

# Предсказание авторства книг в жанре фэнтези

Иванов И. Кузьмин П. Шлюпенков С.

# В чем идея?

В данном исследовании представлена модель, которая классифицирует данные по авторам в жанре фэнтези.

Для нашего исследования мы взяли 5 из самых популярных произведений в жанре фэнтези: Дж. Р. Р. Толкин "Властелин колец", А. Сапковский "Последнее желание", Дж. Мартин "Игра престолов", Р. Желязны "Хроники Амбера" и Н. Гейман "Американские боги".

Мы создадим датасет, в котором у нас будут содержаться данные для решения задачи классификации предложений по авторству из исходных текстов.

# Цели и задачи

**Цель** работы заключается в создании модели-классификатора указания авторства в книгах жанра фэнтези.

## **Задачи:**

1. Собрать необходимые данные
2. Провести препроцессинг данных и подготовить их к анализу
3. Создать модель-классификатор
4. Провести анализ полученных результатов

# Этапы работы

Исследование поделено на две части:

В первой части мы собираем и подготавливаем данные для последующей загрузки в модель.

Во второй части мы проводим необходимый анализ текста для выполнения задачи классификации. Вторая часть состоит из трех независимых друг от друга частей. В каждой используются свои решения задачи классификации.

```

out_data = pd.DataFrame()
out_data['text'] = combined
out_data['author'] = labels
print(out_data.head())
print(out_data.tail())

```

	text	author
0	Тысячи лет пираты с Железных Островов – железн...	Martin
1	И на сей раз на него нахлынул приступ ярости, ...	Gaiman
2	И ответа не было.	Gaiman
3	Любопытство мое было подогрето до такой степен...	Tolkien
4	Пламя лизало камень злыми красными язычками.	Martin
	text	author
49995	Потом я разглядел, что у нее в левой руке, и б...	Zelazny
49996	Знал, наверное, еще как знал!	Tolkien
49997	Он вытер руки о джинсы и протянул Тени могучую...	Gaiman
49998	А потом сказала с неуверенной улыбкой:..	Gaiman
49999	Кто-то заметил на дереве краснохвостого ястреб...	Gaiman

+ Code

+ Markdown

Получившийся необработанный датасет

```
data['text'] = normed_text
data.to_csv('preprocessed_data.csv', index=False)
print(data.head())
print(data.tail())
```

	text	author
0	тысячи лет пираты с железных островов железны...	Martin
1	и на сей раз на него нахлынул приступ ярости г...	Gaiman
2	и ответа не было	Gaiman
3	любопытство мое было подогрето до такой степен...	Tolkien
4	пламя лизало камень злыми красными язычками	Martin
	text	author
49995	потом я разглядел что у нее в левой руке и быс...	Zelazny
49996	знал наверное еще как знал	Tolkien
49997	он вытер руки о джинсы и протянул тени могучую...	Gaiman
49998	а потом сказала с неуверенной улыбкой	Gaiman
49999	кто-то заметил на дереве краснохвостого ястреба...	Gaiman

[+ Code](#)[+ Markdown](#)

Получившийся обработанный датасет

# **БЛОК 1**

## **Сбор и подготовка данных**

;

```
authors = Counter(author)
authors
```

```
Counter({'Martin': 10000,
        'Gaiman': 10000,
        'Tolkien': 10000,
        'Zelazny': 10000,
        'Sapkowski': 10000})
```

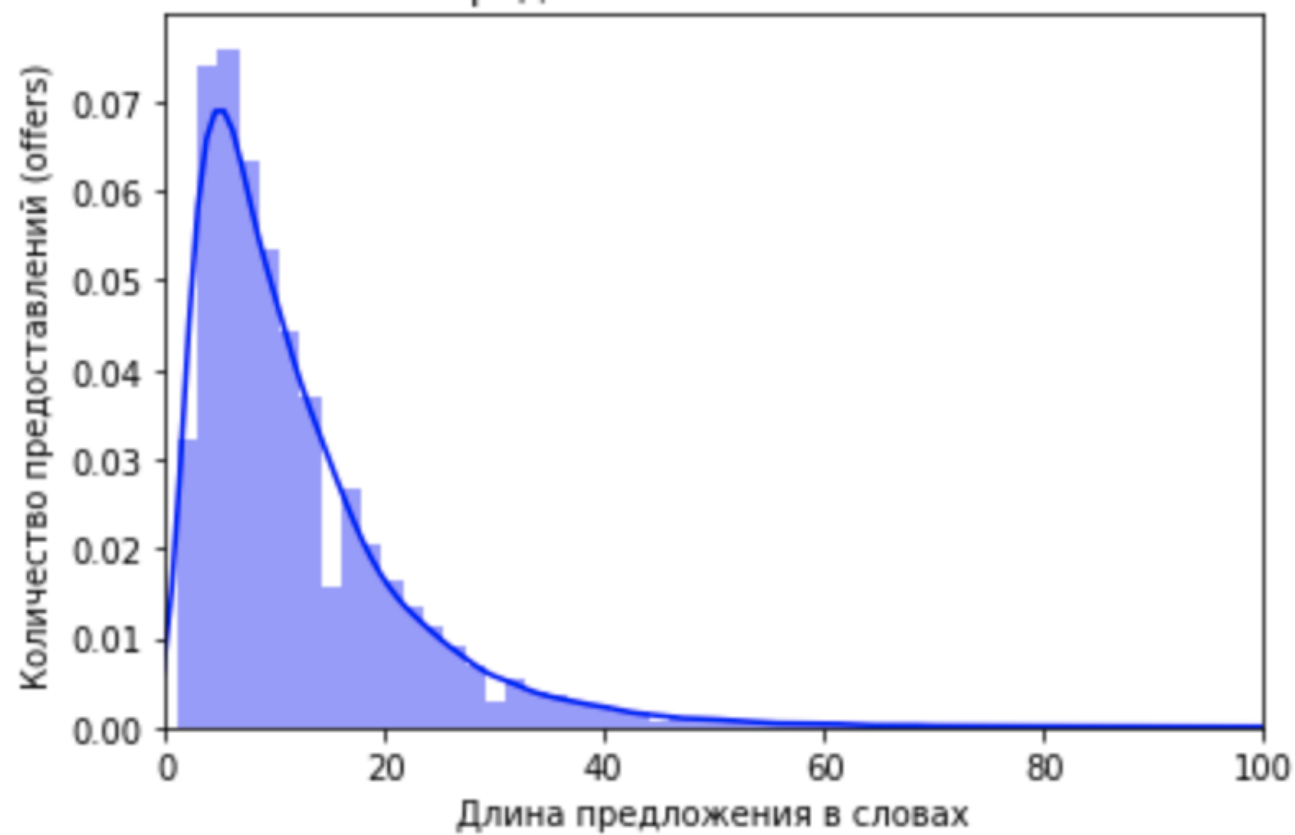
В нашем датасете мы равномерно распределили  
количество предложений по авторам



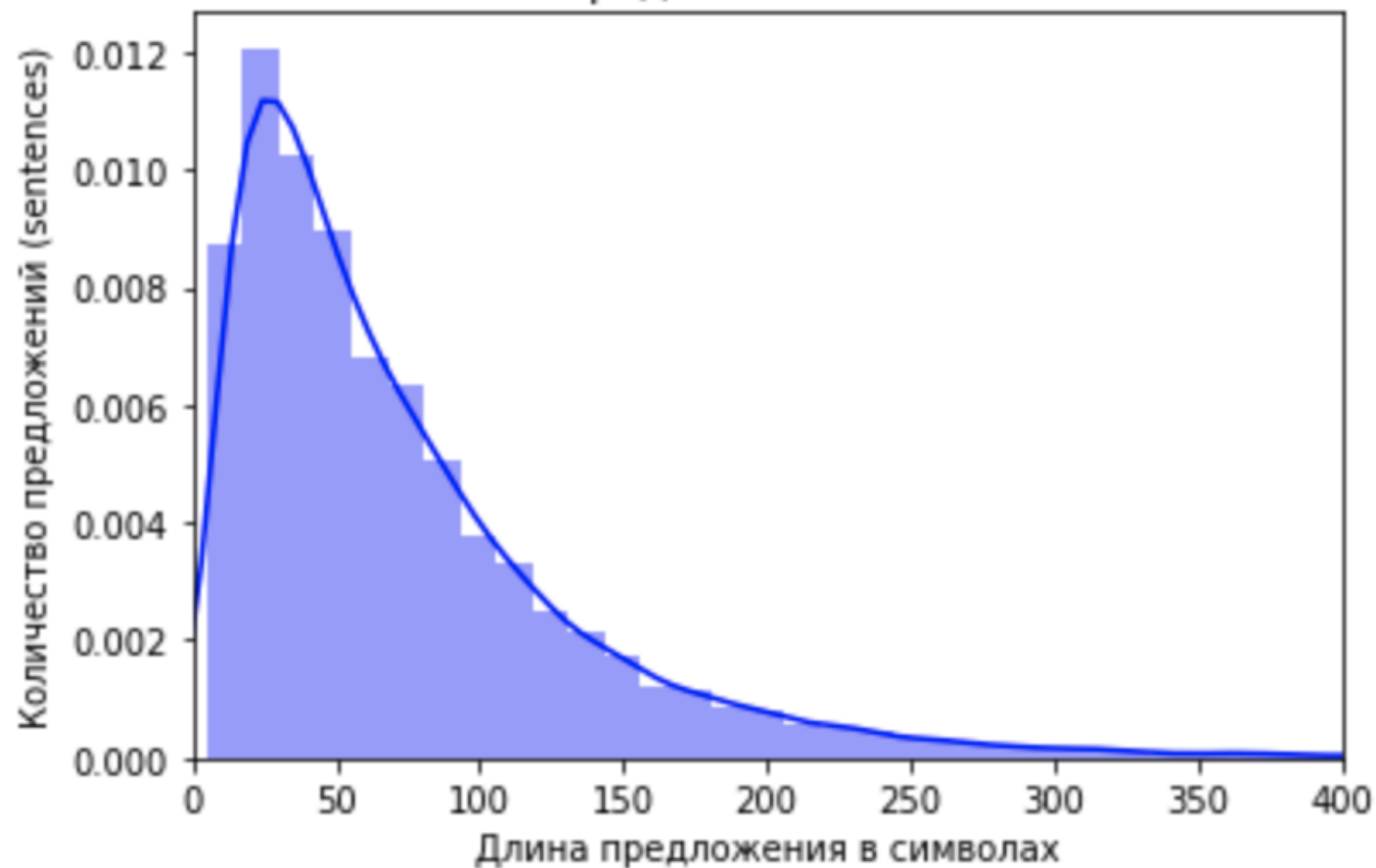
```
word_count = np.array([len(sent.split()) for sent in text])  
char_count = np.array([len(sent) for sent in text])  
ave_length = char_count / word_count
```

Пока мы не обработали данные для работы с ними в  
способах 1 и 2 посчитаем базовую статистику

Распределение количества слов



Распределение символов



Распределение средней длины слова



## **БЛОК 2**

**Анализ текста и выполнение  
задачи классификации**

**Способ 1**

**N-граммы: сложный и долгий  
метод**

---

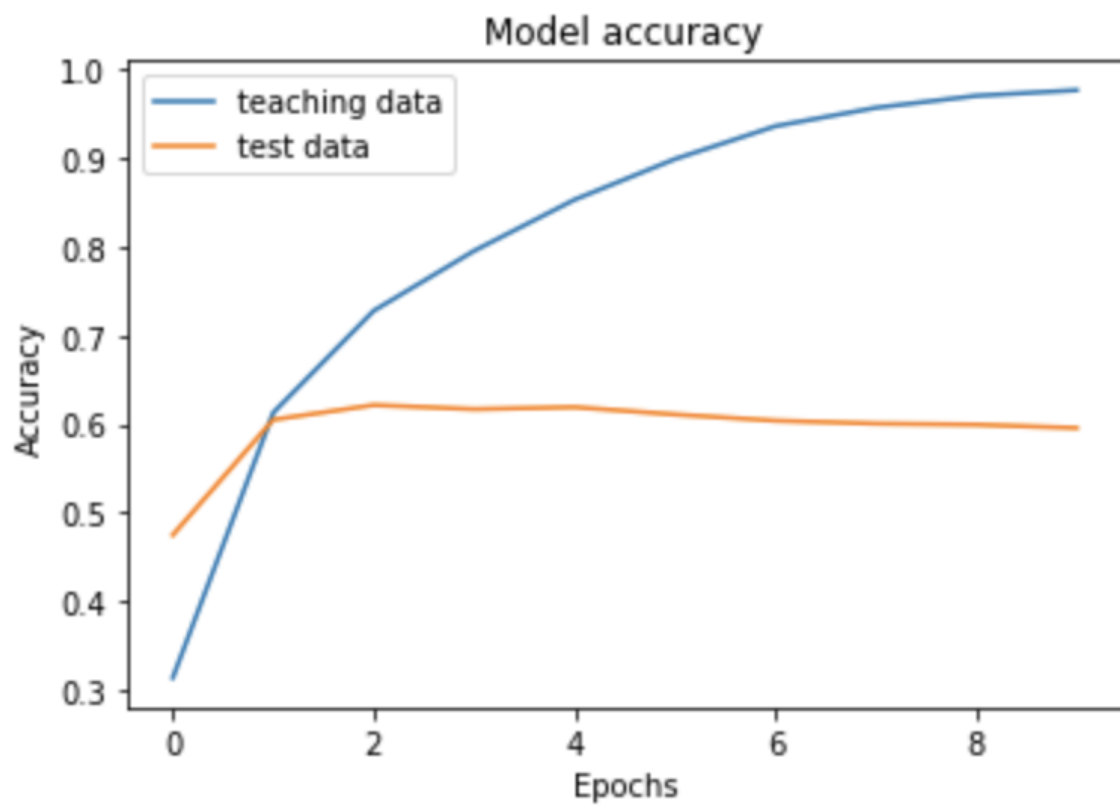
Accuracy: 0.5949  
Average Precision: 0.5960417779145594  
Average Recall: 0.5949  
Average F1 Score: 0.5939566255094416  
Learning time: 577.1681146621704 секунд  
Prediction time: 272.6482763290405 секунд  
Confusion matrix:  
[[1359 121 187 153 191]  
[ 200 1004 230 312 260]  
[ 168 164 1186 224 229]  
[ 138 163 208 1267 224]  
[ 184 186 259 250 1133]]

По итогу расчета наших n-грамм (1-3) мы получаем  
первую модель и следующие метрики

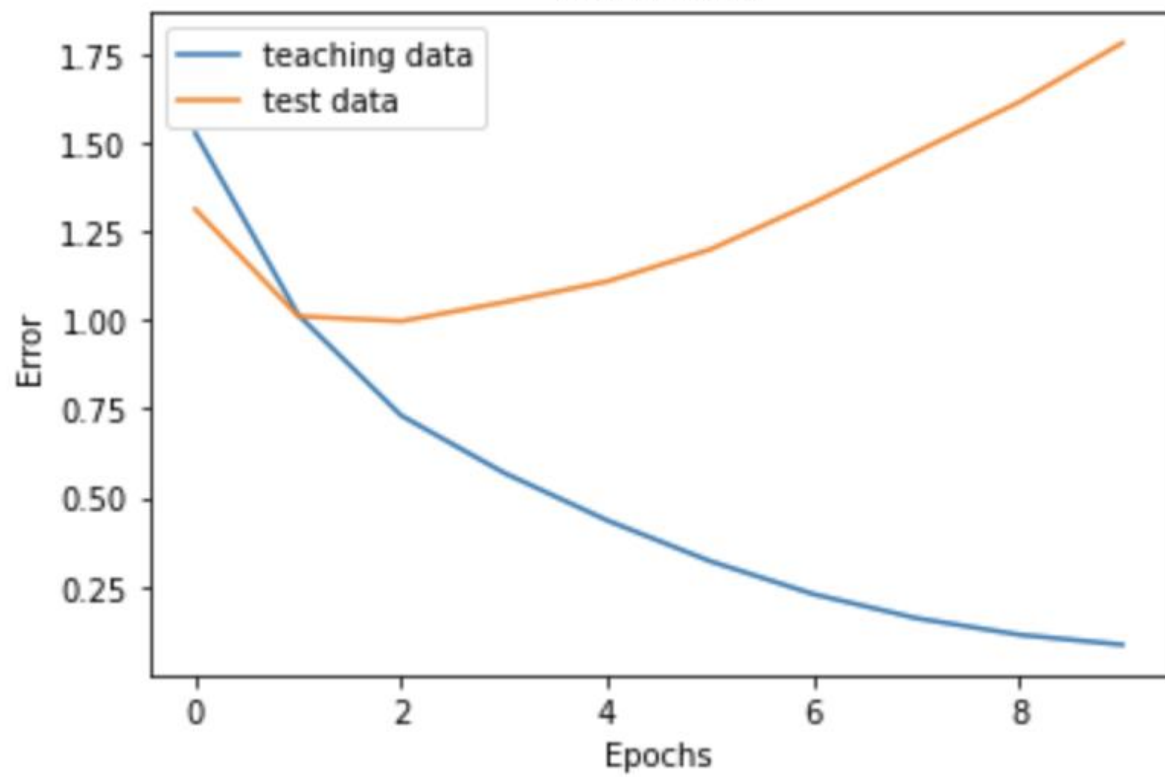
Матрица ошибок Модель 1







Model error

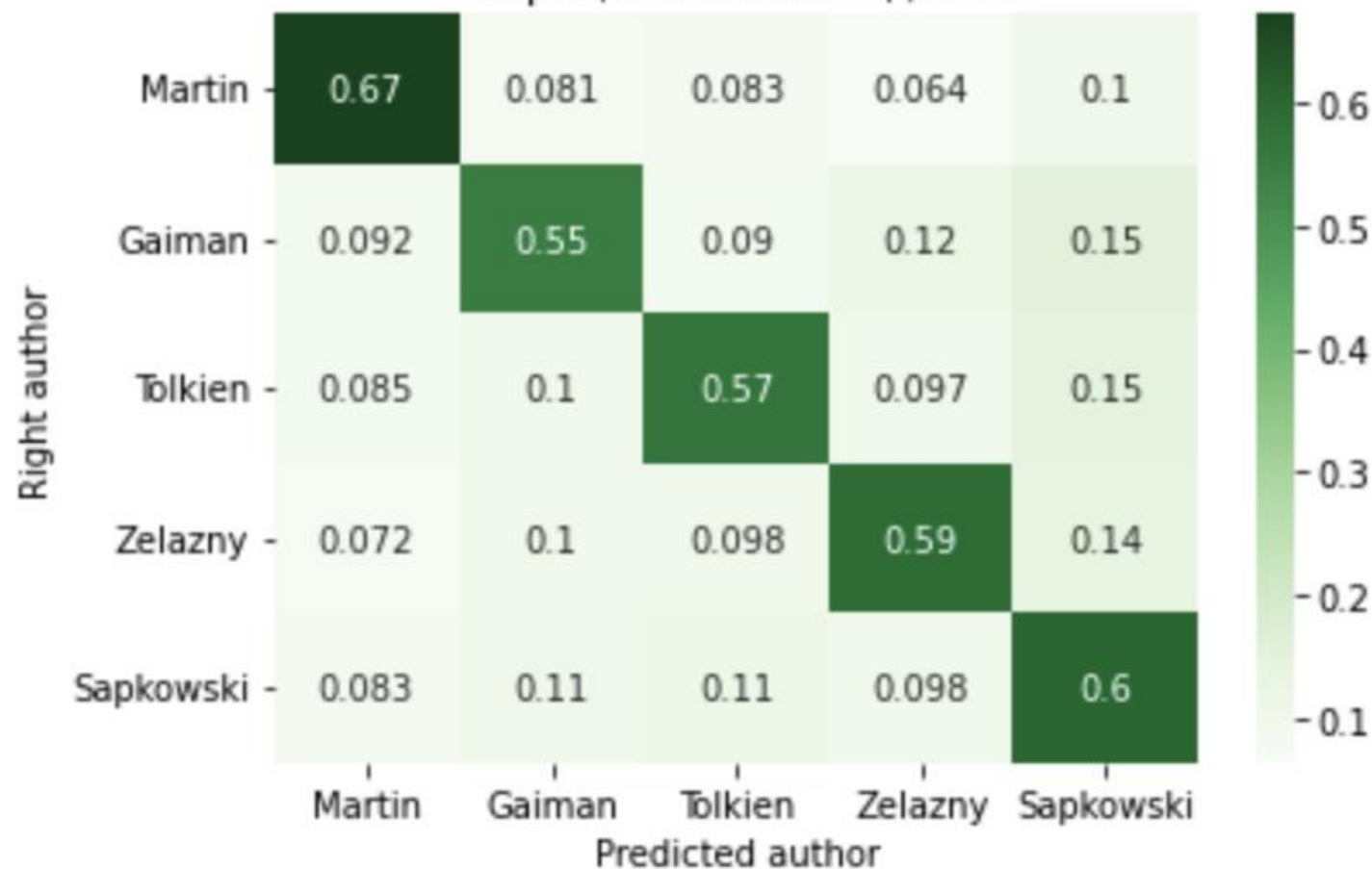


---

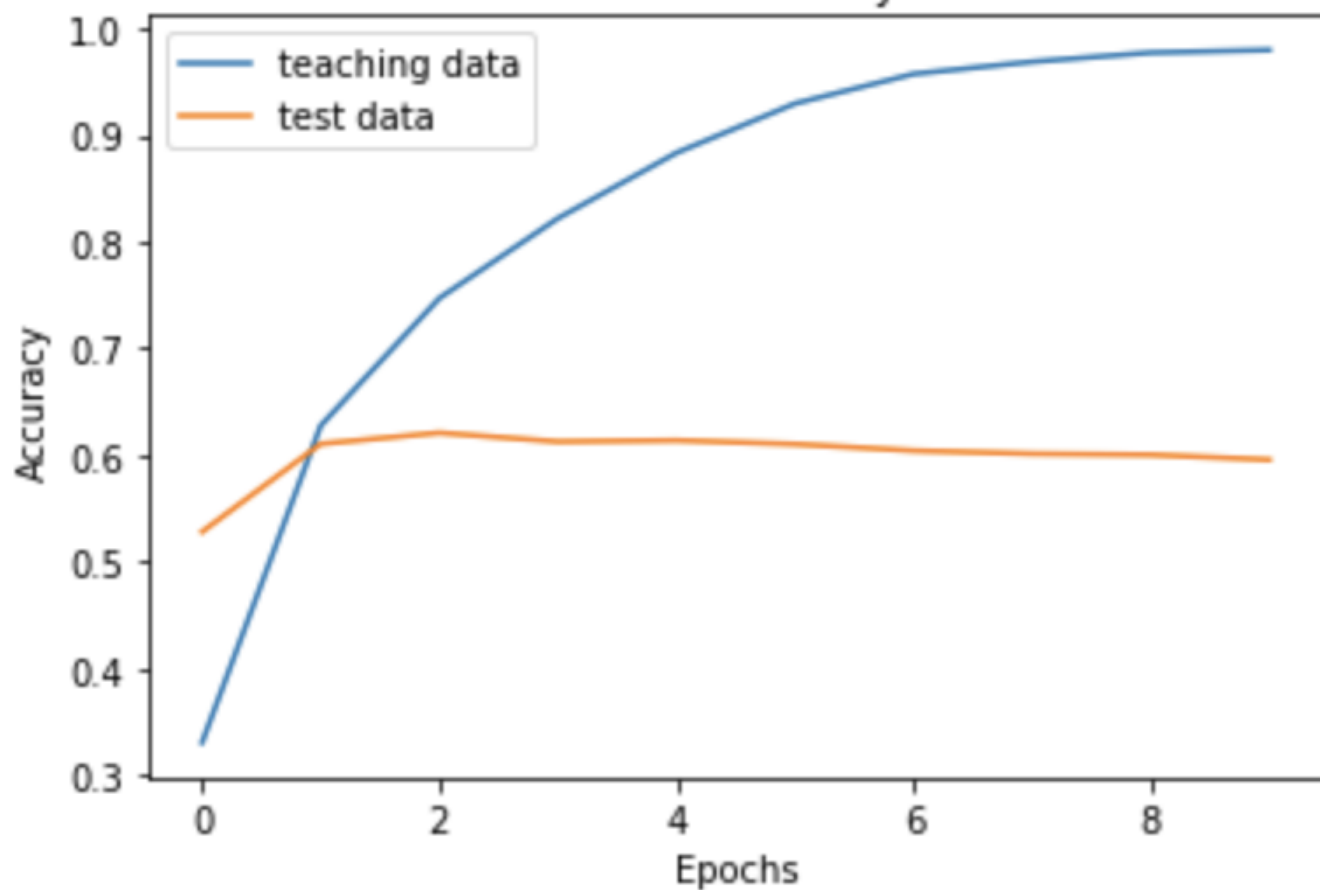
Accuracy: 0.5975  
Average Precision: 0.5992085171792201  
Average Recall: 0.5974999999999999  
Average F1 Score: 0.5976984427392338  
Learning time: 806.9759614467621 секунд  
Predict time: 376.9495553970337 секунд  
Confusion matrix:  
[[1350 162 167 129 203]  
[ 185 1103 180 233 305]  
[ 168 199 1126 192 286]  
[ 144 202 196 1183 275]  
[ 166 212 224 197 1213]]

После апгрейда модели мы получаем вторую модель и  
следующие метрики

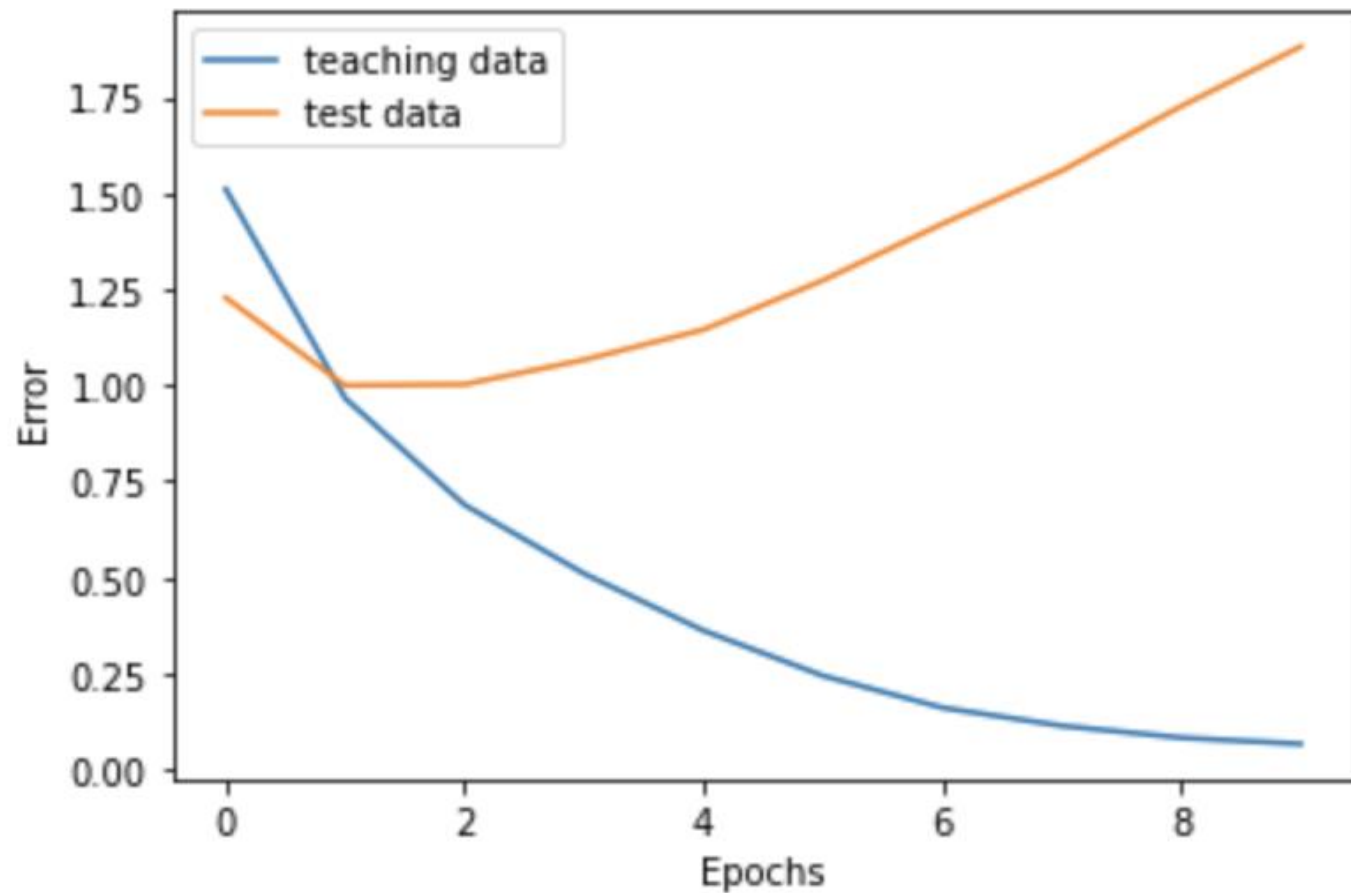
Матрица ошибок Модель 2



Model accuracy



Model error



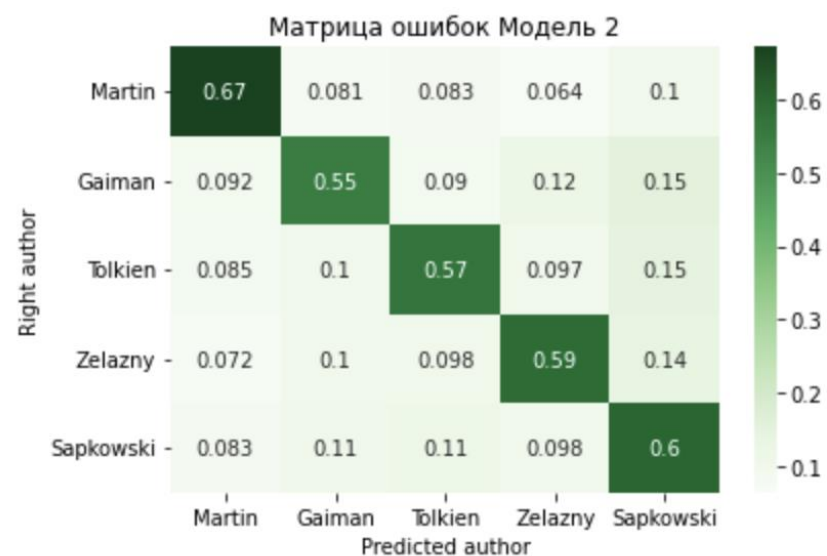
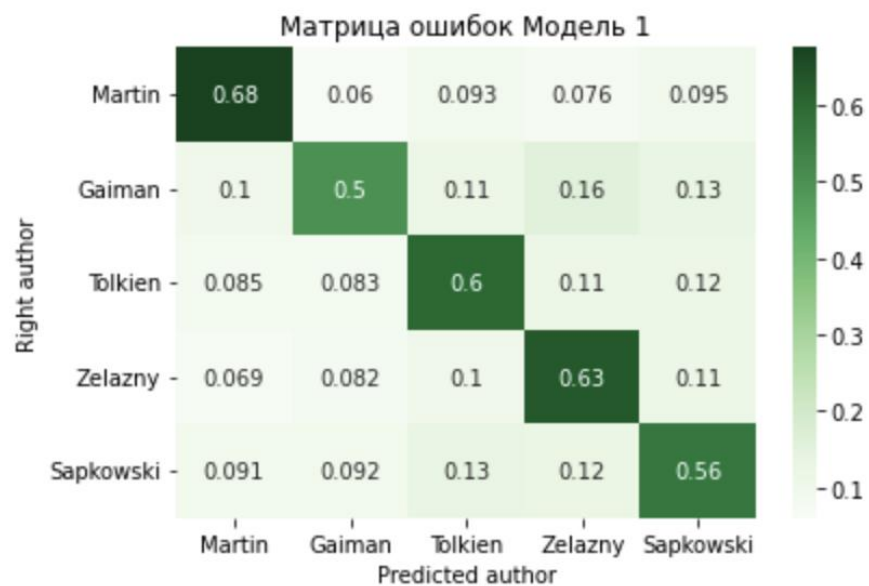
---

Accuracy: 0.5949  
Average Precision: 0.5960417779145594  
Average Recall: 0.5949  
Average F1 Score: 0.5939566255094416  
Learning time: 577.1681146621704 секунд  
Prediction time: 272.6482763290405 секунд  
Confusion matrix:  
[[1359 121 187 153 191]  
[ 200 1004 230 312 260]  
[ 168 164 1186 224 229]  
[ 138 163 208 1267 224]  
[ 184 186 259 250 1133]]

---

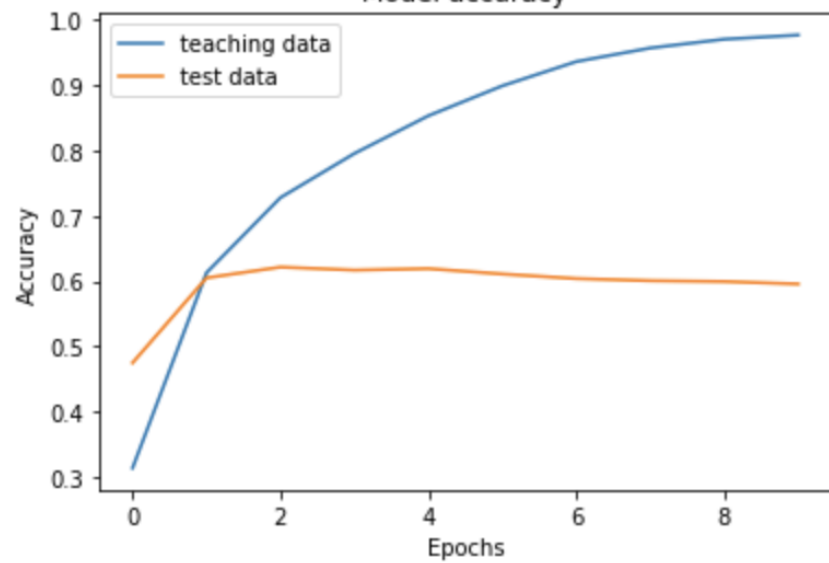
Accuracy: 0.5975  
Average Precision: 0.5992085171792201  
Average Recall: 0.5974999999999999  
Average F1 Score: 0.5976984427392338  
Learning time: 806.9759614467621 секунд  
Predict time: 376.9495553970337 секунд  
Confusion matrix:  
[[1350 162 167 129 203]  
[ 185 1103 180 233 305]  
[ 168 199 1126 192 286]  
[ 144 202 196 1183 275]  
[ 166 212 224 197 1213]]

Наглядно

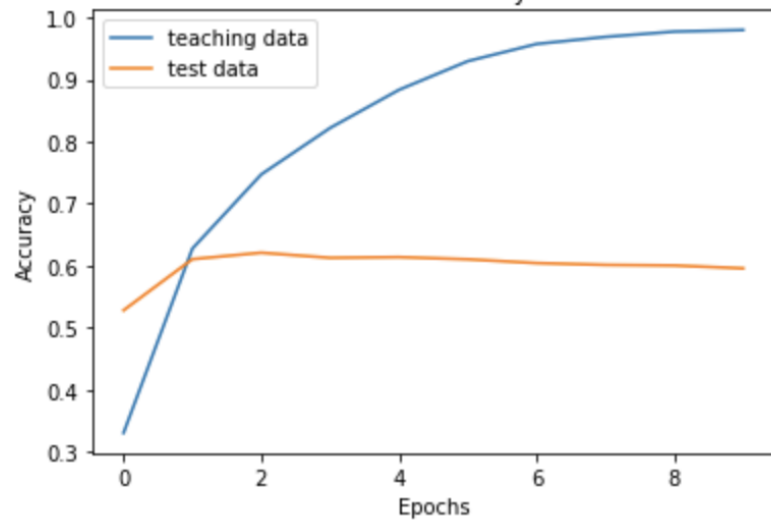




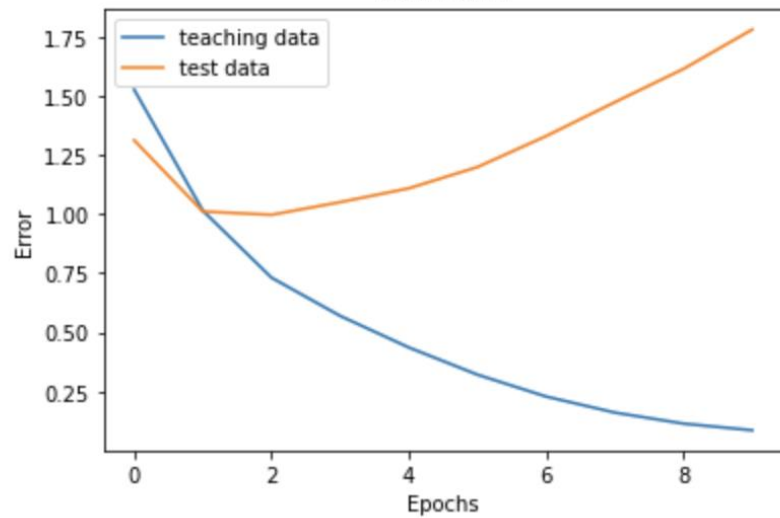
Model accuracy



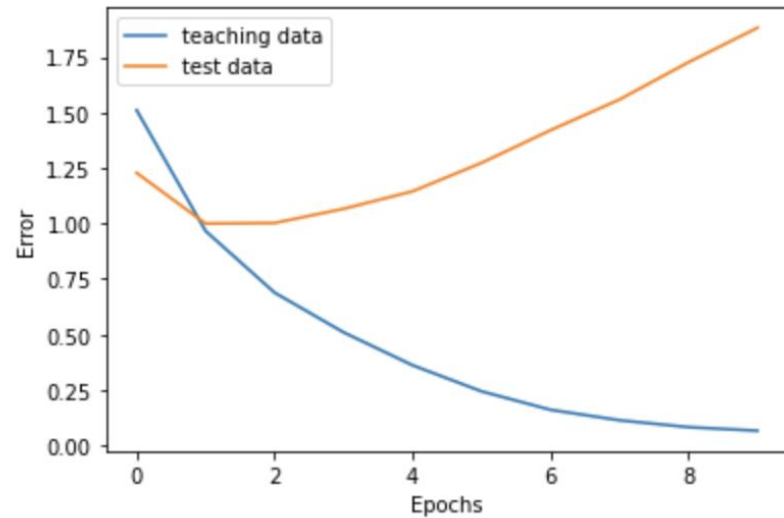
Model accuracy



Model error



Model error



```
for i in range(20):  
    print('Предложение', i, '- Верный ответ =', author_test[i], 'Предсказание Модели 1 =', author_pred1[i], 'Предсказание Модели 2 =', author_pred2[i])  
    print(text_test[i], '\n')
```

Предложение 0 – Верный ответ = Gaiman Предсказание Модели 1 = Sapkowski Предсказание Модели 2 = Sapkowski  
геродота

Предложение 1 – Верный ответ = Sapkowski Предсказание Модели 1 = Sapkowski Предсказание Модели 2 = Gaiman  
держи рот на замке и не встречай

Предложение 2 – Верный ответ = Martin Предсказание Модели 1 = Martin Предсказание Модели 2 = Martin  
с ним был сир марк пайпер они привели сына сира реймена дарри парнишку не старше брана

Предложение 3 – Верный ответ = Gaiman Предсказание Модели 1 = Tolkien Предсказание Модели 2 = Sapkowski  
я слышал чейто голос

Предложение 4 – Верный ответ = Sapkowski Предсказание Модели 1 = Gaiman Предсказание Модели 2 = Gaiman  
и не крутили задом у мужчин перед глазами

Предложение 5 – Верный ответ = Sapkowski Предсказание Модели 1 = Sapkowski Предсказание Модели 2 = Sapkowski  
тебе не уехать геральт

Предложение 6 – Верный ответ = Tolkien Предсказание Модели 1 = Tolkien Предсказание Модели 2 = Tolkien  
я сразу же покинул дзэнтора но еще по дороге на север получил весть из лориэна арагорна побывал там и просил передать что он отыскал существо именуемое голлумом

Предложение 7 – Верный ответ = Sapkowski Предсказание Модели 1 = Sapkowski Предсказание Модели 2 = Sapkowski  
ни одного родственника или жениха дочки

Предложение 8 – Верный ответ = Zelazny Предсказание Модели 1 = Gaiman Предсказание Модели 2 = Gaiman  
в каком смысле

Предложение 9 – Верный ответ = Zelazny Предсказание Модели 1 = Sapkowski Предсказание Модели 2 = Sapkowski  
для этого мне пришлось вызвать логрус который упал между нами словно нож гильотины и дернул меня так словно я прикоснулся к оголенному проводу под напряжением

Предложение 10 – Верный ответ = Gaiman Предсказание Модели 1 = Martin Предсказание Модели 2 = Martin  
я хочу с ней попрощаться

# Способ 2

Стемминг и векторные  
представления: логистическая  
регрессия

		text	author
0	тысяч лет пират железн остров железн людъм наз...		Martin
1	се нахлынул приступ ярост глухо темны собра ку...		Gaiman
2		ответ	Gaiman
3	любопытств мо подогрет тако степен последова н...		Tolkien
4	плам лиза камен злым красн язычк		Martin
		text	author
49995	разглядел лево рук быстр посмотрел джул		Zelazny
49996		знал наверн знал	Tolkien
49997	вытер рук джинс протянул тен могуч лапищ слыша...		Gaiman
49998		сказа неуверенно улыбка	Gaiman
49999	замет дерев краснохв ястреб сказа эт скор соко...		Gaiman

---

(42500, 246287) (7500, 246287)

/opt/conda/lib/python3.7/site-packages/sklearn/linear\_model/\_logistic.py:818: ConvergenceWarning: lbfgs failed to converge (status=1):  
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

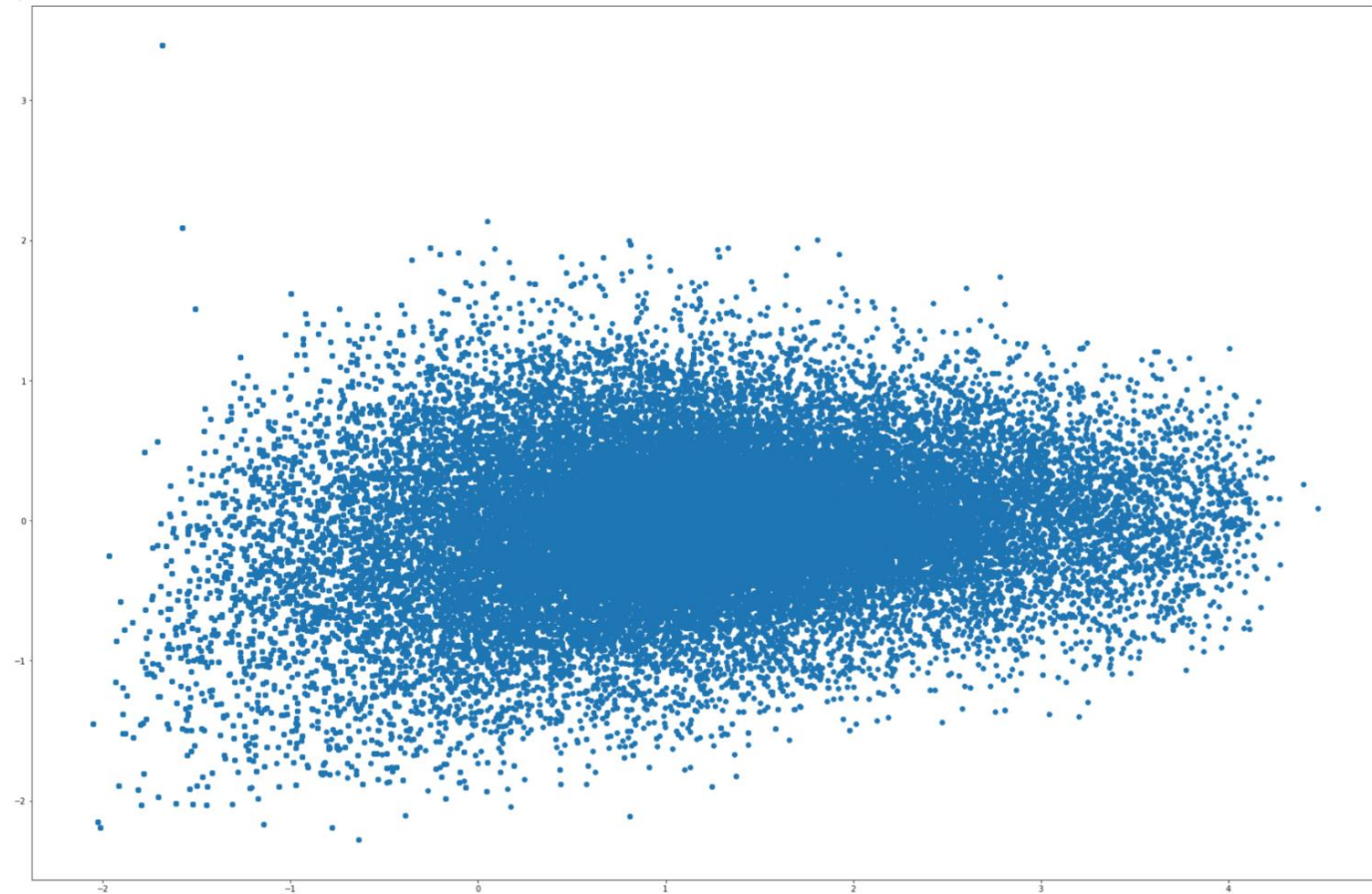
extra\_warning\_msg=\_LOGISTIC\_SOLVER\_CONVERGENCE\_MSG,

0.6726666666666666

## Способ 3

Визуализация эмбеддингов для  
каждой книги

Время начала: 2023-01-22 21:00:21.098808

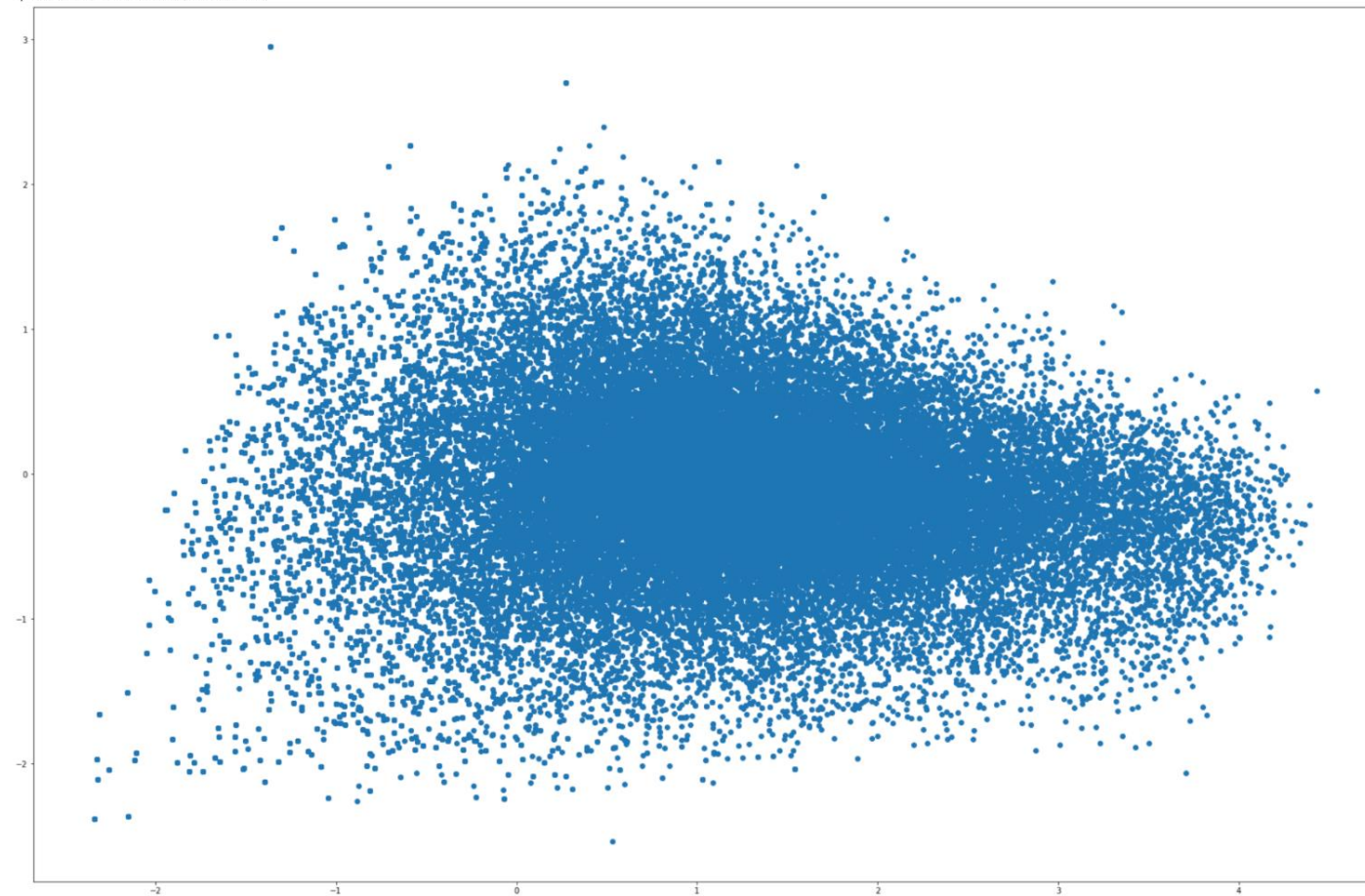


Время начала: 2023-01-22 21:00:21.098808  
Время окончания: 2023-01-22 21:09:20.000910

Считаем эмбединги для "Американские боги" Н.  
Геймана.



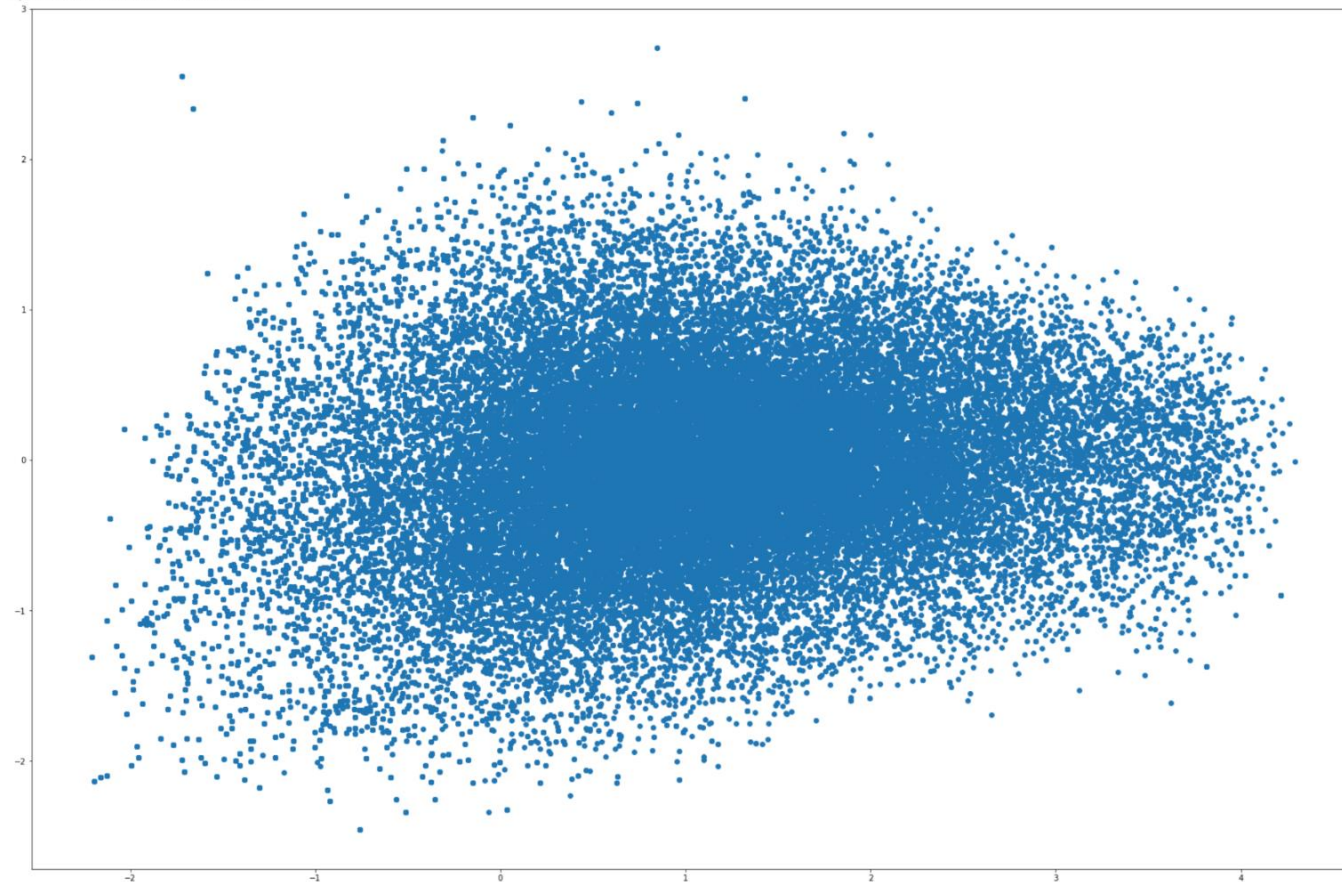
Время начала: 2023-01-22 21:11:36.573349



Время начала: 2023-01-22 21:11:36.573349  
Время окончания: 2023-01-22 21:26:33.176794

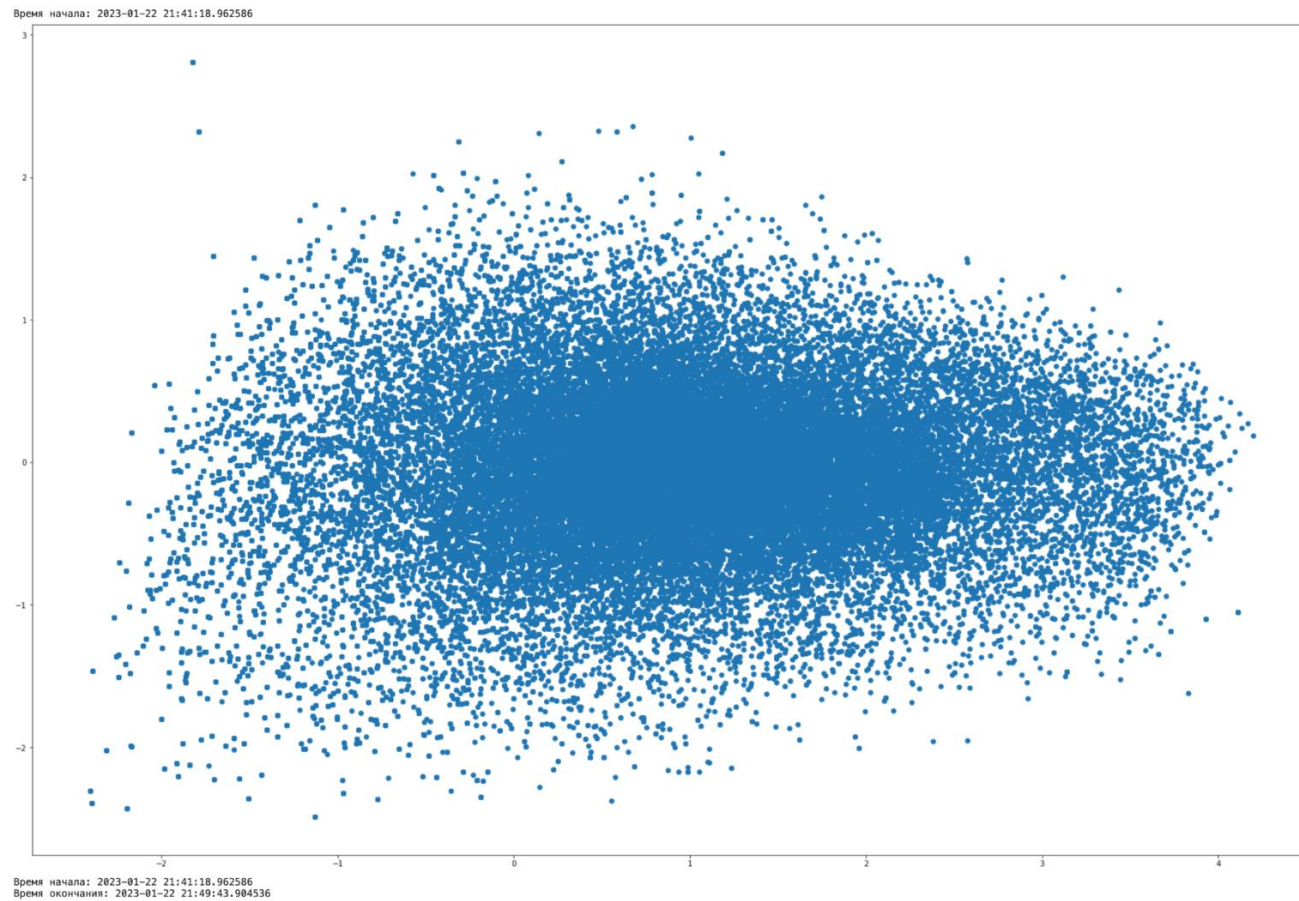
Считаем эмбединги для "Хроники Амбера" Р. Желязны.

Время начала: 2023-01-22 21:27:49.703321



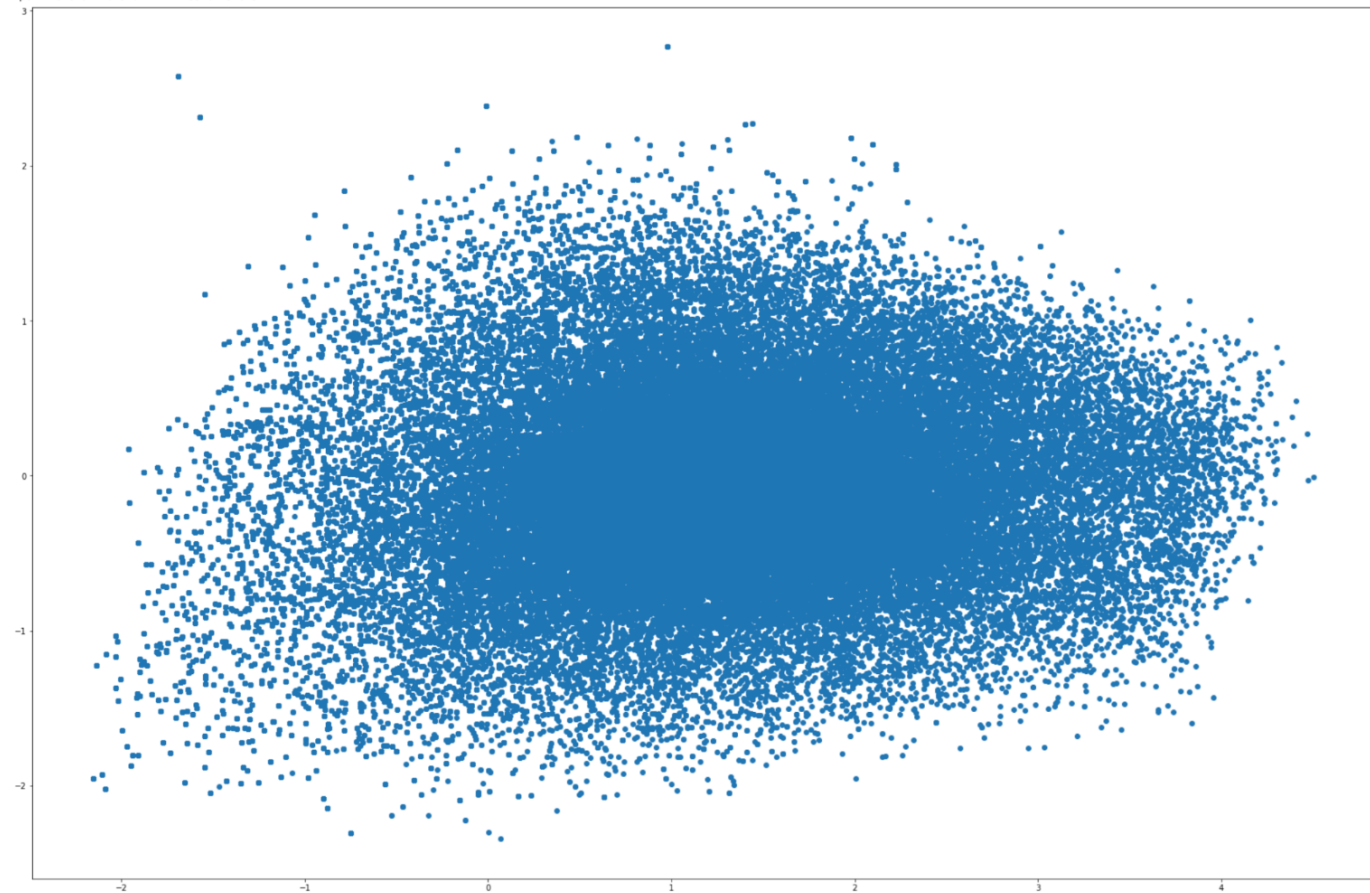
Время начала: 2023-01-22 21:27:49.703321  
Время окончания: 2023-01-22 21:40:23.490132

Считаем эмбединги для "Игра престолов" Дж. Мартина.



Считаем эмбединги для "Последнее желание. Меч предназначения" А.  
Сапковского.

Время начала: 2023-01-22 21:55:02.520434



Время начала: 2023-01-22 21:55:02.520434  
Время окончания: 2023-01-22 22:15:20.848591

Считаем эмбединги для "Властелин колец" Дж. Р. Р. Толкина.