

Федеральное государственное бюджетное
образовательное учреждение высшего образования
«Санкт-Петербургский национальный исследовательский
Академический университет Российской академии наук»
Центр высшего образования

Кафедра математических и информационных технологий

Подгузов Никита Владимирович

Определение наилучшего ответа на StackOverflow

Бакалаврская работа

Допущена к защите.
Зав. кафедрой:
д. ф.-м. н., профессор Омельченко А. В.

Научный руководитель:
Курбанов Р. Э.

Рецензент:
Артамонов А. С.

Санкт-Петербург
2018

SAINT-PETERSBURG ACADEMIC UNIVERSITY
Higher education centre

Department of Mathematics and Information Technology

Nikita Podguzov

Best answer detection on StackOverflow

Graduation Thesis

Admitted for defence.
Head of the chair:
professor Alexander Omelchenko

Scientific supervisor:
Rauf Kurbanov

Reviewer:
Alexey Artamonov

Saint-Petersburg
2018

Оглавление

| | |
|---|-----------|
| Введение | 5 |
| 1. Обзор предметной области | 7 |
| 1.1. Обзор существующих решений | 7 |
| 1.2. Векторное представление слов | 8 |
| 1.2.1. Bag of Words | 8 |
| 1.2.2. Word2Vec | 9 |
| 1.2.3. Fasttext | 10 |
| 1.3. Нейронные сети | 11 |
| 1.3.1. Рекуррентная нейронная сеть | 11 |
| 1.3.2. Сверточная нейронная сеть | 12 |
| 2. Цель и задачи | 14 |
| 3. Описание данных | 15 |
| 3.1. Сбор данных | 15 |
| 3.2. Анализ данных | 15 |
| 3.3. Предобработка данных | 16 |
| 3.4. Постобработка данных | 17 |
| 4. Методы и реализация | 18 |
| 4.1. Извлечение признаков | 18 |
| 4.1.1. Векторное представление вопроса и ответа | 18 |
| 4.1.2. Лингвистические признаки | 19 |
| 4.1.3. Словарные признаки | 21 |
| 4.1.4. Метаинформация | 22 |
| 4.1.5. Нормализация | 22 |
| 4.2. Методы | 23 |
| 4.3. Реализация модели | 24 |
| 4.3.1. Сверточная нейронная сеть | 24 |
| 4.3.2. Рекуррентная нейронная сеть | 24 |
| 4.3.3. Детали реализации | 25 |
| 4.3.4. Репозиторий | 25 |
| 5. Тестирование | 27 |
| 5.1. Используемые метрики | 27 |
| 5.2. Производительность | 27 |
| 5.3. Результаты | 28 |
| 5.4. Анализ результатов | 29 |

| | |
|--------------------------|-----------|
| Заключение | 31 |
| Список литературы | 33 |

Введение

С распространением и развитием интернета стали популярны системы вопросов и ответов — вид сайтов, позволяющий пользователям задавать вопросы и отвечать на вопросы, заданные другими пользователями. Среди наиболее популярных подобных систем можно выделить многопрофильные (такие, как *Yahoo Answers*, *Quora*) и узкоспециализированные (например, *Stackoverflow* [31]). В рамках данной работы мы сконцентрируемся на системе StackOverflow, которая в первую очередь специализируется на вопросах по программированию и смежным областям.

На сайтах системы Stackoverflow у автора есть возможность после того, как он задал интересующий его вопрос, отметить ответ, решивший его проблему, как правильный.

Несмотря на это, достаточно большая часть вопросов остается без подобной отметки даже при наличии хороших ответов-кандидатов, что может усложнить поиск нужной информации по данному вопросу другим пользователями. Кроме того, пользователю, задавшему вопрос по какой-то не самой популярной теме, может быть сложно оценить правильность и релевантность полученных ответов.

Таким образом перед нами появляется задача определения правильного ответа для вопросов с Stackoverflow.

Мотивация

В связи с многократно возросшими вычислительными мощностями в последнее время в области машинного обучения приобрели популярность нейронные сети. Сфера применимости моделей, использующих нейронные сети, очень широка: от распознавания образов на изображениях до предсказания курсов валют на бирже.

Несмотря на наличие работ, исследующих задачу предсказания лучшего ответа для систем вопросов и ответов, в основном в них используются более классические методы машинного обучения, хотя нейронные сети и показывают достойные результаты в задачах классификации текстов.

В данной работе исследуется применимость нескольких моделей, основанных на нейронных сетях и использующих такие особенности сайта StackOverflow, как технический язык вопросов и наличие сниппетов кода в вопросах и ответах, к поставленной задаче.

Кратко о последующих главах

В главе 1 производится обзор имеющихся решений и исследований по данной теме, их преимуществ и недостатков.

В главе 2 описывается формальная постановка задачи, а также рассматриваются ее особенности и ставятся основные цели и задачи работы.

В третьей главе описана работа по предварительной обработке и очистке данных с сайта StackOverflow, а также их анализ.

Глава 4 содержит описание основных методов и подходов к поставленной задаче, а также детали их реализаций.

В последней главе приведены результаты экспериментов с различными конфигурациями реализованных моделей, а также проведен анализ полученных результатов и сделаны выводы о применимости модели к задаче определения правильного ответа.

1. Обзор предметной области

1.1. Обзор существующих решений

Существует множество работ, исследующих различные аспекты систем вопросов и ответов: определение наилучшего ответа, определение качества ответа, выбор ответов к новым вопросам из уже существующих, определение качества вопроса, а также вероятности получения ответа на него.

Тем не менее, в рамках темы нашей работы стоит обратить внимание на исследования, касающиеся первого обозначенного аспекта: определения наилучшего ответа. Все статьи, фокусирующиеся на этой теме, имеют общую структуру: они извлекают различные признаки из текстов ответов, а также метаинформации, а затем используют некоторый бинарный классификатор, обученный на размеченных данных.

Были подробно рассмотрены особенности следующих работ:

- В одной из первых работ [32], анализирующих данную проблему в разрезе сайта StackOverflow, использовался классификатор "случайный лес" по признакам трех типов: соответствующих содержанию ответа, измеряющих похожесть вопроса и ответа, а также использующих метаинформацию.
- Также достаточно важна работа [15], в которой, во-первых, проведено масштабное исследование 21 различного сайта системы StackExchange, а во-вторых представлена так называемая идея дискретизации признаков: процесс получения нового признаком путем сортировки одного из старых и замены его значений на порядковый номер в получившемся списке (эта идея будет раскрыта более подробно далее).
- Наконец, нужно упомянуть исследование [5], проводившееся на трех датасетах сайтов системы StackExchange. Подход, представленный в работе, использует большое количество различных признаков, добиваясь наилучших результатов с помощью такой метаинформации, как рейтинги ответа или пользователя, отвечавшего на вопрос, которые не использовались в нашей работе. Тем не менее, авторы предоставляют полученные метрики экспериментов для различных подмножеств признаков, поэтому мы все равно можем сравнить полученные нами результаты с представленными в статье. В качестве классификатора авторы использовали чередующиеся решающие деревья.

Проанализировав имеющиеся исследования, можно отметить некоторые недостатки существующих решений:

- Лишь одна из рассмотренных работ использует текст вопроса, хотя это, безусловно, важная составляющая при определении правильности ответа на вопрос.

- Те работы, в которых анализируется содержание текста ответа, во-первых, никак не учитывают контекст слова в тексте, а во-вторых, не пытаются работать с семантикой языка, в частности, не имеют возможности видеть близкие по смыслу слова.
- Некоторые работы также не используют специфику сайта StackOverflow, вопросы на котором написаны техническим языком и с достаточно большим количеством терминов, а также могут содержать фрагменты исходного кода.

1.2. Векторное представление слов

При работе с текстами в задачах обработки естественного языка требуется представлять слова и тексты в виде векторов, чтобы с ними было удобнее работать в дальнейшем. Существует несколько подходов к векторному представлению слов.

1.2.1. Bag of Words

Bag of Words [13] — это одно из самых простых представлений, первые упоминания о котором датируются еще 1950-ми годами.

Основная идея состоит в том, что каждый документ представляется в виде вектора, размерность которого равна размеру словаря, а i -я компонента равна 1, если i -ое слово из словаря присутствует в документе, и 0 иначе. При этом существуют различные модификации этого метода: могут не учитываться так называемые ”стоп-слова”, которые встречаются слишком часто и не несут в себе никакой информации о содержании документа; вместо бинарных значений компонент могут также использоваться количество вхождений слова в документ либо мера TF-IDF [17], когда значение i -й компоненты определяется по формуле:

$$w_i = tf_i \cdot \log \frac{N}{df_i} \quad (1)$$

где tf_i — количество вхождений слова в документ,

N — количество документов в нашем корпусе данных,

df_i — количество документов, в которых встречается данное слово.

Подобная модификация позволяет учитывать насколько то или иное слово часто встречается в других документах и давать больший вес более редким словам.

Пример применения базовой модели Bag of Words со стоп-словами вы можете увидеть на рисунке 1.

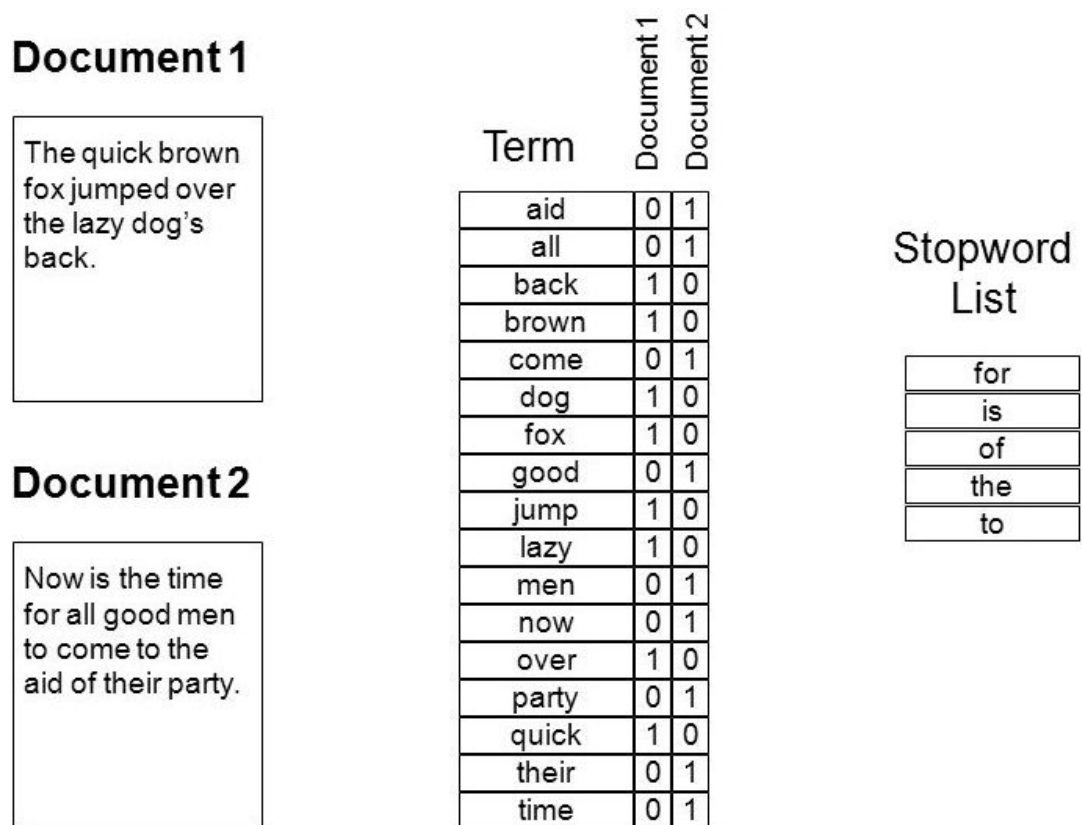


Рис. 1: Пример работы модели Bag of Words

1.2.2. Word2Vec

В 2013 году был представлен Word2Vec [8] — другой способ векторизации слов, основанный на нейронных сетях. Основная идея метода состоит в предположении о том, что слова, находящиеся в похожих контекстах, скорее всего будут значить похожие вещи, то есть будут семантически близкими. В данном случае под контекстом подразумевается набор слов из окна фиксированный ширины вокруг текущего слова. Тем не менее, контекст этот можно использовать для обучения двумя разными способами, поэтому в Word2Vec используется один из двух типов алгоритма:

- **Модель CBOW** (рисунок 2а)

В данной модели используется нейронная сеть с одним скрытым полносвязным слоем с линейной функцией активации. Обучение происходит следующим образом: по корпусу имеющихся текстов просматриваются все окна фиксированной ширины, на вход сети подается контекст, представленный в виде вектора размерности, равной длине словаря, полученный по модели мешка слов, описанной ранее, а на выходе ожидается вектор, в котором i -я компонента обозначает вероятность того, что в центре окна i -е слово. В качестве итогового векторного представления слов используется матрица весов скрытого слоя.

Таким образом, мы пытаемся по данному контексту научиться предсказывать слово.

- **Модель Skip-gram** (рисунок 2b)

Альтернативой первой модели является модель Skip-gram. Она имеет схожую архитектуру, но при этом основывается на обратной идее: предсказании контекста по слову, то есть получая на вход слово, как единичный вектор, она пытается предсказать вероятности вхождения каждого из слов в контекст текущего.

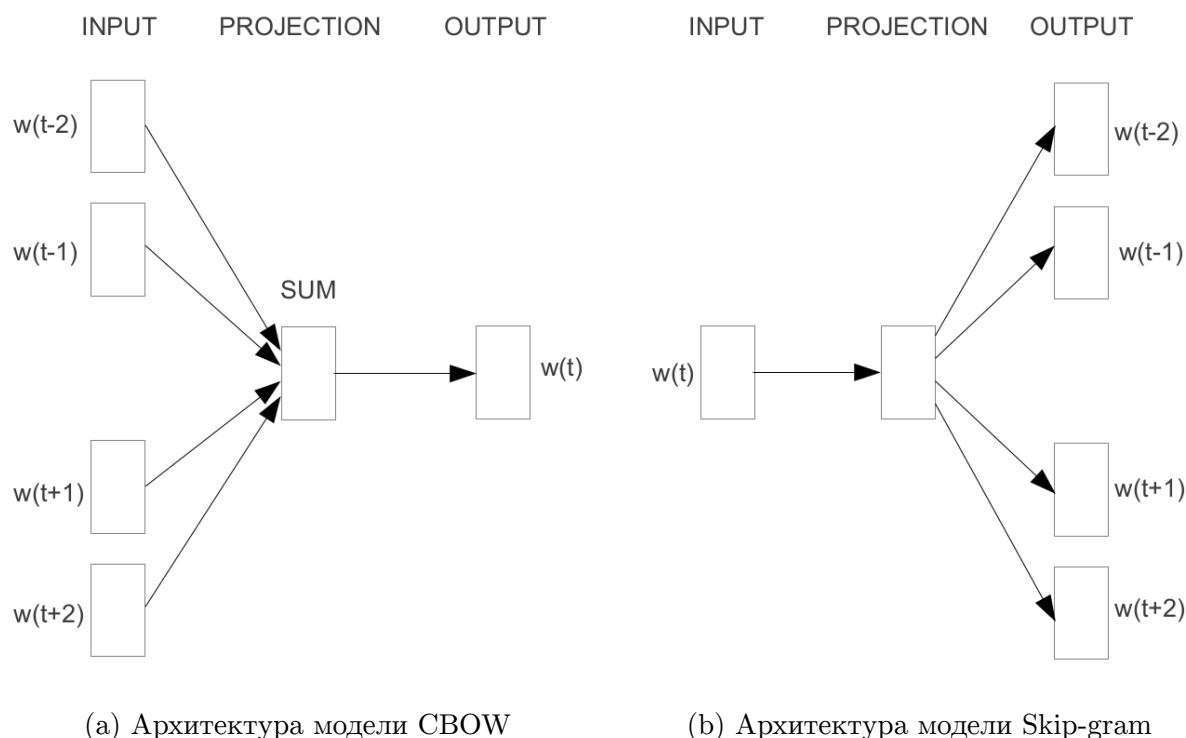


Рис. 2: Модели Word2Vec

Стоит отметить, что для получения векторных представлений хорошего качества требуется обучение Word2Vec на большом корпусе неразмеченных текстов (больше миллиарда слов). Интересной особенностью полученных векторизаций является не только близкое расположение векторов, отвечающих за близкие по смыслу слова, но и сохранение семантических отношений между словами: например, разность векторов, соответствующих словам "король" и "королева" очень близка к разности векторов "мужчина" и "женщина".

1.2.3. Fasttext

Тем не менее, у предыдущей модели есть один достаточно важный недостаток: она не позволяет получать векторные представления слов, которые не встречались в обучающей выборке.

Для того, чтобы решить эту проблему, используется модификация модели Word2Vec — Fasttext [9]. Если раньше в качестве представления слова мы использовали единичный вектор, то в модели Fasttext из слова также выделяются все n -граммы — непрерывные подстроки длины n (в базовой версии для всех n от 3 до 6), и в качестве представления слова берется взвешенная сумма векторов, соответствующих n -граммам.

Подобный трюк позволяет получать векторизации для слова не из словаря, используя векторизации его n -грамм.

1.3. Нейронные сети

За последние годы в области обработки естественных языков, в частности для задачи классификации текстов, нейронные сети стали показывать блестящие результаты. Ключевым эффектом является векторное представление целого текста, которое отражает его семантику, то есть смысл, заложенный в нем.

Для этого используется две принципиально разные архитектуры, описание которых представлено ниже.

1.3.1. Рекуррентная нейронная сеть

Базовые идеи архитектуры рекуррентных нейронных сетей (RNN), были разработаны еще в 1980-е годы. В основе этой архитектуры лежит использование рекуррентных слоев (рисунок 3), которые позволяют работать с данными, которые являются некоторой последовательностью. Для этого в качестве входных данных на каждую следующую ячейку RNN поступает не только новый элемент последовательности, но и некоторая информация из состояния предыдущей ячейки.

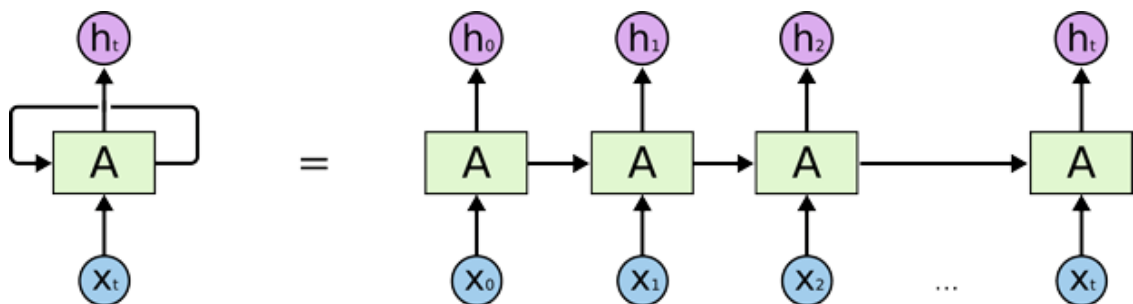


Рис. 3: Рекуррентная нейронная сеть

Данный механизм может помочь и в решении задачи классификации текстов — последовательности слов, позволяя учитывать связь различных слов и предложений друг с другом. Для векторизации слов может быть использован один из методов из подраздела 1.2.

Тем не менее, в текстах часто важен достаточно широкий контекст, а обычные RNN на практике оказываются неспособны захватывать подобные связи. Чтобы избежать

подобной проблемы были придуманы LSTM-ячейки [14] (рисунок 4), которые способны сохранять информацию в течение длительного времени.

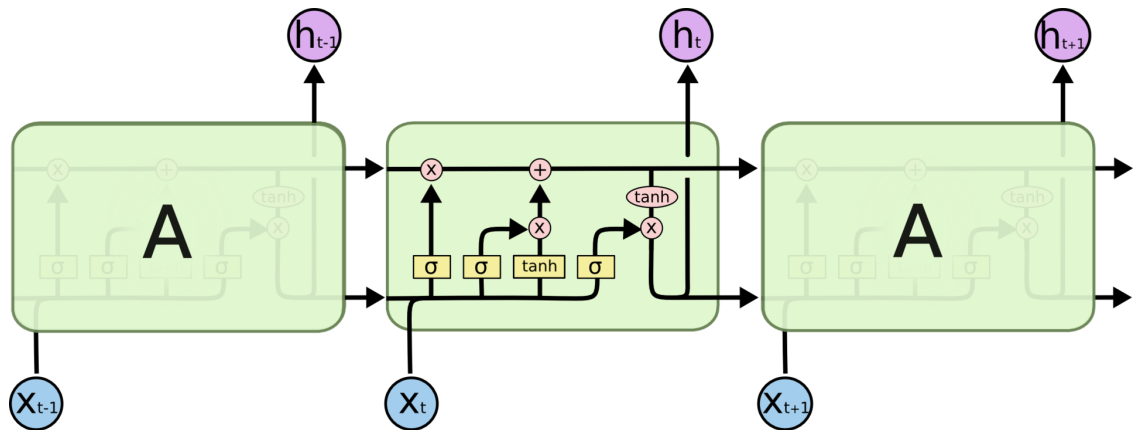


Рис. 4: LSTM-ячейка

Подобная архитектура позволяет захватывать контекст слева, хотя при анализе текста важен контекст, окружающий слово с обеих сторон. Для этого придумана модификация обычного рекуррентного слоя — двунаправленный рекуррентный слой [27], который представляет из себя два слоя, в один из которых данные подаются на вход в прямом порядке, а в другой — в обратном, а в качестве выхода используется конкатенация выходов с этих двух слоев.

1.3.2. Сверточная нейронная сеть

Основой для сверточных нейронных сетей [2] служат сверточные слои, которые работают следующим образом: на вход этому слою подается матрица, после чего каждый фрагмент этой матрицы скалярно умножается на ядро свертки, и получившаяся матрица передается следующему слою сети.

Несмотря на то, что изначально этот механизм применялся для работы с изображениями, последние исследования [19] показали, что он также эффективен при работе с текстом.

На рисунке 5 показан пример архитектуры, используемой для классификации текстов. В качестве входной матрицы подаются векторизованные представления слов текста, после чего к ним применяется свертка размером $k \cdot D$, где D — размерность векторного представления, а k — параметр сети. Обычно используется несколько сверток с различными ядрами, при этом используются k , равные 2, 3, 5, 7, что позволяет захватить контекст вокруг слова. После сверточного слоя используется глобальный Max-Pooling слой, который оставляет максимальное значение из каждой свертки, тем самым позволяя уменьшить размерность данных. Получившиеся значения конкатенируются и отдаются на вход следующему слою. В данной архитектуре используется

обычный полносвязный слой, с помощью которого определяется принадлежность текста к одному из двух классов.

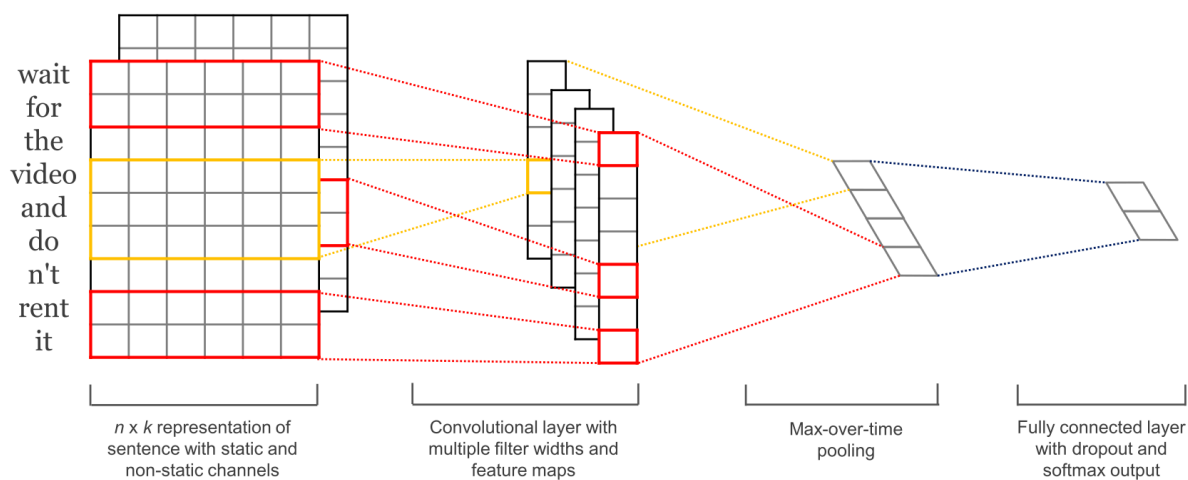


Рис. 5: Сверточная нейронная сеть для классификации текстов

2. Цель и задачи

Рассмотрев особенности уже существующих решений, постараемся избавиться от их недостатков, используя для этого передовой инструмент последних лет, использующийся для классификации текстов, — нейронные сети.

Таким образом, целью данной работы является разработка системы, позволяющей определять правильность ответа на StackOverflow, используя тексты вопроса и ответов.

Для достижения описанной выше цели нам потребуется решить следующие задачи:

- Собрать и проанализировать данные, которые будут использоваться в дальнейшем для обучения и тестирования моделей.
- Реализовать несколько моделей бинарных классификаторов на основе нейронных сетей, использующих текст ответа и предсказывающих его правильность.
- Модифицировать имеющиеся модели таким образом, чтобы они также использовали текст вопроса.
- Произвести обучение и тестирование различных конфигураций моделей.
- Сравнить полученные результаты с базовыми алгоритмами и результатами других исследований.

3. Описание данных

Для задач, решаемых методами машинного обучения с учителем, к числу которых принадлежат и нейронные сети, важно наличие данных, на которых можно обучать, а затем и тестировать алгоритм.

3.1. Сбор данных

Так как большая часть имеющихся работ проводило тестирование на датасетах одного из сайтов системы StackExchange — Server Fault [29], то для удобства сравнения результатов мы также решили использовать эти данные. Кроме того, этот датасет удобен тем, что содержит достаточное количество данных для тренировки, при этом с ним удобно работать во время предобработки данных, в отличие от базы данных вопросов сайта StackOverflow, размер которой превосходит 50 гигабайт.

Для сбора данных использовалась база вопросов Server Fault, размещенная в открытом доступе [24] и содержащая данные с апреля 2009 года по март 2018 года. База представляет из себя XML-файл размером ≈ 0.9 гигабайт, каждая строка которого соответствует либо вопросу, либо ответу. В таблице 1 кратко описаны атрибуты имеющегося датасета.

| Атрибут | Описание |
|------------------|---|
| Id | Идентификатор поста |
| PostTypeId | Тип поста |
| AcceptedAnswerId | Идентификатор выбранного автором ответа (если есть) |
| ParentId | Идентификатор вопроса, на который был дан текущий ответ |
| CreationDate | Дата создания поста |
| Score | Рейтинг поста |
| ViewCount | Количество просмотров поста |
| Body | Текст поста в виде HTML-текста |
| OwnerUserId | Идентификатор автора поста |
| Title | Название поста |
| Tags | Теги вопроса |
| AnswerCount | Количество ответов на вопрос |
| CommentCount | Количество комментариев к посту |

Таблица 1: Описание полей датасета

3.2. Анализ данных

После анализа было получено, что в имеющемся датасете содержится:

- 684 тысячи постов, из них 257 тысяч вопросов и 427 тысяч ответов.

- 130 тысяч вопросов (51%) без выбранного автором вопроса ответа.
- 28 тысяч вопросов (11%), у которых нет ни одного ответа.

Кроме того, не учитывались вопросы, у которых нет ни одного ответа (так как они не давали никаких новых данных в рамках нашей задачи), а также вопросы с отрицательным рейтингом, так как вероятнее всего это шумовые данные, в которых заданный вопрос либо плохо сформулирован, либо не соответствует тематике сайта.

После подобной фильтрации осталось 80 тысяч вопросов и 180 тысяч ответов.

Также было проанализировано содержание текстов вопросов и ответов, в ходе которого выяснилось, что тексты содержат большое количество терминов и слов с опечатками.

На рисунке 6 можно увидеть распределения длин текстов вопросов и ответов, которые понадобятся нам в дальнейшем для построения моделей классификаторов.

Медианным значением длины текста является 500 для вопросов и 350 для ответов.

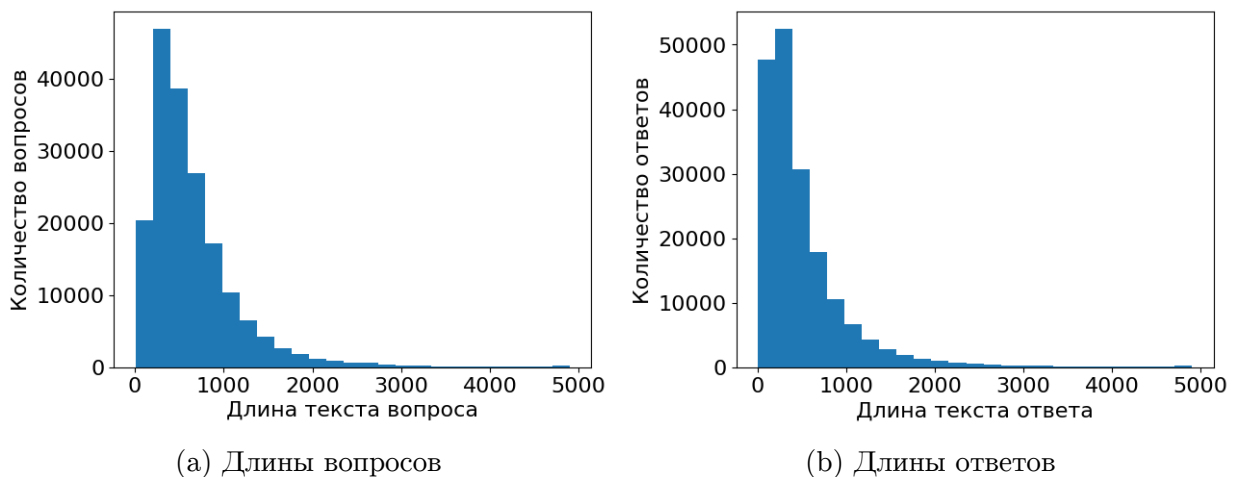


Рис. 6: Распределения длин текстов

3.3. Предобработка данных

Описанный выше датасет нужно привести к форме, с которой будет удобно работать в дальнейшем, для этого был написан скрипт на языке Python, который сохраняет данные в формате CSV-файла, каждая запись в котором содержит:

- Признаки, относящиеся к ответу: идентификатор и текст ответа, был ли ответ выбран автором вопроса, как правильный.
- Признаки, относящиеся к вопросу: идентификатор, название и текст вопроса, на который был дан текущий ответ.

- Признаки, являющиеся метаинформацией: количество ответов на вопрос, возраст ответа, то есть разница в миллисекундах между датами создания ответа и вопроса.

Так как в исходном датасете текст хранился в виде HTML кода, был написан обработчик дерева DOM [34], который:

- Очищает текст от тегов.
- Игнорирует содержимое тегов типа `del`, `s` и `strike`.
- Заменяет содержимое тега `code` на специальное слово `CODE_LEXEM`.

Обработанный подобным образом текст также сохраняется в CSV-файле.

3.4. Постобработка данных

Для решения задачи машинного обучения итоговая выборка ответов была подразделена на 3: обучающую, проверочную и тестовую.

Сначала из всех ответов была выделена тестовая подвыборка размером 20% от всей выборки, а затем с таким же соотношением оставшаяся выборка была разделена на тренировочную и проверочную. Размер подвыборки 20% был выбран, как один из наиболее распространенных для задач машинного обучения.

Итого имеющийся датасет был разбит на три выборки:

- Обучающая выборка: 51000 ответов.
- Проверочная выборка: 13000 ответов.
- Тестовая выборка: 16000 ответов.

4. Методы и реализация

4.1. Извлечение признаков

Напомним, что после предобработки датасета мы получили данные, в которых присутствуют тексты ответа и вопроса, к которому относится текущий ответ, а также некоторая метайнформация.

Нейронные сети работают с числовыми признаками, поэтому в следующих подразделах будет, во-первых, описан использованный подход к векторизации слов, а во-вторых, перечислены дополнительные признаки, которые могут помочь в поставленной перед нами задаче классификации.

4.1.1. Векторное представление вопроса и ответа

Как уже говорилось в подразделе 1.2, существует несколько способов векторизации слов, но нас будут интересовать Word2Vec и Fasttext, которые помимо прочего хорошо описывают семантику слов.

Для обучения этих алгоритмов требуется большой корпус неразмеченных данных, в нашей работе использовалось обучение на базе данных [33], состоящей из статей Wikipedia до октября 2017 года. В этом корпусе присутствует ≈ 4.9 миллиона статей, в которых содержатся ≈ 132 миллиона предложений и ≈ 2.5 миллиарда слов. При обучении игнорировались слишком короткие предложения (длиной менее 5), а также использовалась модель CBOW.

В связи со спецификой Server Fault в текстах вопросов и ответов встречается множество технических слов и терминов, кроме того, как и на любом форуме, часто пользователи пишут слова с ошибками или опечатками. Это приводит к тому, что модель Word2Vec, обученная по корпусу Wikipedia, не содержит множества слов, присутствующих в текстах вопросов и ответов, в том числе терминов, которые являются важной информацией. В среднем, модель не содержала векторных представлений ≈ 10 слов на каждую пару вопроса и ответа.

Чтобы это исправить, было предпринято два шага:

- Так как часто новые слова являются конкатенацией слов из словаря, например, *texthtml* или *powercenter*, а слова с опечатками отличаются от слов из словаря всего одним символом, то есть множество n -грамм у них сильно пересекается с множеством n -грамм слов из словаря, то вместо использования Word2Vec мы решили использовать модификацию, которая работает с n -граммами — Fasttext.
- Также у нас есть специализированные тексты, которые можно использовать для обучения нашей модели, поэтому модель была дотренирована на текстах вопросов и ответов из обучающей подвыборки.

Английский — аналитический язык, что значит, что в нем грамматические отношения передаются через синтаксис, то есть через отдельные служебные слова. Например, в отличие от синтетического русского языка в нем нет склонений.

А исходя из того, что корпус текстов Wikipedia большой, то отпадает потребность производить стемминг, который наоборот может лишь ухудшить качество векторизации, склеивая несколько различных слов в одно.

4.1.2. Лингвистические признаки

Кроме содержания, часто важно и грамотное оформление ответа: логичное разбиение на предложения и параграфы, проставленные ссылки на источники и удобочитаемость. Поэтому в качестве дополнительных лингвистических были извлечены следующие признаки, описывающие либо структуру HTML-кода, либо текста ответа:

- Количество ссылок в тексте, то есть количество тегов `<a>`.
- Количество фрагментов кода в тексте, то есть количество тегов `<code>`.
- Количество параграфов в тексте, то есть количество тегов `<p>`.
- Доля заглавных букв в тексте, то есть отношение количества заглавных букв к длине текста.
- Доля строчных букв в тексте, то есть отношение количества строчных букв к длине текста.
- Доля пробелов в тексте, то есть отношение количества пробельных символов к длине текста.
- Длина текста.
- Длина самого длинного предложения.
- Количество слов в самом длинном предложении.
- Среднее количество слов в предложении.
- Средняя длина слова в тексте.
- Количество предложений.
- Автоматический индекс удобочитаемости [30]:

$$ARI = 4.71 \cdot \frac{characters}{words} + 0.5 \cdot \frac{words}{sentences} - 21.43 \quad (2)$$

- Индекс удобочитаемости Флеша [10]:

$$FRE = 206.835 - 1.015 \cdot \frac{words}{sentences} - 84.6 \cdot \frac{syllables}{words} \quad (3)$$

- Индекс SMOG [21]:

$$SMOG = 1.043 \cdot \sqrt{\frac{30 \cdot polysyllables}{sentences}} + 3.1291 \quad (4)$$

- Индекс Флеша-Кинкейда [7]:

$$FK = 0.39 \cdot \frac{words}{sentences} + 11.8 \cdot \frac{syllables}{words} - 15.59 \quad (5)$$

- Индекс Колман-Лиау [6]:

$$CLI = 5.88 \cdot \frac{characters}{words} - 29.6 \cdot \frac{sentences}{words} - 15.8 \quad (6)$$

- Фог-индекс [12]:

$$FOG = 0.4 \cdot \left(\frac{words}{sentences} + 100 \cdot \frac{polysyllables}{words} \right) \quad (7)$$

- Индекс LIX [3]:

$$LIX = \frac{words}{sentences} + 100 \cdot \frac{longwords}{words} \quad (8)$$

В упомянутых выше формулах использованы следующие обозначения:

- *characters* — количество символов в тексте,
- *words* — количество слов в тексте,
- *sentences* — количество предложений в тексте,
- *syllables* — количество слогов в тексте,
- *polysyllables* — количество слов с хотя бы тремя слогами,
- *longwords* — количество слов длины хотя бы 7.

На рисунке 7 можно увидеть распределение значений индекса удобочитаемости LIX правильных и неправильных ответов, из которого видно, что подобные индексы действительно влияют на решение пользователя о выборе правильного ответа.

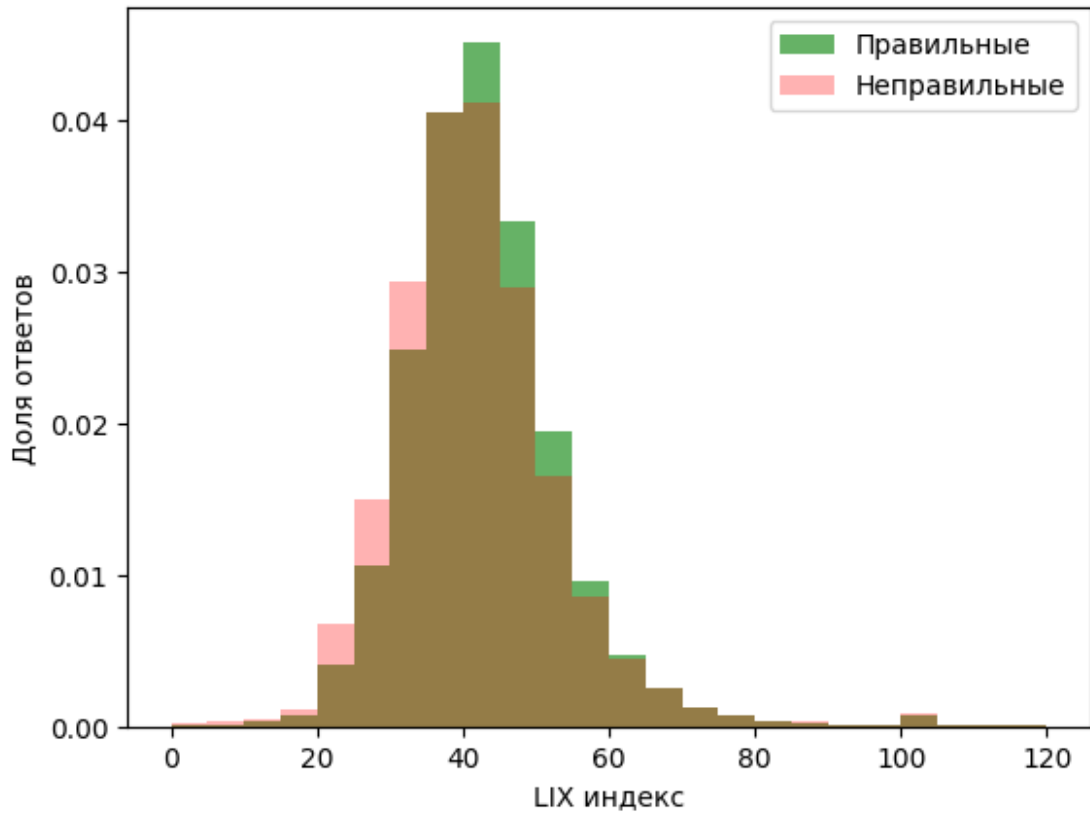


Рис. 7: Распределение значений индекса LIX

4.1.3. Словарные признаки

Как было замечено ранее, в текстах встречается множество терминов, поэтому для оценки того, насколько ответ соответствует заданному вопросу было также добавлено 4 дополнительных словарных признака:

- Коэффициент Жаккара [16] для текстов вопроса и ответа:

$$J = \frac{|V_Q \cap V_A|}{|V_Q \cup V_A|} \quad (9)$$

где V_Q — набор слов вопроса, а V_A — набор слов ответа.

- Коэффициент Жаккара для текстов вопроса и ответа без учета стоп-слов.
- Модификация первого словарного признака с учетом IDF, то есть теперь вес слова не единичный, а обратно пропорционален частоте употребления этого слова в данном датасете.
- Аналогичная модификация второго словарного признака.

4.1.4. Метаинформация

В качестве метаинформации используется всего 2 признака:

- Количество ответов на заданный вопрос.
- Возраст ответа, то есть количество времени в миллисекундах, которое прошло от момента создания вопроса до публикации текущего ответа.

Стоит отметить, что перед вычислением векторных представлений вопросов и ответов, а также словарных признаков, тексты были очищены от любых небуквенных символов, приведены к нижнему регистру и разбиты на отдельные слова.

Также хочется подчеркнуть, что данная работа ставит своей целью определение правильного ответа без использования какой-либо внешней информации о рейтинге, а лишь используя содержания и даты публикации вопросов и ответов, так как часто (например, для новых или непопулярных вопросов) подобная информация может быть недоступна. Поэтому в качестве признаков не извлекаются рейтинги ответов или репутация пользователей-авторов, хотя это и могло бы повысить качество классификатора.

4.1.5. Нормализация

Перед использованием извлеченных признаков необходимо произвести их нормализацию (то есть приведению их значений к отрезку $[0, 1]$). Были посчитаны три вида нормализации:

- **Глобальная нормализация.**

Обычная нормализация, когда максимальное значение признака переводится в 1, а минимальное — в 0.

- **Нормализация по вопросу.** Так как признаки ответов могут сильно отличаться (например, длина ответа может быть от нескольких десятков до нескольких тысяч символов), а ответы для конкретного вопроса выбираются из небольшого подмножества, то была добавлена нормализация по вопросу: процесс, когда все ответы группируются по вопросу, к которому они относятся, а затем проводится нормализация внутри каждой подгруппы.
- **Дискретизация.** Кроме того, могут быть важны не конкретные значения признаков, а лишь их порядок относительно признаков других ответов на текущий вопрос. Для учета этого, были добавлены дискретизированные признаки: если обозначить множество ответов в одной подгруппе за A , то все признаки ответов внутри A сортируются и переводятся в множество $\left\{ \frac{1}{|A|}, \frac{2}{|A|}, \dots, \frac{|A|-1}{|A|}, 1 \right\}$ с учетом

порядка, то есть минимальное значение признака переходит в $\frac{1}{|A|}$, а максимальное — в 1.

Например, если применить дискретизацию к признаку ”возраст ответа”, то мы получим позицию ответа, что является важным признаком при классификации.

4.2. Методы

Для векторизаций текстов вопросов и ответов использовались архитектуры, представленные в подразделе 1.3: либо сверточные слои с ядрами различных размеров для учета контекстов различной длины, либо двунаправленный рекуррентный слой. Стоит отметить, что увеличение количества слоев (например, добавление еще одного рекуррентного слоя) сильно замедляло время обучения, при этом не улучшая метрики на проверочной выборке, поэтому в итоговой модели использовались именно такие конфигурации слоев.

Было рассмотрено несколько различных архитектур, которые по-разному используют признаки первого типа — векторные представления текста вопроса и ответа.

- Первый подход не использует эти признаки вообще: сеть представляет из себя двуслойную нейронную сеть со скрытым полносвязным слоем и линейной функцией активацией, на вход которой подаются лингвистические, словарные и метаинформационные признаки, а на выходе с помощью сигмоидной функции активации получается вероятность принадлежности к классу правильных ответов.
- Тем не менее, мы хотим учитывать в том числе и тексты, поэтому в следующем подходе для учитывания текста ответа использовались сверточные либо рекуррентные слои, после чего полученный вектор, сконкатенированный вместе с остальными признаками использовался в использованной в первом подходе нейронной сети.
- Однако, зачастую просто векторизация текста ответа сама по себе — ни о чем не говорящий признак.

При добавлении векторного представления вопроса, возникает два важных момента, которые хочется рассмотреть подробнее:

- Должны ли параметры рекуррентных или сверточных слоев обучаться совместно для текстов вопросов и ответов или это должны быть одни параметры для вопросов и другие для ответов?

Для ответа на этот вопрос были опробованы оба способа, анализ результатов которых можно найти в подразделе 5.4.

- Как совместно учитывать семантические представления вопроса и ответа?

На выходе из предыдущих слоев мы получаем два вектора Q и A , соответствующие вопросу и ответу соответственно. Были рассмотрены несколько функций преобразования этих векторов в один вектор-признак:

- * Конкатенация:

$$Out = Concat(Q, A) \quad (10)$$

- * Сумма:

$$Out = Q + A \quad (11)$$

- * Косинусный коэффициент:

$$Out = \frac{QA^T}{||Q|| ||A||} \quad (12)$$

- * GESD [1]:

$$Out = \frac{1}{1 + ||Q - A||} \cdot \frac{1}{1 + \exp(-1 - QA^T)} \quad (13)$$

- * AESD [1]:

$$Out = \frac{0.5}{1 + ||Q - A||} + \frac{0.5}{1 + \exp(-1 - QA^T)} \quad (14)$$

Для того, чтобы избежать эффекта переобучения, после получения векторных представлений текстов использовался Dropout-слой.

4.3. Реализация модели

В этом подразделе будут продемонстрированы две финальные архитектуры, которые были использованы в сравнении с результатами других работ.

В обоих случаях в качестве функции потерь использовалась бинарная кросс-энтропия, а в качестве оптимизатора использовался Adam [20]. Для определения схожести вопроса и ответа использовалась функция *AESD*, введенная в предыдущем подразделе.

4.3.1. Сверточная нейронная сеть

На рисунке 8 указана архитектура, основанная на сверточной сети. В качестве оптимальных параметров были использованы 128 фильтров: по 32 с размерами 2, 3, 5 и 7, размер скрытого слоя равен 128, а коэффициент Dropout равен 0.5.

4.3.2. Рекуррентная нейронная сеть

На рисунке 9 указана архитектура, основанная на рекуррентной сети. В качестве оптимальных параметров были использованы двунаправленная рекуррентная сеть с

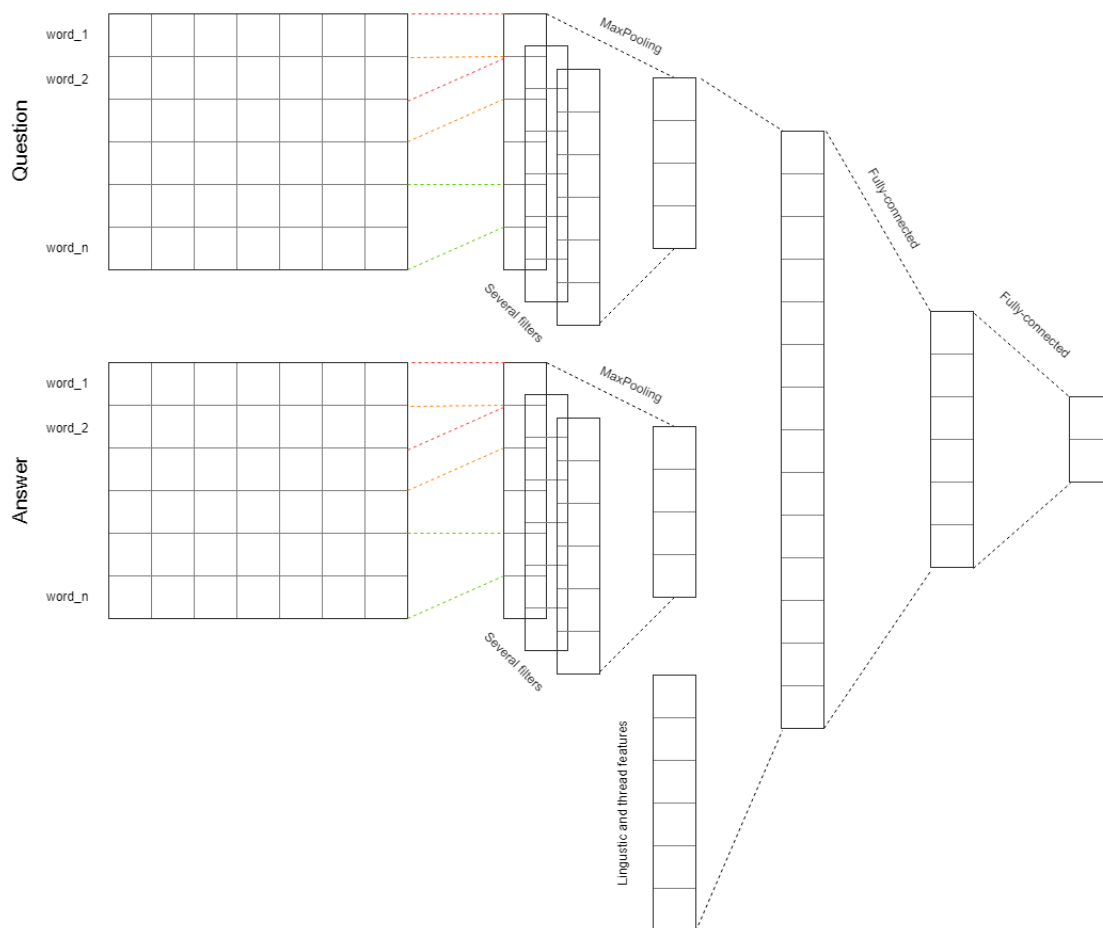


Рис. 8: Архитектура, основанная на сверточных слоях

выходной размерностью, равной 128, размер скрытого слоя равен 256, а коэффициент Dropout равен 0.5.

Стоит отметить, что эта модель требует в разы больше времени на обучение, что объясняется наличием двунаправленного рекуррентного слоя.

4.3.3. Детали реализации

Все вышеобозначенные модели и их модификации были реализованы с помощью фреймворка Keras [18].

Для обучения и работы с моделями Word2Vec и Fasttext использовался фреймворк Gensim [11].

Для предобработки текста использовалась библиотека NLTK [22], а для извлечения различных лингвистических признаков была использована модификация библиотеки textstat [35].

4.3.4. Репозиторий

Исходный код доступен по ссылке [25].

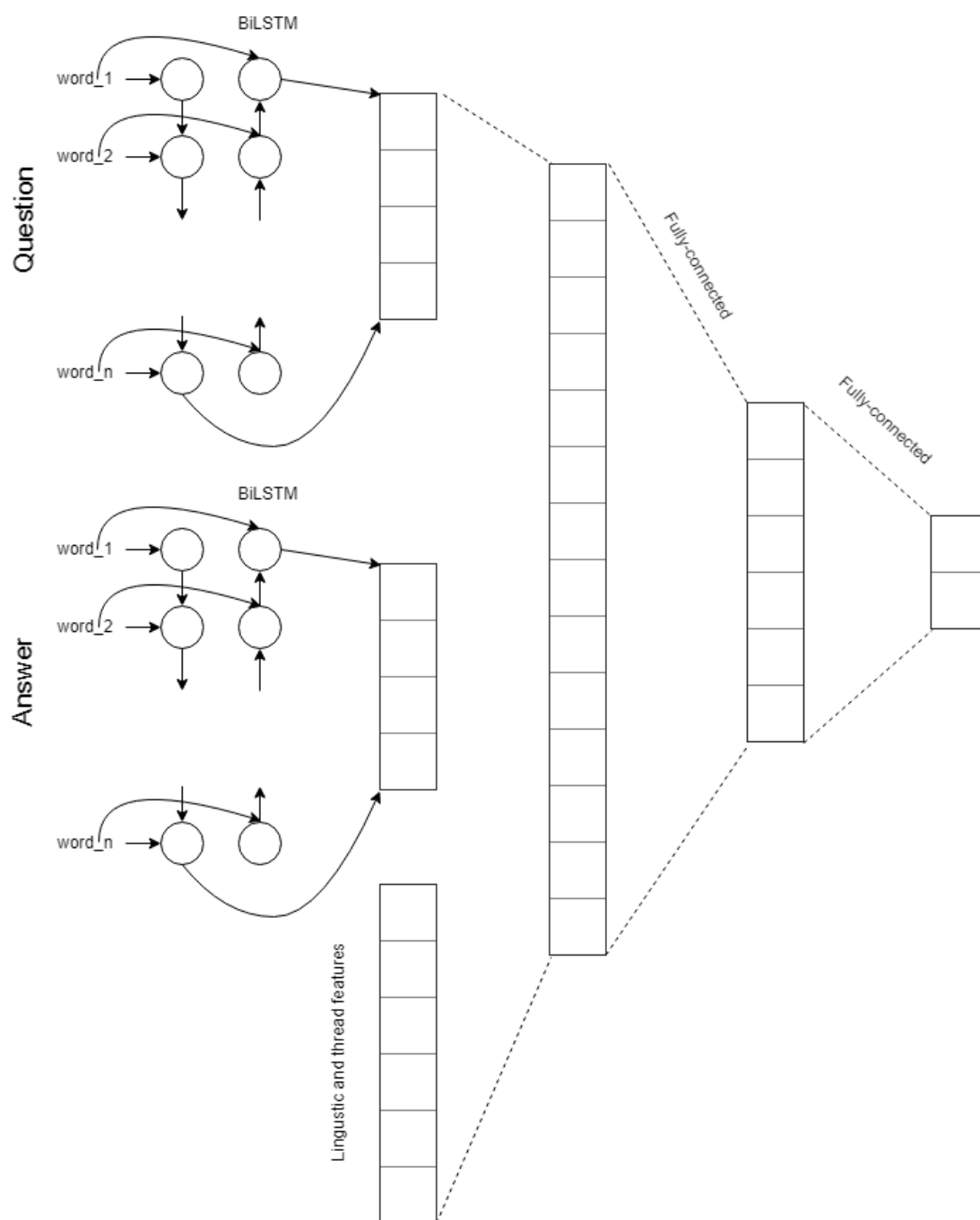


Рис. 9: Архитектура, основанная на рекуррентных слоях

5. Тестирование

Для оценки качества различных модификаций было произведено тестирование обученных нейронных сетей на тестовой выборке.

5.1. Используемые метрики

Для сравнения моделей использовалось несколько метрик, классических для задачи бинарной классификации:

- Аккуратность:

$$Acc = \frac{TP + TN}{TP + FP + FN + TN} \quad (15)$$

- Точность:

$$P = \frac{TP}{TP + FP} \quad (16)$$

- Полнота:

$$R = \frac{TP}{TP + FN} \quad (17)$$

- F -мера, позволяющая объединить точность и полноту в единую метрику:

$$F_1 = \frac{2 \cdot P \cdot R}{P + R} \quad (18)$$

- Площадь под ROC -кривой (AUC) — графиком, показывающим зависимость количества верно классифицированных положительных примеров от количества неверно классифицированных отрицательных примеров при изменении порога решающего правила.

В упомянутых выше формулах использованы следующие обозначения:

- TP — количество верно классифицированных правильных ответов.
- FP — количество неверно классифицированных правильных ответов.
- FN — количество неверно классифицированных неправильных ответов.
- TN — количество верно классифицированных неправильных ответов.

5.2. Производительность

Обучение и тестирование производилось на GPU Tesla V100 с 16 гигабайтами памяти. Обучение одной модели занимало несколько часов (от 4 в случае сверточной сети до 15 в случае рекуррентной).

5.3. Результаты

Было произведено тестирование и сравнение моделей, описанных в разделе 4.2, с показателями исследований, упоминавшихся в подразделе 1.1, а также с двумя наивными алгоритмами:

- Наивный байесовский классификатор, на вход которому подается векторное представление ответа по модели Bag of Words с мерой TF-IDF. Для тестирования использовалась реализация из фреймворка Scikit-Learn [28].
- Классификатор, который считает правильным ответ, который был дан раньше остальных.

| Модель | Виды признаков | Acc | P | R | F ₁ | AUC |
|-------------------------|--|-------------|------|------|----------------|-------------|
| First-answer | метаинформация | 0.66 | 0.69 | 0.60 | 0.64 | 0.66 |
| Naive-Bayes with TF-IDF | словарные | 0.6 | 0.58 | 0.54 | 0.56 | 0.59 |
| Burel et al. [5] | лингвистические, словарные, пользовательские, метаинформация | - | 0.77 | 0.77 | 0.76 | 0.83 |
| Tian et al. [32] | лингвистические, метаинформация | 0.72 | - | - | - | - |
| Gkotsis et al. [15] | лингвистические, словарные, метаинформация | - | 0.83 | 0.66 | 0.74 | 0.85 |
| FC | лингвистические, словарные, метаинформация | ? | ? | ? | ? | ? |
| CNN | текстовые, лингвистические, словарные, метаинформация | 0.77 | 0.81 | 0.62 | 0.71 | 0.86 |
| RNN | текстовые, лингвистические, словарные, метаинформация | 0.78 | 0.82 | 0.64 | 0.72 | 0.87 |

Таблица 2: Сравнения качества работы различных моделей

Был произведен перебор таких параметров модели, как размерность выхода LSTM-ячеек в рекуррентной архитектуре или количество фильтров в сверточной, а также коэффициент Dropout и размер скрытого слоя, и выбраны в качестве используемых в итоге те, которые показывают лучшие результаты на проверочной выборке.

Результаты сравнения приведены в таблице 2.

5.4. Анализ результатов

Как мы видим из таблицы, лучше всего себя показала модель, использующая двунаправленную рекуррентную архитектуру и функцию *AESD* для анализа сходства текстов вопроса и ответа.

На рисунке 10 представлена ROC-кривая, соответствующая этой модели, из которой, в частности, видно, что благодаря метаинформации наш классификатор выдает высокие оценки вероятности в большинстве своем правильным ответам (кривая смещена влево).

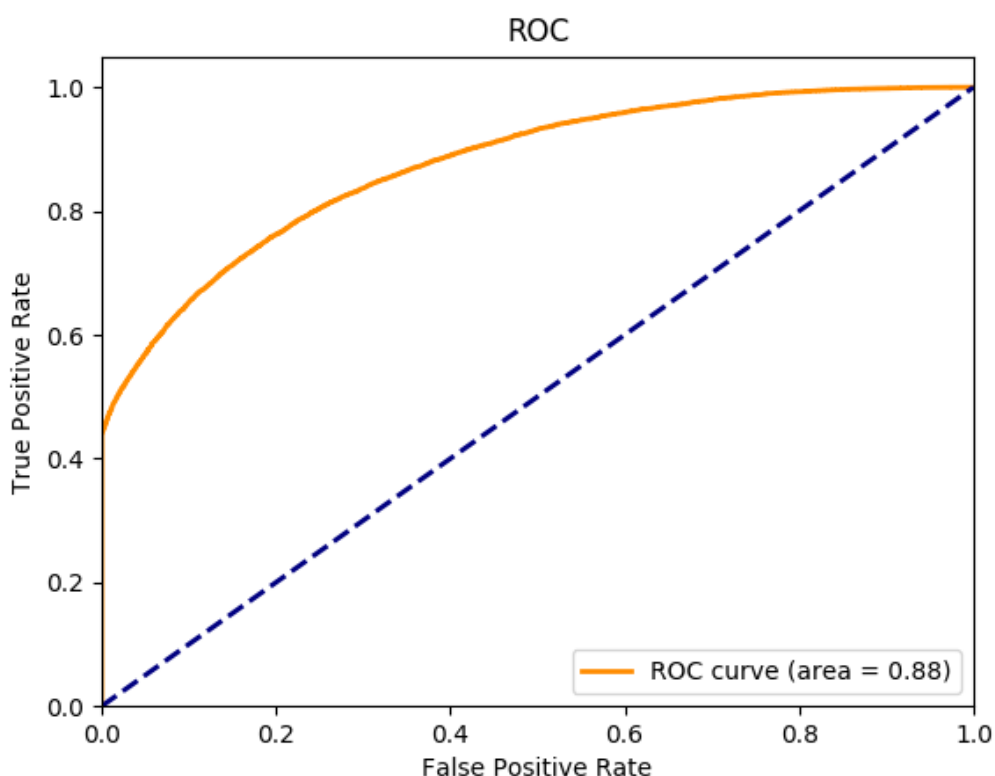


Рис. 10: ROC-кривая

Анализируя полученные результаты, можно отметить следующие вещи:

- Результаты модели, не использующей семантическое представление вопроса совпадают с результатами модели, которая использует только лингвистические и словарные признаки, а также метаинформацию. То есть векторное представление текста ответа в отрыве от вопроса не помогает определить, насколько ответ хорош. В целом, это ожидаемый результат, ведь подобная векторизация создана для передачи смысла текста, что неприменимо в дальнейшем без текста вопроса.
- Были протестированы модели, обучающие параметры сверточных либо рекуррентных слоев и совместно для вопросов и ответов, и по отдельности.

Оказалось, что архитектура с совместными параметрами, во-первых, быстрее обучается, а во-вторых, показывает лучшие результаты на проверочной выборке. Это сходится с показателями других работ по тематике выбора ответа, например [1].

- Были протестированы модели, использующие различные функции определения похожести вопроса и ответа.

При использовании функций конкатенации и суммы время обучения увеличилось по сравнению с тремя остальными, кроме того, они давали худшие результаты при тестировании на проверочной выборке.

Несмотря на то, что в работе [1] функция *AESD* показывала более хорошие результаты, в нашем случае все три остальные функции (косинусный коэффициент, *GESD* и *AESD*) давали в среднем одинаковые метрики.

- В целом, использование текстовых признаков дало незначительный прирост по сравнению с моделью, которая использует признаки остальных типов. Вероятно, это связано с тем, что большая часть вопросов в датасете ожидает описание некоторого процесса, а не какой-то короткий фактологический ответ. Кроме того, даже для человека поставленная задача не является простой: только в 87% случаев ответ с максимальным рейтингом (рейтинг рассчитывается с помощью оценок других пользователей) отмечается автором вопроса как наилучший.

Тем не менее, реализованная модель превосходит как базовые алгоритмы, так и уже существующие решения, использующие другие методы машинного обучения, по метрикам аккуратности и площади под *ROC*-кривой, которые наиболее полно описывают качество бинарного классификатора.

Заключение

В рамках данной работы были достигнуты следующие результаты:

- Рассмотрены общие подходы к задаче классификации текстов и существующие решения для поставленной в данной работе задачи, проанализированы их достоинства и недостатки.
- Подготовлен корпус данных, представляющий из себя ответы на вопросы с сайта Server Fault.
- Предложены и посчитаны дополнительные признаки, которые в дальнейшем используются для классификации.
- Реализовано несколько моделей бинарных классификаторов на основе рекуррентных и сверточных нейронных сетей.
- Произведено тестирование и сравнение полученных результатов с имеющимися в области исследованиями.

Представленная модель комбинирует в себе как семантические признаки, так и более классические лингвистические и словарные, в результате чего достигает хороших результатов в задаче классификации наилучших ответов, которые сравнимы по одним метрикам, либо превосходят по другим уже имеющиеся исследования в этой области.

В качестве дальнейшего развития имеющихся моделей можно выделить следующие аспекты:

- Использование более сложной архитектуры для определения того, решает ли ответ обозначенную в вопросе проблему или нет: например, совместное использование рекуррентных и сверточных слоев, как в работе [26]. Тем не менее, для обучения подобного вида архитектуры может потребоваться больше вычислительных ресурсов, что может служить некоторым ограничением.
- Как уже было отмечено ранее, в связи с технической специализацией вопросов в текстах присутствует множество терминов, поэтому еще одним направлением развития может стать выделение подобных сущностей и отношений между ними [23] и работе с ними отдельно.

В частности, это может помочь понять, насколько ответы, описывающие какой-то процесс, например установки чего-либо, действительно соответствуют тематике вопроса.

- Если пытаться расширить имеющуюся модель на основной сайт StackOverflow, то стоит также учитывать содержание фрагментов кода, которые присутствуют более, чем в половине постов (в то время как на ServerFault код присутствует не в таком большом количестве постов и имеет не такую большую важность). Для этого может быть применен механизм векторизации для кода [4], который имеет что-то общее с принципом работы Word2Vec.

Тем не менее, стоит подчеркнуть, что подобная модификация является действительно сложной задачей, так как зачастую даже человек не может по фрагменту коду без использования дополнительных средств понять, насколько тот хорошо решает проблему автора вопроса.

Список литературы

- [1] Applying Deep Learning to Answer Selection: A Study and An Open Task / Minwei Feng, Bing Xiang, Michael R. Glass et al. // arXiv preprint arXiv:1508.01585. — 2015.
- [2] Back-propagation applied to handwritten zip code recognition / Y. LeCun, B. Boser, J. S. Denker et al. // Neural Computation. — 1989. — Vol. 1, no. 4. — P. 541–551.
- [3] Björnsson Carl-Hugo. Läsbarhet. — Stockholm: Liber, 1968.
- [4] Building Program Vector Representations for Deep Learning / Lili Mou, Ge Li, Yuxuan Liu et al. // arXiv preprint arXiv:1409.3358. — 2014.
- [5] Burel Gregoire, He Yulan, Alani Harith. Automatic Identification of Best Answers in Online Enquiry Communities // ESWC'12 Proceedings of the 9th international conference on The Semantic Web: research and applications. — 2012. — P. 514–529.
- [6] Coleman Meri, Liao T. L. A computer readability formula designed for machine scoring // Journal of Applied Psychology. — 1975. — Vol. 60. — P. 283–284.
- [7] Derivation of new readability formulas (Automated Readability Index, Fog Count and Flesch Reading Ease Formula) for Navy enlisted personnel / J. Peter Kincaid, Robert P. Fishburne, Richard L. Rogers, Brad S. Chissom // Research Branch Report. — 1975. — Vol. 8-75.
- [8] Efficient Estimation of Word Representations in Vector Space / Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean // arXiv preprint arXiv:1301.3781. — 2013.
- [9] Enriching Word Vectors with Subword Information / Piotr Bojanowski, Edouard Grave, Armand Joulin, Tomas Mikolov // arXiv preprint arXiv:1607.04606. — 2016.
- [10] Flesch Rudolf. A new readability yardstick // Journal of Applied Psychology. — 1948. — Vol. 32. — P. 221–233.
- [11] Gensim. — 2018. — Режим доступа: <https://radimrehurek.com/gensim/> (дата обращения: 01.03.2018).
- [12] Gunning Robert. The Technique of Clear Writing. — McGraw-Hill. — P. 36–37.
- [13] Harris Zellig. Distributional structure // Word. — 1954. — Vol. 10. — P. 146–162.
- [14] Hochreiter Sepp, Schmidhuber Jürgen. Long short-term memory // Neural computation. — 1997. — Vol. 9, no. 8. — P. 1735–1780.

- [15] It's all in the Content: State of the art Best Answer Prediction based on Discretisation of Shallow Linguistic Features / George Gkotsis, Karen Stepanyan, Carlos Pedrinaci, John Domingue // Proceedings of the ACM conference on Web science. — 2014. — P. 202–210.
- [16] Jaccard Paul. Étude comparative de la distribution florale dans une portion des Alpes et des Jura // Bulletin de la Société Vaudoise des Sciences Naturelles. — 1901. — Vol. 37. — P. 547–579.
- [17] Jones Karen Spärck. A statistical interpretation of term specificity and its application in retrieval // Journal of Documentation. — 1972. — Vol. 28. — P. 11–21.
- [18] Keras. — 2018. — Режим доступа: <https://keras.io/> (дата обращения: 01.03.2018).
- [19] Kim Yoon. Convolutional neural networks for sentence classification // Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing. — 2014. — P. 1746–1751.
- [20] Kingma Diederik P., Ba Jimmy. Adam: A Method for Stochastic Optimization // arXiv preprint arXiv:1412.6980. — 2014.
- [21] McLaughlin G. Harry. SMOG grading: A new readability formula // Journal of Reading. — 1969. — Vol. 12. — P. 639–646.
- [22] NLTK. — 2018. — Режим доступа: <https://www.nltk.org/> (дата обращения: 01.03.2018).
- [23] Neural Architectures for Named Entity Recognition / Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian et al. // arXiv preprint arXiv:1603.01360. — 2016.
- [24] Overflow Stack. Stack Exchange Data Dump // Internet Archive. — 2018. — Режим доступа: <https://archive.org/details/stackexchange> (дата обращения: 13.03.2018).
- [25] Podguzov Nikita. StackOverflow Correct Answer Detection // Github. — 2018. — Режим доступа: <https://github.com/Nikitosh/StackOverflow-Correct-Answer-Detection> (дата обращения: 01.03.2018).
- [26] Recurrent convolutional neural networks for text classification / Siwei Lai, Liheng Xu, Kang Liu, Jun Zhao // Proceedings of the 2015 Conference of the Association for the Advancement of Artificial Intelligence. — 2015. — P. 2267–2273.

- [27] Schuster Mike, Paliwal Kuldeep K. Bidirectional recurrent neural networks // IEEE Transactions on Signal Processing. — 1997. — Vol. 45, no. 11. — P. 2673–2681.
- [28] Scikit-Learn. — 2018. — Режим доступа: <http://scikit-learn.org/stable/index.html> (дата обращения: 01.03.2018).
- [29] Server Fault // Stack Overflow. — 2018. — Режим доступа: <https://serverfault.com> (дата обращения: 01.03.2018).
- [30] Smith E. A., Senter R. J. Automated readability index // AMRL-TR-66-220. — 1967.
- [31] Stack Overflow // Stack Overflow. — 2018. — Режим доступа: <https://stackoverflow.com> (дата обращения: 01.03.2018).
- [32] Tian Qiongjie, Zhang Peng, Li Baoxin. Towards Predicting the Best Answers in Community-Based Question-Answering Services // Seventh International AAAI Conference on Weblogs and Social Media. — 2013.
- [33] Wikimedia. Wikimedia database dump of the English Wikipedia on October 01, 2017 // Internet Archive. — 2017. — Режим доступа: <https://archive.org/details/enwiki-20171001> (дата обращения: 10.01.2017).
- [34] Wikipedia. Document Object Model // Wikipedia, the free encyclopedia. — 2016. — Режим доступа: https://ru.wikipedia.org/wiki/Document_Object_Model (дата обращения: 03.01.2016).
- [35] textstat. — 2018. — Режим доступа: <https://github.com/shivam5992/textstat> (дата обращения: 01.03.2018).