

ИНСТИТУТ ИНТЕЛЛЕКТУАЛЬНЫХ КИБЕРНЕТИЧЕСКИХ СИСТЕМ

КАФЕДРА КИБЕРНЕТИКИ (№ 22)

09.03.04 «Программная инженерия»

НАПРАВЛЕНИЕ ПОДГОТОВКИ

Учебно-исследовательская работа на тему:
**Разработка и сравнительный анализ моделей
машинного обучения для ответа на связанные с
медициной контекстно-зависимые вопросы**

Студент: Плотников В. И.

Группа: Б21-514

Научный руководитель: Сбоев А. Г.



Реферат

Общий объем основного текста, без учета приложений - 21 страницы, с учетом приложений 29.

Количество использованных источников 11.

Количество приложений 2.

Количество рисунков 7.

Количество таблиц 2.

Ключевые слова: NLP, text classification, BERT, transfer learning.



Актуальность работы

С ростом объема медицинских данных возникает необходимость в эффективных методах их обработки и анализа. Одним из направлений в этой области является задача оценки способности моделей "извлекать информацию" из текста и правильно отвечать на уточняющие вопросы, в том числе связанные с медициной.

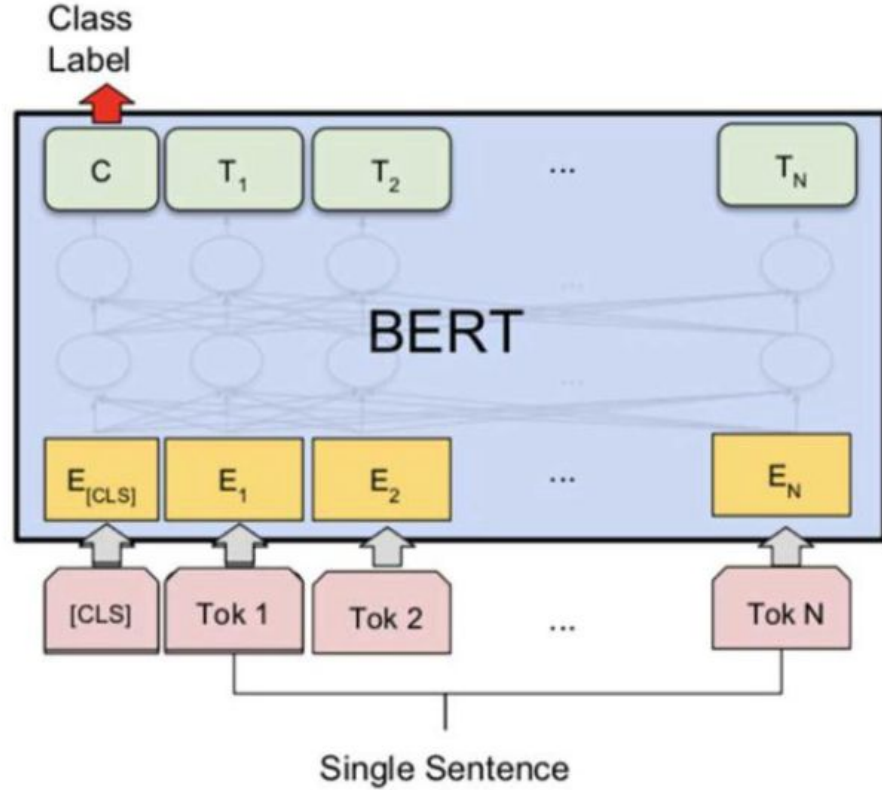
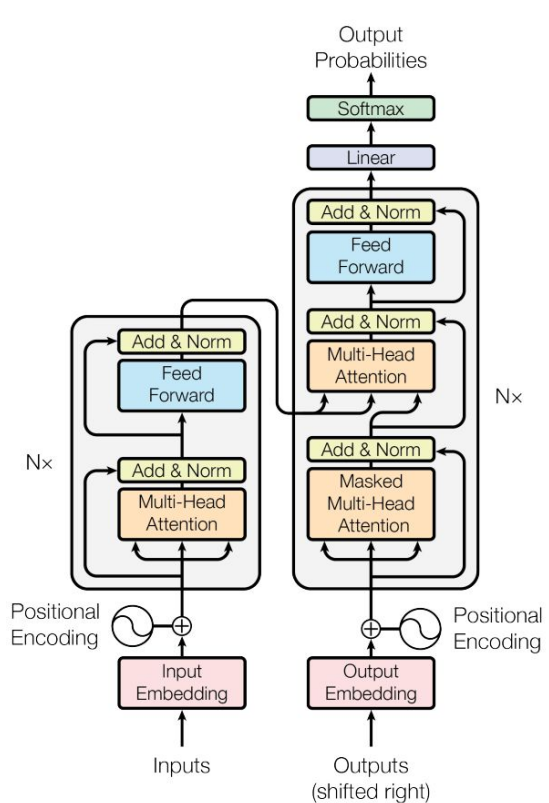


Цель УИР

- Разработка моделей для решения задачи ответа на связанный с медициной контекстно-зависимый вопрос



Подходы к решению задачи классификации в условиях недостатка данных и вычислительных мощностей



Задачи УИР

- 1) Нахождение модели, натренированной на похожих данных
- 2) Нахождение размеченных данных для дообучения модели
- 3) Обучение на новых данных



Моделирование алгоритма для дообучения сети

Для реализации необходимы следующие функции:

- Функция векторизации вопроса и контекста для подачи в модель
- Функция начальной инициализации (загрузки) модели
- Функция дообучения модели
- Функция для оценки результатов дообучения
- Функционал для взаимодействия пользователя с дообученной моделью



Список системных и пользовательских требований

Системные требования: в процессе обучения сети должно использоваться не более 16 ГБ GPU (объем памяти GPU P100)

Пользовательские требования: пользователь должен иметь возможность ввести контекст с вопросом, затем получить ответ



Инструментальные средства реализации

В качестве языка программирования был выбран Python из-за обилия библиотек, связанных с машинным обучением, и в частности NLP.

Модели брались из huggingface.co/models. Выбор обусловлен большим количеством доступных для дообучения моделей и возможностью загрузки напрямую с сервера

Сами вычисления производились в Kaggle notebook. Выбор обусловлен наличием возможности использования GPU P100



Метрики

Predicted class		
Positive	Negative	
TP	FN	Positive
FP	TN	Negative

True class

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

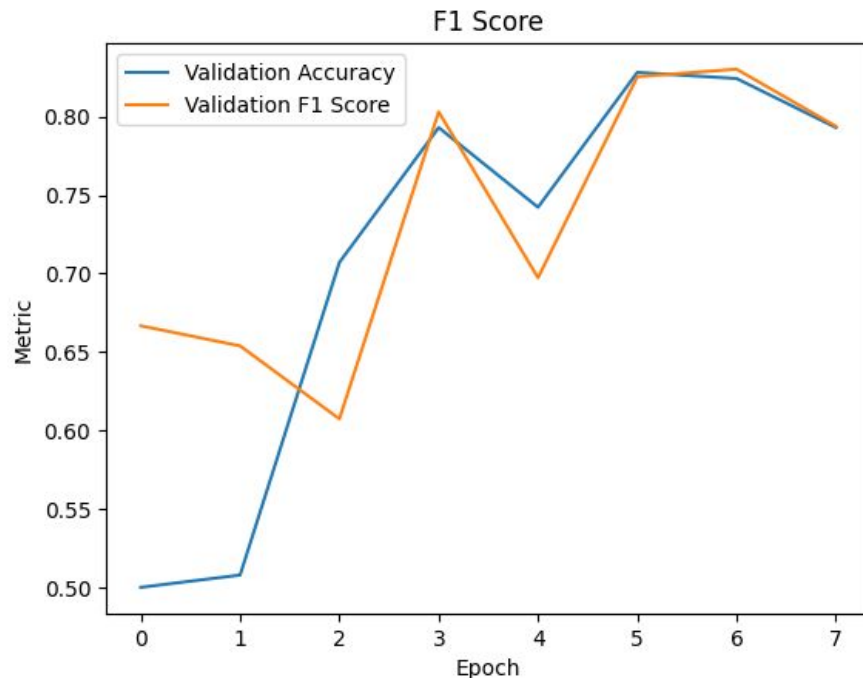
$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F_1 = \frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$$

Выбор модели

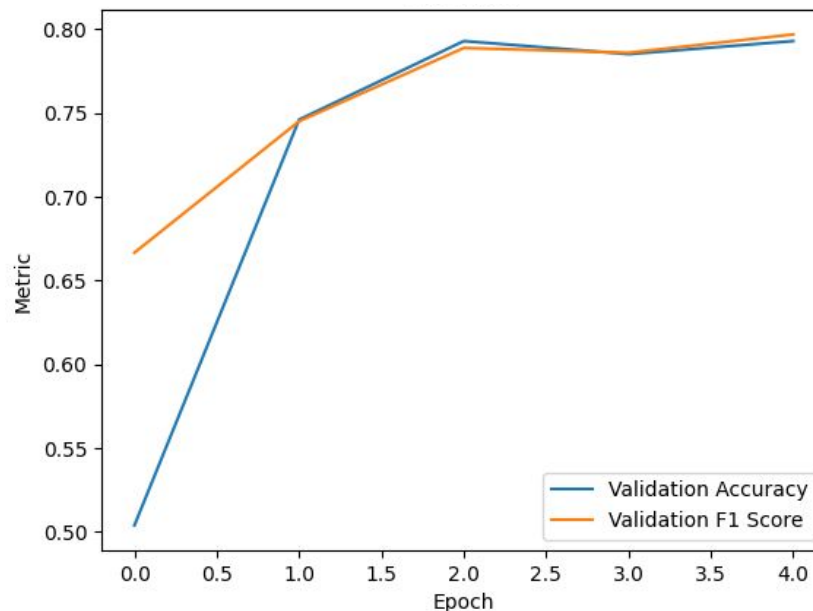
Была выбрана модель xlm-roberta-large-sag, потому что она показала лучшие значения метрик

Модель	accuracy	F1
RuBioRoBERTa	0.79	0.80
rubert-base-cased-mcn	0.59	0.59
xlm-roberta-large-sag	0.83	0.83
RuBioRoBERTa (команда SAI junior)	0.77	



Сравнение с известными аналогами

Команда SberAI
использовала модель
RuBioRoBERTa и дообучала
на тех же данных. accuracy
их модели равно 0,77. Я
тоже экспериментировал с
данной моделью сумел
немного улучшить
результат: accuracy моей
модели равно 0,79



010101
000100
100100
101010
101010
100000
010101
010111
111010
101010
101111
111111
100100
100101
010001
010010
101010
010100
101010
010100
100000
010000
111001
001001
000010
100001
011111
111111
101010
101001
001011
1

Пользовательские возможности

Пользователь может ввести вопрос с контекстом и получить ответ

Введите контекст: Корь крайне опасная для здоровья и заразная болезнь
Введите вопрос: Карантин из-за кори является крайне излишней мерой?
Предсказание: нет



Заключение

В ходе работы по теме “Разработка и сравнительный анализ моделей машинного обучения для ответа на связанные с медициной контекстно-зависимые вопросы” были сделаны:

1. Обзор литературы, содержащей описание различных способов для решения задачи классификации текстов. После сравнительного анализа было принято решение использовать BERT-подобные модели.
2. Работа по подбору гиперпараметров для модели RuBioRoBERTa. При скорости обучения равной 10^{-5} были достигнуты следующие результаты: accuracy = 0,79, F1 = 0,80.
3. Сравнение различных BERT моделей. Лучший результат достигнут при использовании модели xlm-roberta-large-sag: accuracy = 0,83, F1 = 0,83.



Список литературы

1. Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean / “Efficient Estimation of Word Representations in Vector Space”
2. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin / “Attention Is All You Need” 2017
3. Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova / “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”, 2018



010101
000100
100100
101010
101010
100000
010101
010111
111010
101010
101111
111111
100100
100101
010001
010010
101010
010100
101010
010100
100000
010000
111001
001001
000010
100001
011111
111111
101010
101001
001011
1

Спасибо за внимание!



НИИИВ

www.mephi.ru

Кафедра №22

www.kaf22.ru

