

УДК 004.912

В. И. Бондаренко, В. О. Елисеев, Т. В. Ермоленко
Федеральное государственное бюджетное образовательное учреждение
высшего образования «Донецкий государственный университет»
283001, Донецкая Народная Республика, г. Донецк, ул. Университетская, 24

АНАЛИЗ ЭФФЕКТИВНОСТИ ГЛУБОКИХ ЯЗЫКОВЫХ МОДЕЛЕЙ ДЛЯ ЗАДАЧИ ОПРЕДЕЛЕНИЯ ТОНАЛЬНОСТИ РУССКОЯЗЫЧНЫХ ТЕКСТОВ*

V. I. Bondarenko, V. O. Eliseev, T. V. Yermolenko
Federal State Budgetary Educational Institution of Higher Education "Donetsk State University"
283001, Donetsk People's Republic, Donetsk, University st, 24

ANALYZING THE EFFECTIVENESS OF DEEP LANGUAGE MODELS FOR THE TASK OF TONE DETECTION IN RUSSIAN-LANGUAGE TEXTS

В статье проведено исследование качества работы глубоких нейросетей, полученных в результате тонкой настройки, для задачи определения тональности текстов на русском языке. В качестве базовых языковых моделей используются RuGPT-3 и RuBERT. Тонкая настройка осуществляется путем замены последнего линейного слоя на линейный классифицирующий с числом выходов, соответствующим числу классов (нейтральный, позитивный, негативный). Для анализа эффективности моделей в качестве метрик использованы точность классификации и средневзвешенная F1-мера. Численные исследования показали, что после тонкой настройки RuGPT-3 имеет более высокое качество классификации: в среднем точность на 8.45% больше по сравнению RuBERT, а усредненная F1-мера – на 10.95%. Однако, модели, использующие архитектуру GPT, имеют более низкую скорость работы.

Ключевые слова: языковая модель, обработка естественного языка, анализ тональности текстов, тонкая настройка нейронных сетей, GPT, BERT.

The article describes the process of solving the task of sentiment analysis across texts of varying lengths, such as customer reviews and news articles. A methodology involving fine-tuning machine learning models based on RuGPT-3 and RuBERT is proposed, achieved through the substitution of the last linear layer with a classification layer having outputs corresponding to the number of classes (neutral, positive, negative). Research indicates the advantages of utilizing RuGPT-3-based models, revealing a notable increase in predictive quality despite their lower operational speed. Additionally, a comparison of models trained on one text type to predict sentiments in another was conducted. The results show that models trained on news articles exhibit slightly superior classification of reviews. However, the resulting accuracy falls short for the multimodal application of trained models.

Key words: language model, natural language processing, sentiment analysis, fine-tuning, GPT, BERT.

* Работа выполнена при финансовой поддержке Минобрнауки России в рамках научной темы "Разработка и совершенствование интеллектуальных методов классификации и прогнозирования для задач распознавания образов и моделирования информационных процессов" FREM-2024-0001 (Регистрационный номер 1023111000141-9-1.2.1)

Введение

Анализ тональности текста (*Sentiment Analysis, SA*) – одна из задач обработки естественного языка (*Natural Language Processing, NLP*), связанная с определением эмоциональной оценки авторов (мнений) относительно некоторых объектов, о которых идет речь в тексте. SA-задача сводится к задаче классификации текста, классами тональных оценок, как правило, являются позитивная, негативная, нейтральная (не имеет эмоциональной окраски). Анализ тональности текста имеет широкое практическое применение – от анализа общественной реакции на всевозможные события по публикациям пользователей в социальных сетях, до ранжирования товаров в интернет-магазинах и торговых площадках по отзывам покупателей и формирования выдачи новостных записей с учетом их эмоциональной окраски. Большую заинтересованность в создании SA-систем проявляют и представители среднего и малого бизнеса, которые заинтересованы в автоматическом мониторинге реакции клиентов на предоставляемые услуги.

Наиболее точную и надежную классификацию текста позволяют получить методы глубокого обучения, которые могут автоматически извлекать неявные признаки из большого объема текстовых данных. На сегодняшний день самые эффективные модели для задач обработки текста основаны на архитектуре Transformer, лидерами являются GPT (*Generative Pretraining Transformer*) [1], [2] и BERT (*Bidirectional Encoder Representations from Transformers*) [3].

Главными проблемами, с которыми сталкиваются разработчики NLP-систем, использующие глубокое обучение, являются отсутствие: 1) необходимых вычислительных ресурсов; 2) корректно аннотированных наборов данных для построения качественных моделей [4]. Частично эти трудности решаются с помощью тонкой настройки (*fine-tuning*) уже существующих моделей машинного обучения.

Настоящая статья посвящена вопросам трансферного обучения (*transfer learning*) в задачах анализа тональности текстов больших (новостной корпус) и сравнительно малых (корпус отзывов покупателей) размеров, в частности, возможности тонкой настройки модели GPT и релевантности данного решения.

Цель данного исследования – осуществить тонкую настройку моделей RuGPT-3 и RuBERT для задачи анализа тональности русскоязычных текстов различного объема, а также провести анализ эффективности полученных моделей.

Для достижения поставленной цели необходимо решить следующие **задачи**:

- 1) сформировать для обучения моделей размеченный текстовый корпус, содержащий сбалансированные классы;
- 2) используя методы машинного обучения и группировки данных, провести тонкую настройку моделей RuGPT-3 и RuBERT для задачи определения тональности русскоязычных текстов;
- 3) сравнить эффективность полученных моделей, используя в качестве метрик точность классификации и средневзвешенную F1 меру (*macro F1 score*) [5].

Подготовка данных для обучения моделей

В данной работе используются два набора аннотированных по тональности текстовых корпусов, предоставленных ресурсом Kaggle:

- 1) корпус отзывов покупателей [6];
- 2) новостной корпус [7].

Оба корпуса содержат русскоязычные тексты трех классов: нейтральные, позитивные, негативные.

Предварительный анализ данных проводится с целью выявления и исключения повторяющихся данных, проверки сбалансированности распределения классов [8], а также выявления возможных стратегий оптимизации процесса обучения модели [9]. Результатом данного этапа являются предположения по организации процесса обучения и разделения данных с учетом ограниченности вычислительных ресурсов.

Текстовый корпус отзывов покупателей содержит 90 тыс. записей, при этом содержит ряд повторяющихся записей. После удаления дубликатов, объем набора данных составил 87321 запись.

Была построена гистограмма частот появления классов отзывов в очищенном от повторяющихся записей наборе данных. Полученная гистограмма приведена на рис. 1; на основе её анализа можно сделать вывод, что классы являются сбалансированными, т.е. нет необходимости в применении методов семплирования данных при несбалансированных классах [10].

После удаления повторяющихся записей, набор данных с покупательскими отзывами содержит большое количество записей. Для обработки корпуса (выборки) такого размера требуется значительный объем оперативной памяти, для сокращения ресурсоемкости размер обучающей выборки уменьшен путем отбора случайным образом 15% от всего датасета [11]. Объем полученной таким образом подвыборки составил 13098 записей.

Классы маркировались следующим образом: негативному классу текста соответствует значение 0, нейтральному – 1, позитивному – 2.

Для проверки сбалансированности классов построены гистограммы распределения классов в сформированных наборах данных отзывов (рис. 1).

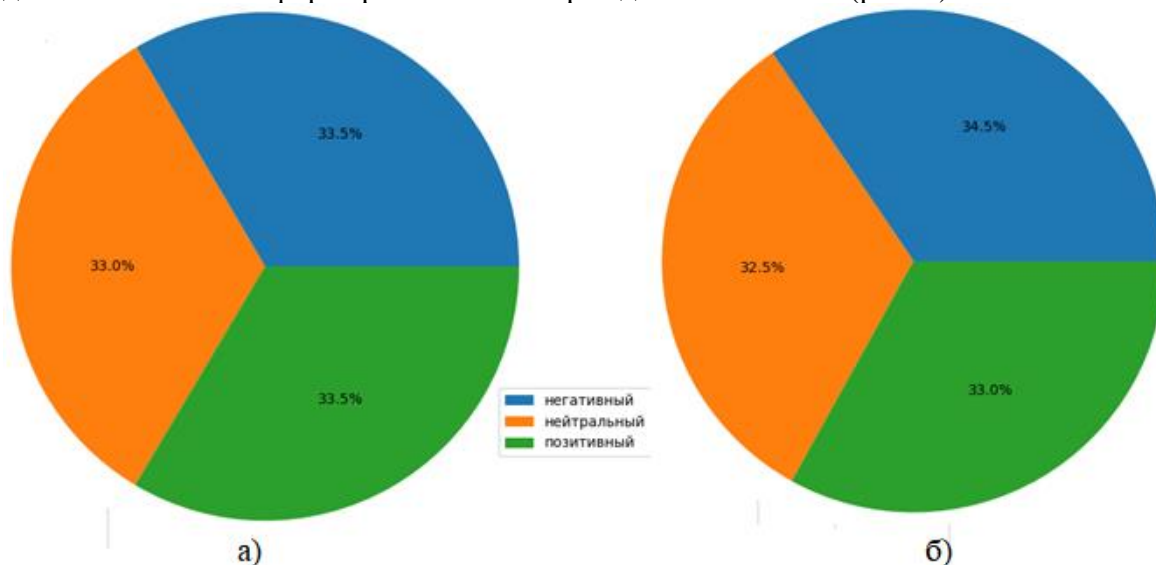


Рисунок 1 – Круговые диаграммы распределения классов на корпусе отзывов полного объема (а) и на 15%-ной подвыборке корпуса отзывов (б)

Как видно из рис. 1, как в случае полного набора данных, так и в случае использования случайной подвыборки, проблемы несбалансированности классов нет.

Модель ruGPT-3 имеет несколько конфигураций архитектуры: small, medium и large, отличающихся числом параметров (125, 355 и 750 млн соответственно). Их характеристики приведены в табл. 1.

В рамках данного исследования используются конфигурации small и medium. Уменьшение размера обучающей применялся для обучения модели RuGPT-3 medium, поскольку она требует значительно большего числа вычислительных ресурсов, чем модель RuGPT-3 small, которая обучалась на полном наборе данных.

Таблица 1 – Характеристики различных конфигураций RuGPT-3

	small	medium	large
Размер входной последовательности	2048	2048	2048
Размер скрытого слоя	768	1024	1536
Количество голов внимания	12	16	16
Количество скрытых слоёв	12	24	24
Размер словаря	50264	50257	50257
Общее число параметров сети, млн	125	355	760

Текстовый корпус новостей после удаления дубликатов содержит около 8 тыс. записей, гистограмма распределения классов в нем показана на рисунке 2.

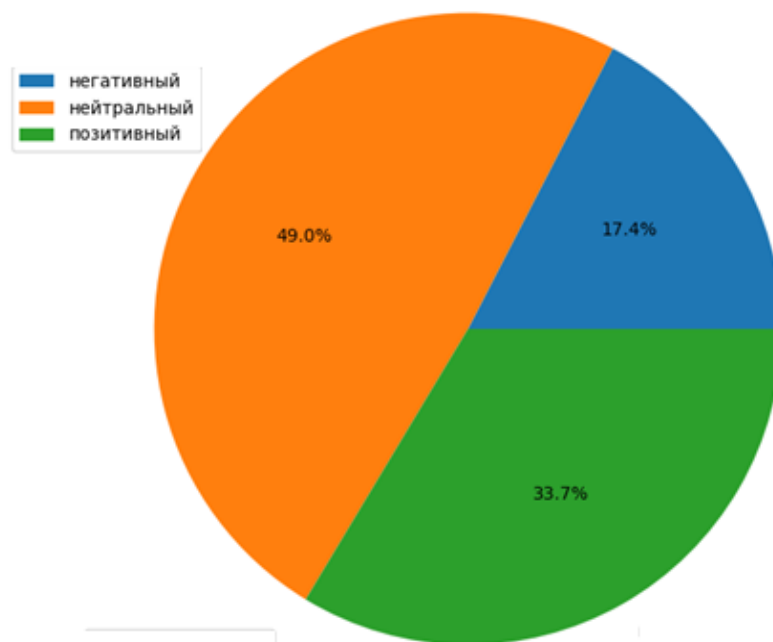


Рисунок 2 – Круговая диаграмма распределения классов на новостном корпусе

Как видно из рис. 2, мажоритарным классом является класс нейтральных отзывов, однако, существенного дисбаланса классов не наблюдается, поэтому нет необходимости в применении методов семплирования данных для выравнивания классов.

Обучение и тонкая настройка моделей классификации текстов

В рамках данного исследования обучены следующие классификаторы:

- модели RuGPT-3 конфигурации small и RuBERT, обученные на полном корпусе отзывов покупателей;
- модель RuGPT-3 конфигурации medium, обученная на подвыборке из 15% от корпуса отзывов покупателей;
- модели RuGPT-3 конфигурации small и RuBERT, обученные на корпусе новостей.

Для обучения моделей классификации текстов необходимо провести процесс преобразования текста в числовой вектор, т.е. провести процесс эмбединга. Первым этапом для осуществления этой процедуры является токенизация (деление текста на слова или подслова) [12]. Современные языковые модели для решения этой сложной задачи используют простой алгоритм архивирования данных – алгоритм BPE (Byte Pair Encoding), созданный в 1994 году.

BPE – это восходящий алгоритм токенизации подслов, который изучает словарь подслов определенного размера (размер словаря является гиперпараметром). Идея BPE-токенизации заключается в том, что более частым словам следует присваивать уникальные идентификаторы, тогда как менее частые слова следует разбивать на единицы подслов, которые лучше всего сохраняют свое значение [13].

В рамках данного исследования классификаторы, основанные на RuBERT, использовали BPE-токенизацию с размером словаря, равным 120138 токенов [6], классификаторы, основанные на RuGPT-3, – BBPE-токенизацию (Byte-level Byte Pair Encoding) [14]. Эта модификация BPE работает не с текстом, а напрямую с его байтовым представлением. Размер словаря для моделей архитектуры GPT составил 50257, что соответствует размеру словаря для GPT-2 [15], на основе которой и была построена модель RuGPT-3.

Процесс тонкой настройки моделей семейства RuGPT-3 заключался в замене выходного линейного слоя (на вход поступает тензор размером равным размеру скрытого слоя, а на выходе размер равен размеру словаря) на классифицирующий линейный слой, где размер выхода соответствует количеству классов для определения тональности (3 класса). Далее для всех слоёв, кроме вновь добавленного, останавливалось обновление градиентов при обучении, т.е. производилось обновление градиентов лишь для последнего классифицирующего слоя на предоставленных наборах данных. Этот процесс называется «заморозкой» («freeze») слоёв [16].

Для тонкой настройки моделей семейства RuBERT использовался механизм «заморозки» всех слоёв, кроме классифицирующего линейного, но вид этого слоя (количество выходов) задавалось при инициализации объекта модели, а не как в случае выше – путем удаления старого слоя и внедрения собственного, с тремя выходами.

В качестве функции потерь была выбрана кросс-энтропия (Cross-Entropy Loss) [17], измеряющая расхождение вероятностных распределений истинных ответов и прогнозов модели, и наиболее точно подходящая для задачи классификации нескольких классов.

В качестве метода оптимизации для обучения моделей был выбран AdamW, где «W» означает «Weight Decay» [18], – вариант оптимизатора Adam, который корректирует реализацию снижения веса. В стандартном Adam перед вычислением скорости адаптивного обучения применяется снижение веса: часть весов вычитается перед обновлением весов, это затухание включено в скользящие средние, что может привести к неправильному обновлению при адаптации скорости обучения. AdamW исправляет это, отделяя снижение веса от скорости адаптивного обучения с помощью применения затухания веса непосредственно к весам, что аналогично традиционному затуханию веса. Такой подход приводит к более точному обновлению [19], что подтверждено работой Meta Research по созданию Llama 2 [20].

Проблема «взрыва градиента» решена при помощи градиентного клиппирования по норме [21].

Сравнительный анализ качества полученных моделей

Для реализации языковых моделей и классификаторов использовалась библиотека глубокого обучения PyTorch [16]. Обучение глубоких нейросетей требует значительных вычислительных мощностей, которые для данного исследования обеспечил облачный сервис, предоставляемый платформой Kaggle (графический процессор P100, содержащий 16 Гб видеопамяти). Датасет для каждой модели разделен на обучающую и тестовую выборку в соотношении 80/20.

Исходный код проекта размещен в открытом доступе [22].

Здесь и далее используются следующие обозначения моделей:

GPT_S_Rev – RuGPT-3 small на корпусе отзывов полного объема;

GPT_M – RuGPT-3 medium на 15%-ной подвыборке корпуса отзывов;

GPT_S_News – RuGPT-3 small на новостном корпусе;

BERT_Rev – RuBERT на 15%-ной подвыборке корпуса отзывов;

BERT_News – RuBERT на новостном корпусе.

В табл. 2 приведены значения гиперпараметров исследуемых моделей и время их обучения.

Таблица 2 – Гиперпараметры исследуемых моделей классификации

	GPT_S_Rev	GPT_M	GPT_S_News	BERT_Rev	BERT_News
количество эпох	5	15	12	15	12
размер пакета	32	32	16	32	16
коэффициент обучения	0.001	0.01	0.01	0.01	0.01
скорость снижения веса	0.01	0.0001	0.0001	0.0001	0.0001

На рис. 3-7 изображены графики зависимости от количества эпох точности и потерь каждой из моделей во время обучения.

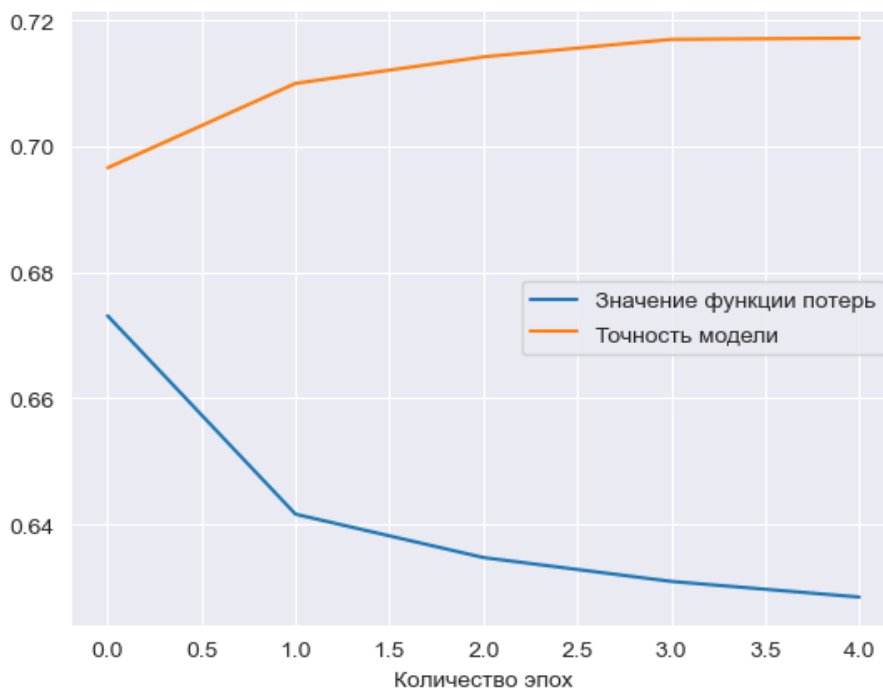


Рисунок 3 – Зависимость значений точности и loss-функции от количества эпох модели GPT_S_Rev во время обучения

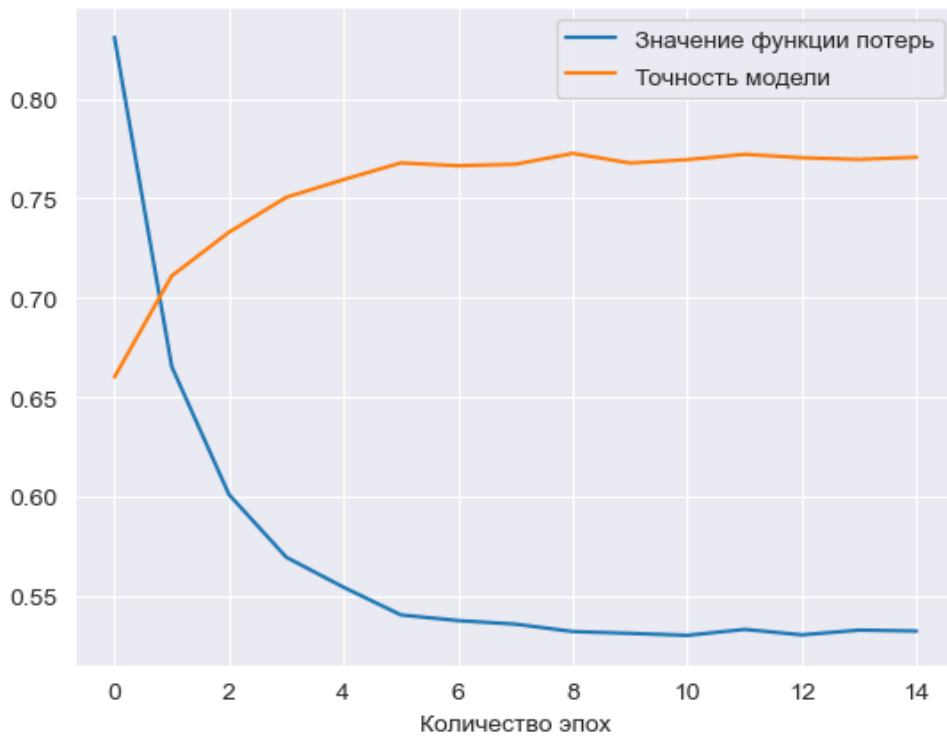


Рисунок 4 – Зависимость значений точности и loss-функции от количества эпох модели GPT_M во время обучения

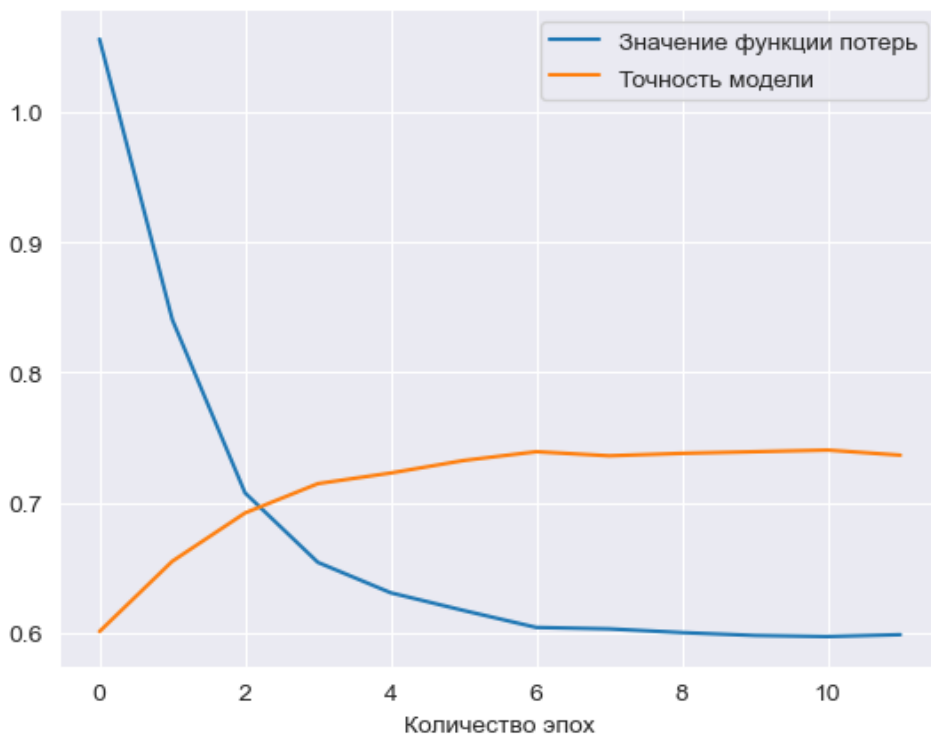


Рисунок 5 – Зависимость значений точности и loss-функции от количества эпох модели GPT_S_News во время обучения

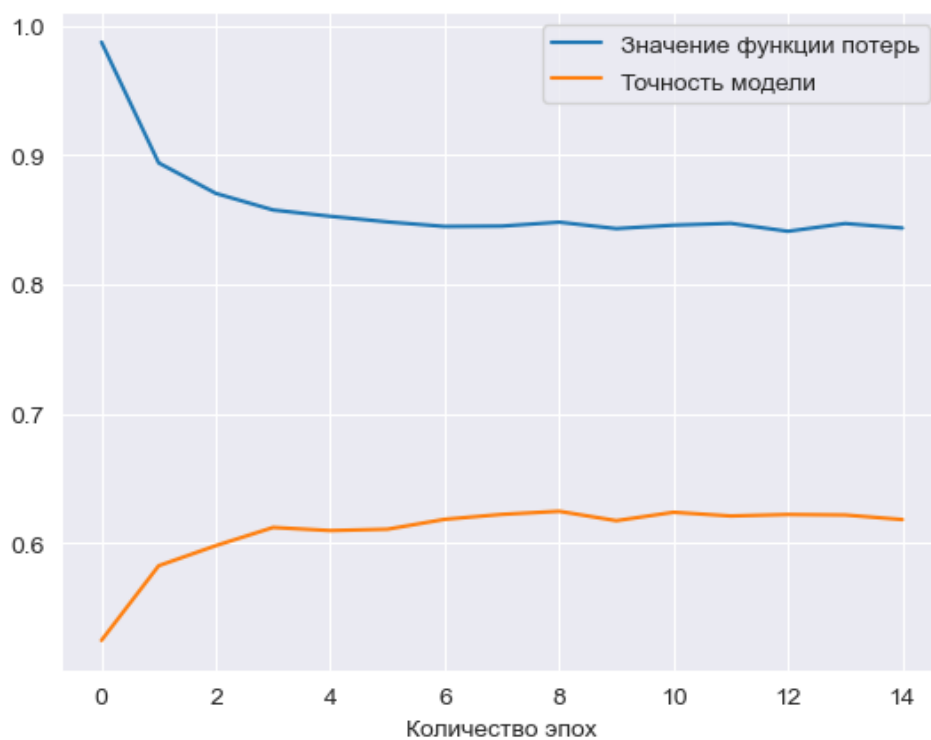


Рисунок 6 – Зависимость значений точности и loss-функции от количества эпох модели BERT_Rev во время обучения

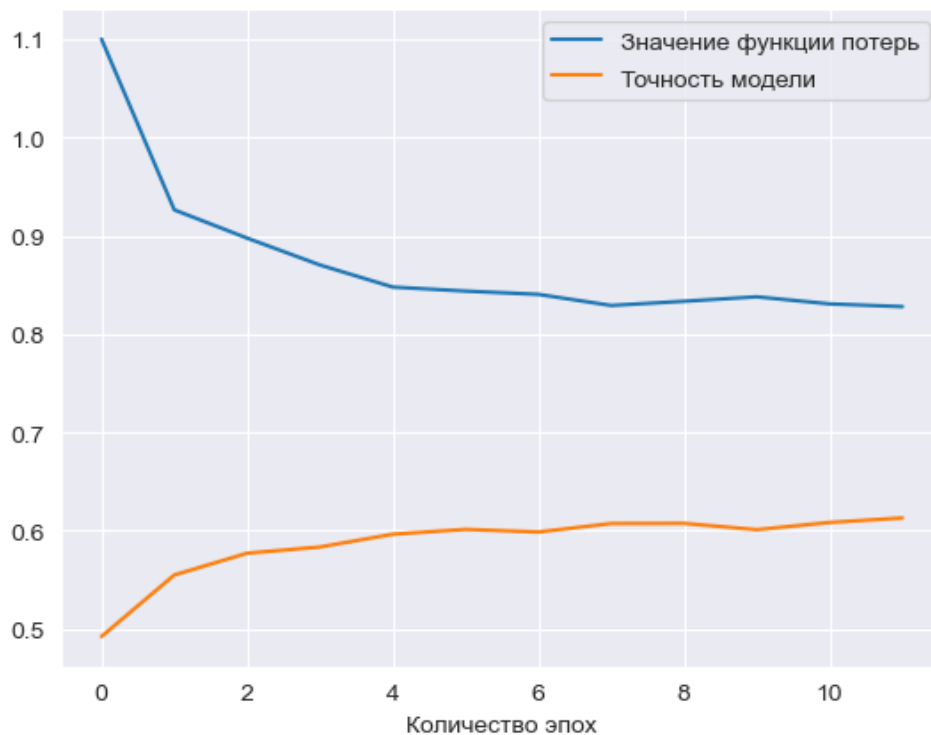


Рисунок 7 – Зависимость значений точности и loss-функции от количества эпох модели BERT_News во время обучения

Для проведения сравнительного анализа эффективности моделей для каждой из них выполнены оценки точности классификации, усредненной F1-меры, а также времени обучения и обработки тестового набора данных. Результаты сведены в табл. 3.

Таблица 3 – Показатели эффективности исследуемых моделей классификации

	GPT_S_Rev	GPT_M	GPT_S_News	BERT_Rev	BERT_News
Время обучения, мин	553	108	197	56	91
Время обработки тестовых данных, с	521	1280	358	331	91
Точность классификации	0.723	0.718	0.733	0.650	0.638
Усредненная F1-мера	0.721	0.715	0.726	0.647	0.582

Как можно заключить из численных исследований, модели с архитектурой GPT-3 превосходят модели BERT более, чем на 8%, хотя скорость обучения BERT значительно выше.

Следует отметить, что значение F1-меры более, чем 0.7, на сегодняшний день является хорошим результатом. По данным соревнования Kaggle для задачи анализа тональности русского текста [4] значение F1-меры классификатора с архитектурой RuGPT-3 составило 0.65138, что соответствует топ-10 результатов по данному критерию.

Выводы

В результате исследования проведен сравнительный анализ эффективности работы тонко настроенных глубоких языковых моделей архитектуры на RuGPT-3 и RuBERT для задачи определения тональности текстов. Численные исследования показали превосходство в точности классификации моделей, основанных на RuGPT-3: средний прирост в точности для текстов разного объема составил **8.45%** и в значении усредненной F1-меры – **10.95%**, что может говорить о том, что тонкая настройка именно этих моделей наиболее оправдана. Стоит отметить, что скорость работы моделей, основанных на RuBERT в некоторых применениях значительно выше.

Если сравнить качество работы моделей, основанных на RuGPT-3 small и RuGPT-3 medium, обученных на полном наборе отзывов и его части соответственно, то заметных различий в их точности не наблюдается, но модель конфигурации medium работает примерно в 2 раза медленнее, что говорит о нецелесообразности ее практического применения в этой области задач. Необходимо провести дополнительные эксперименты с увеличением объема тренировочных данных для модели RuGPT-3 с конфигурацией medium, чтобы убедиться в изменении метрик в положительную сторону в этом случае.

Дальнейшие эксперименты по улучшению итоговых метрик связаны с обучением всей сети, а не последнего линейного слоя, а также построением собственной классифицирующей сети, состоящей из сверточных и рекуррентных слоев.

Список литературы

1. Radford, A. Improving Language Understanding by Generative Pre-Training [Электронный ресурс]. URL: <https://gwern.net/doc/www/s3-us-west-2.amazonaws.com/d73fdc5ffa8627bce44dcda2fc012da638ffb158.pdf>.
2. Zmitrovich, D. A Family of Pretrained Transformer Language Models for Russian / D. Zmitrovich, A. Abramov, A. Kalmykov, M. Tikhonova, E. Taktasheva, D. Astafurov, M. Baushenko, A. Snegirev, T. Shavrina, S. Markov, V. Mikhailov, A. Fenogenova. 2023.
3. Kuratov Y. Adaptation of deep bidirectional multilingual transformers for Russian language / Y. Kuratov, M. Arkhipov. *Komp'yuternaja Lingvistika i Intellektual'nye Tehnologii*. 2019. Тт. 2019-Май.

4. Ермоленко, Т.В. Разработка алгоритмов и языковых моделей мультиязычной системы автоматического аннотирования текстов разных жанров / Т.В. Ермоленко, В.И. Бондаренко, Я.С. Пикалёв. *Вестник ДонНУ. Серия Г: Технические науки*. 2023. Т. 2. С. 22-43.
5. Humphrey A. Machine-learning classification of astronomical sources: estimating F1-score in the absence of ground truth / A. Humphrey, W. Kuberski, J. Bialek, N. Perrakis, W. Cools, N. Nuyttens, H. Elakhrass, P.A.C. Cunha. *Monthly Notices of the Royal Astronomical Society: Letters*. 2022. Т. 517. № 1.
6. Russian-language reviews | Kaggle [Электронный ресурс]. URL: <https://www.kaggle.com/datasets/laytsw/reviews>.
7. Sentiment Analysis in Russian | Kaggle [Электронный ресурс]. URL: <https://www.kaggle.com/competitions/sentiment-analysis-in-russian/data>.
8. Ермоленко, Т.В. Классификация ошибок в тексте на основе глубокого обучения / Т.В. Ермоленко. *Проблемы искусственного интеллекта*. 2019. Т. 3. – № 14. – С. 47-57.
9. Пикалёв Я.С. Разработка системы нормализации текстовых корпусов / Я.С. Пикалёв. *Проблемы искусственного интеллекта*. 2022. Т. 2. № 25. С. 64-78.
10. Ryan Hoens T. Imbalanced datasets: From sampling to classifiers / T. Ryan Hoens, N. V. Chawla. *Imbalanced Learning: Foundations, Algorithms, and Applications*. 2013.
11. Бондаренко, В.И. Классификация научных текстов с помощью методов глубокого машинного обучения / В.И. Бондаренко. *Вестник Донецкого национального университета. Серия Г. Технические науки*. 2021. Т. 3. С. 69-77.
12. Webster J.J. *Tokenization as the initial phase in NLP* / J.J. Webster, C. Kit. 1992.
13. Пикалёв Я.С. Адаптация нейросетевой модели ALBERT для задачи языкового моделирования / Я.С. Пикалёв, Т.В. Ермоленко // *Проблемы искусственного интеллекта*. 2020. № 3(18). С. 111-122.
14. Wang, C. Neural machine translation with byte-level subwords / C. Wang, K. Cho, J. Gu. *AAAI 2020 - 34th AAAI Conference on Artificial Intelligence*. 2020.
15. Radford Alec. Language Models are Unsupervised Multitask Learners | Enhanced Reader [Электронный ресурс]. URL: <https://life-extension.github.io/2020/05/27/GPT技术初探/language-models.pdf>
16. Пикалёв, Я.С. Разработка автоматической системы трансформации английских вставок в русских текстах с применением глубокого обучения / Я.С. Пикалёв, Т.В. Ермоленко. *Проблемы искусственного интеллекта*. 2019. № 2 (13). С. 74-86.
17. Mao A. Cross-Entropy Loss Functions: Theoretical Analysis and Applications / A. Mao, M. Mohri, Y. Zhong. 2023.
18. Loshchilov I. Decoupled weight decay regularization / I. Loshchilov, F. Hutter. *7th International Conference on Learning Representations, ICLR 2019*. 2019.
19. Xie Z. On the Overlooked Pitfalls of Weight Decay and How to Mitigate Them: A Gradient-Norm Perspective / Z. Xie, Z. Xu, J. Zhang, I. Sato, M. Sugiyama. 2020.
20. Llama 2: Open Foundation and Fine-Tuned Chat Models / H. Touvron и др. 2023.
21. Takagi S. On the Effect of Pre-training for Transformer in Different Modality on Offline Reinforcement Learning / S. Takagi. *Advances in Neural Information Processing Systems*. 2022 Т. 35.
22. Eliseev Vadim. Sentiment Analysis Fine Tuned [Электронный ресурс]. URL: <https://github.com/EliseevVadim/sentiment-analysis-fine-tuned>.

References

1. Radford A. Improving Language Understanding by Generative Pre-Training [Electronic resource]. URL: <https://gwern.net/doc/www/s3-us-wes2.amazonaws.com/d73fde5ffa8627bce44dcda2fc012da638ffb158.pdf>.
2. Zmitrovich D. A Family of Pretrained Transformer Language Models for Russian / D. Zmitrovich, A. Abramov, A. Kalmykov, M. Tikhonova, E. Taktasheva, D. Astafurov, M. Baushenko, A. Snegirev, T. Shavrina, S. Markov, V. Mikhailov, A. Fenogenova. 2023.
3. Kuratov Y. Adaptation of deep bidirectional multilingual transformers for Russian language / Y. Kuratov, M. Arkhipov // *Komp'yuternaja Lingvistika i Intellektual'nye Tehnologii*. 2019. Тт. 2019-Май.
4. Yermolenko T.V. Development of algorithms and language models for a multi-language system of automatic summary of texts of different genres / T.V. Yermolenko, V.I. Bondarenko, Ya.S. Pikalyov // *Vestnik of the Donetsk National Univesity. Series D: Technical sciences*. 2023. – № 2. P. 22-43.
5. Humphrey A. Machine-learning classification of astronomical sources: estimating F1-score in the absence of ground truth / A. Humphrey, W. Kuberski, J. Bialek, N. Perrakis, W. Cools, N. Nuyttens, H. Elakhrass, P.A.C. Cunha // *Monthly Notices of the Royal Astronomical Society: Letters*. 2022. Т. 517. № 1.
6. Russian-language reviews | Kaggle [Electronic resource]. URL: <https://www.kaggle.com/datasets/laytsw/reviews>.
7. Sentiment Analysis in Russian | Kaggle [Electronic resource]. URL: <https://www.kaggle.com/competitions/sentiment-analysis-in-russian/data>.

8. Yermolenko T.V. Classification of errors in the text based on deep learning / T.V. Yermolenko // Problems of Artificial Intelligence. – 2019. – № 3(14). – P. 47-57.
9. Pikalyov, Ya. S. The development of a text corpora normalization system / Ya. S. Pikalyov // Problems of Artificial Intelligence. – 2022. – № 2(25). – P. 64-78.
10. Ryan Hoens T. Imbalanced datasets: From sampling to classifiers / T. Ryan Hoens, N. V. Chawla // Imbalanced Learning: Foundations, Algorithms, and Applications. – 2013.
11. Bondarenko V.I. Classification of scientific texts using deep machine learning methods / В.И. Бондаренко // Vestnik of the Donetsk National Univesity. Series D. Technical sciences. – 2021. – № 3. – С. 69-77.
12. Webster J.J. Tokenization as the initial phase in NLP / J.J. Webster, C. Kit. – 1992.
13. Pikalyov, Ya. S. Adaptation of ALBERT neural network model for language modeling problem/ Ya. S. Pikalyov, T. V. Yermolenko// Problems of Artificial Intelligence. – 2020. – №3(18). – С. 111-122.
14. Wang C. Neural machine translation with byte-level subwords / C. Wang, K. Cho, J. Gu // AAAI 2020 - 34th AAAI Conference on Artificial Intelligence. – 2020.
15. Radford Alec. Language Models are Unsupervised Multitask Learners | Enhanced Reader [Электронный ресурс]. URL: <https://lile-extension.github.io/2020/05/27/GPT技术初探/language-models.pdf>
16. Pikalyov, Ya. S. The development of the automatic transformation of english accents in russian texts with the application of deep learning / Ya. S. Pikalyov, T.V. Yermolenko // Problems of Artificial Intelligence. – 2019. – № 2(13). – P. 74-86.
17. Mao A. Cross-Entropy Loss Functions: Theoretical Analysis and Applications / A. Mao, M. Mohri, Y. Zhong. – 2023.
18. Loshchilov I. Decoupled weight decay regularization / I. Loshchilov, F. Hutter // 7th International Conference on Learning Representations, ICLR 2019. – 2019.
19. Xie Z. On the Overlooked Pitfalls of Weight Decay and How to Mitigate Them: A Gradient-Norm Perspective / Z. Xie, Z. Xu, J. Zhang, I. Sato, M. Sugiyama. – 2020.
20. Touvron H. et al. Llama 2: Open foundation and fine-tuned chat models //arXiv preprint arXiv:2307.09288. – 2023.
21. Takagi S. On the Effect of Pre-training for Transformer in Different Modality on Offline Reinforcement Learning / S. Takagi // Advances in Neural Information Processing Systems. 2022. T. 35.
22. Eliseev Vadim. Sentiment Analysis Fine Tuned [Electronic resource]. URL: <https://github.com/EliseevVadim/sentiment-analysis-fine-tuned>.

RESUME

V. I. Bondarenko, V. O. Eliseev, T. V. Yermolenko

Analyzing the effectiveness of deep language models for the task of tone detection in Russian-language texts

The task of sentiment analysis, distinguishing between positive, neutral, and negative emotions in textual content, holds significant relevance across social media platforms for gauging public mood, in business sectors for monitoring user reviews of various products, and in the domain of recommendation systems. While primary implementation resides in the back-end of diverse B2C services, these classifiers also find utility within applied parsing programs for the former task.

This paper proposes a methodology involving fine-tuning a generative neural network based on RuGPT-3 by replacing the output linear layer with a classification layer. This approach is grounded on the assumption that the embedding preceding the network output aptly describes the input text, allowing a trainable linear layer to effectively handle the task with relatively modest training efforts.

Neural networks were re-trained using datasets of customer reviews and news articles. Results demonstrated satisfactory classification capabilities of both neural network types, achieving accuracy over 70% for small text and above 65% for large text, positioning them within the top 10 rankings in the respective Kaggle competition. Comparative analysis with fine-tuning a neural network based on RuBERT revealed that fine-tuning RuGPT-3 exhibited an increase in accuracy by 8.45% and a 10.95% enhancement in the macro F1-score.

РЕЗЮМЕ

В. И. Бондаренко, В. О. Елисеев, Т. В. Ермоленко
Анализ эффективности глубоких языковых моделей для задачи
определения тональности русскоязычных текстов

Задача анализа тональности текста (его эмоциональной окраски, которая может быть позитивной, нейтральной и негативной) является востребованной как в социальной среде (мониторинг настроения населения по их публикациям в социальных сетях), в бизнесе (мониторинг отзывов пользователей на различные товары), так и в сфере рекомендательных систем. Основным местом внедрения таких классификаторов является бэкенд различных B2C-сервисов, однако они могут служить и в прикладных программах-парсерах для первой задачи.

Предложен метод тонкой настройки генеративной нейронной сети, основанной на RuGPT-3, путем замены выходного линейного слоя на классифицирующий в связи с предположением о том, что эмбединг, находящийся перед выходом сети настолько хорошо описывает входной текст, что обычный тренируемый линейный слой справится с задачей хорошо, при сравнительно небольших затратах на процесс обучения.

Нейронные сети дообучались на данных о покупательских отзывах и новостных записях, результаты показали достаточную классифицирующую способность нейросетей обоих типов (точность для small text – свыше 70%, для large text – свыше 65%, что соответствует топ-10 в соответствующем соревновании на Kaggle).

При сравнении с тонкой настройкой нейронной сети, основанной на RuBERT, тонкая настройка RuGPT-3 дает прирост в точности равный 8.45% и 10.95% в значении усредненной F1-меры.

Бондаренко Виталий Иванович – кандидат технических наук, доцент кафедры компьютерных технологий физико-технического факультета, Федерального государственного бюджетного образовательного учреждения высшего образования "Донецкого Государственного Университета".

Область научных интересов: искусственный интеллект, интеллектуальный анализ данных, машинное обучение, математическое моделирование гидро- и теплофизических процессов, разработка пользовательских интерфейсов для прикладных программ моделирования.

Елисеев Вадим Олегович – стажер-исследователь лаборатории интеллектуальных систем Федерального государственного бюджетного научного учреждения Институт прикладной математики и механики.

Область научных интересов: искусственный интеллект, машинное обучение, нейронные сети, обработка естественного языка, генеративные и большие языковые модели.

Ермоленко Татьяна Владимировна – кандидат технических наук, доцент кафедры компьютерных технологий физико-технического факультета, Федерального государственного бюджетного образовательного учреждения высшего образования "Донецкого Государственного Университета".

Область научных интересов: Цифровая обработка сигналов, анализ данных, дискретная математика, теория алгоритмов, распознавание образов, обработка естественного языка, компьютерное зрение, машинное обучение, нейронные сети

Статья поступила в редакцию 14.02.2024.