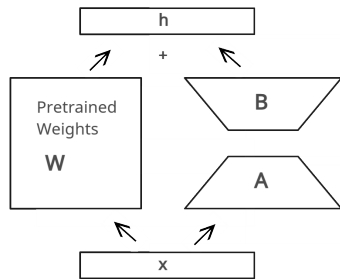


Рис.: LoRA

Эксперимент SI

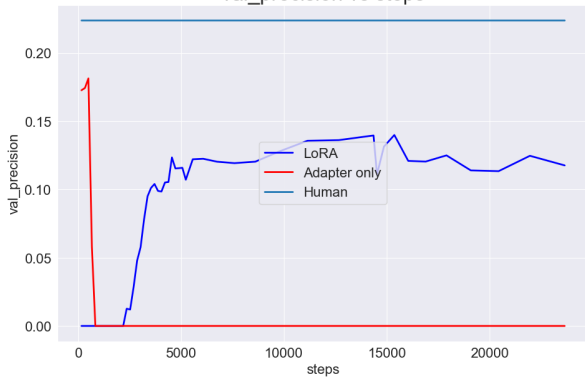
Модель для задачи выделения фрагментов обучалась в двух режимах. Первый - заморозка слоев XLM-RoBERTa-base и обучения только адаптера. Второй - XLM-RoBERTa-base обучалась с использованием малорангового приближения весов и адаптер в стандартном режиме.

Для обеих задач использовался следующий набор гиперпараметров:

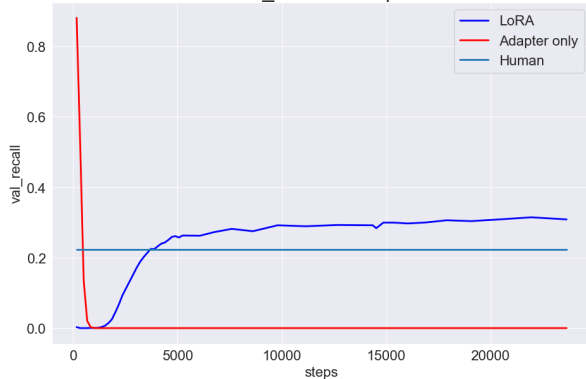
1. Оптимизатор AdamW
2. Темп обучения $1 \cdot 10^{-5}$
3. Шедулер: ExponentialLR с параметром затухания 0.99
4. Параметр Dropout-a: 0.2
5. 12 Encoder - слоев в XLM-RoBERTa-base

Визуализация эксперимента SI

val_precision vs steps

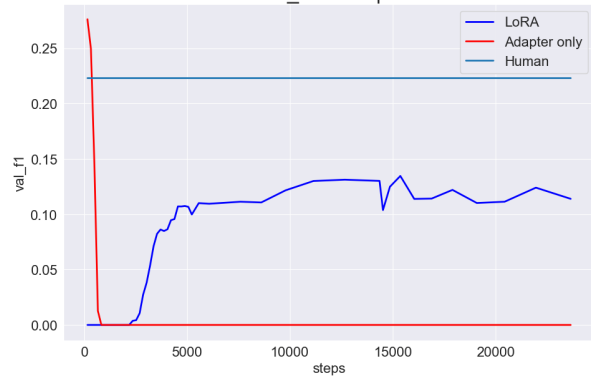


val_recall vs steps

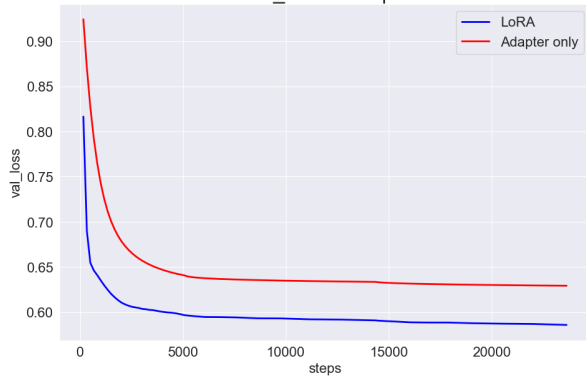


Визуализация эксперимента SI

val_f1 vs steps



val_loss vs steps



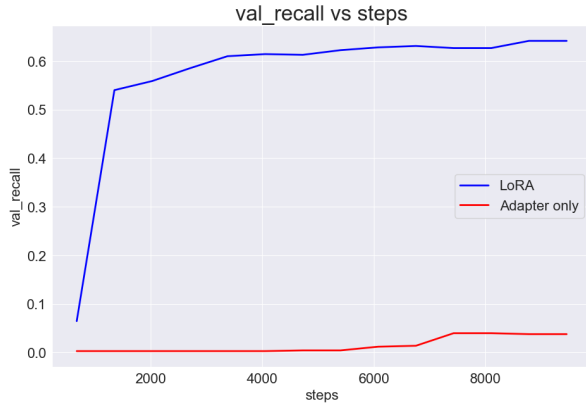
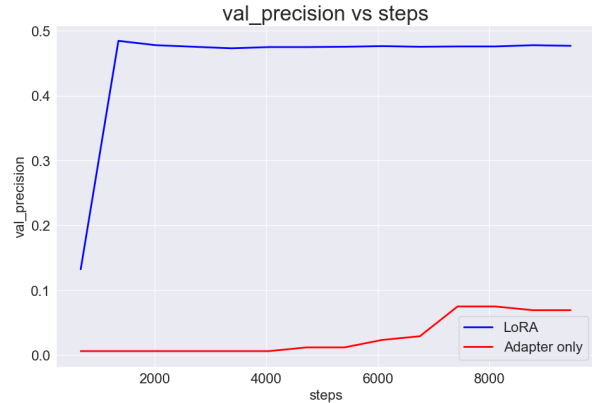
Эксперимент ST

Модель для задачи связи поиска мишени манипуляции обучалась в двух режимах. Первый - заморозка слоев XLM-RoBERTa-base и обучения только адаптера. Вторым - XLM-RoBERTa-base обучалась с использованием малорангового приближения весов и адаптер в стандартном режиме.

Для обеих задач использовался следующий набор гиперпараметров:

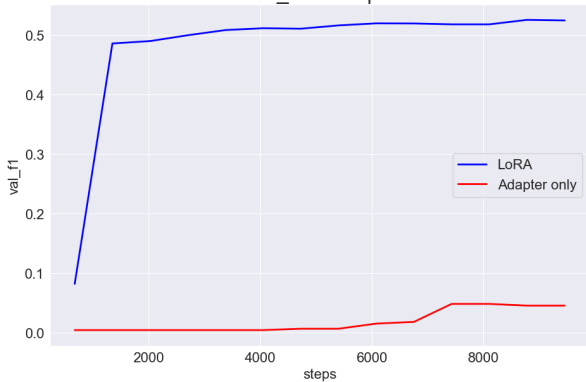
1. Оптимизатор AdamW
2. Темп обучения $1 \cdot 10^{-5}$
3. Шедулер: ExponentialLR с параметром затухания 0.99
4. Параметр Dropout-a: 0.2
5. 12 Encoder - слоев в XLM-RoBERTa-base

Визуализация эксперимента ST

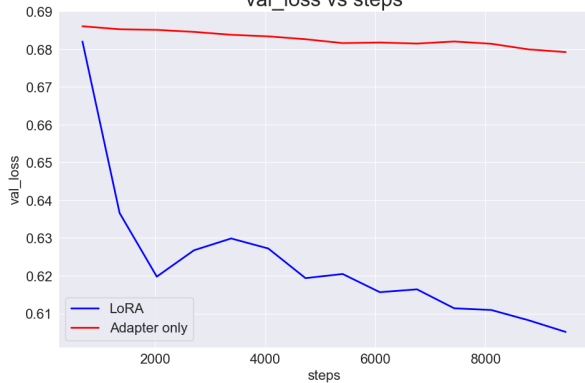


Визуализация эксперимента ST

val_f1 vs steps



val_loss vs steps



Анализ свойств полученных моделей

1. Модель с замороженными слоями после сходимости процесса обучения обладает низкой обобщающей способностью.
2. Модель, в которой происходило дообучение энкодера XLM-RoBERTa-base, после сходимости процесса обучения, обладает заметно лучшей способностью выделять фрагменты манипуляций. Достигнутая полнота выделения манипулятивных фрагментов превышает человеческие показатели, посчитанные между разметчиками обучающего корпуса.

Вывод:

Выделение манипулятивных фрагментов относится к классу задач NLP, в которых исследуемые тексты обладают сложным семантическим устройством языка.

Дообучения языковых моделей является необходимым для решения задач связанных с выявлением манипуляций.

Выносятся на защиту

Дипломная работа:

1. Разработаны и предложены архитектуры базовых моделей для задач выявления фрагментов манипуляций и поиска их мишеней.
2. Продемонстрирована эффективность обучения языковых моделей с использованием малорангового приближения.
3. Реализован и опубликован код для воспроизведения экспериментов из представленной работы⁴.

Публикации на конференциях:

1. IWANN-2023: Long-Term Hail Risk Assessment with Deep Neural Networks⁵

⁴<https://github.com/intsystems/Lukyanenko-BS-Thesis>

⁵<https://arxiv.org/abs/2209.01191>