

УДК 004.855.5

ИСПОЛЬЗОВАНИЕ ТРАНСФОРМЕРОВ ДЛЯ ЛИНГВИСТИЧЕСКОЙ ИДЕНТИФИКАЦИИ АВТОРА

Хазагаров А.А. (Университет ИТМО), Воробьева А.А. (Университет ИТМО)

Научный руководитель – кандидат технических наук, Воробьева А.А.
(Университет ИТМО)

Введение. Задача лингвистической идентификации автора является актуальной во многих областях, в том числе в информационной безопасности. Лингвистическая идентификация позволяет на основе некоторых лингвистических признаков установить, кому из авторов принадлежит данный текст. Для этого существует множество способов. Одним из последних и наиболее перспективных является использование нейронных сетей типа трансформер, которые себя хорошо зарекомендовали при обработке естественного языка.

Основная часть. Наибольшей популярностью для работы с текстом пользуются рекуррентные нейронные сети. Среди них выделяют RNN, LSTM и GRU [1, 2]. Они себя хорошо зарекомендовали в лингвистическом анализе, в том числе определение тональности текста, машинном переводе и многих других направлениях. Недавно представленная сеть трансформер, также предназначена для обработки текстов на естественном языке. Данная сеть получила широкое распространение в задачах машинного перевода, где показывает результаты сравнимые, а в некоторых моментах и превосходящие, чем рекуррентные нейронные сети. Также, важную часть при работе с текстом является векторизация и способ извлечения признаков. Во многом, от такого, какие из этих параметров будут выбраны, зависит точность и производительность алгоритма. В данной работе мы сравним использование трансформеров в задаче лингвистической идентификации автора с рекуррентными нейронными сетями. Определим наилучший способ векторизации текста и оптимальную структуру сети.

Выводы. В качестве результатов можно выделить, что нейронная сеть типа трансформер и n-граммы показали наилучший результат по точности и потреблению ресурсов для лингвистической идентификации автора.

Список использованных источников:

1. Dmitrin Y.V. Comparison of deep neural network architectures for authorship attribution of russian social media texts / Dmitrin Y.V., Botov D.S., Klenin J.D., Nikolaev I.E. // *Komp'yuternaja lingvistika i intellektual'nye tehnologii* – 2018.
2. Khazarov, A. Preventing Hidden Information Leaks Using Author Attribution Methods and Neural Networks / A. Khazarov, A. Vorobeva, V. Korzhuk // *Proceedings of the 29th Conference of Open Innovations Association FRUCT*. – 2021.