

Министерство образования и науки Российской Федерации
Федеральное государственное автономное образовательное учреждение высшего
образования
«КАЗАНСКИЙ (ПРИВОЛЖСКИЙ) ФЕДЕРАЛЬНЫЙ УНИВЕРСИТЕТ»

ИНСТИТУТ ВЫЧИСЛИТЕЛЬНОЙ МАТЕМАТИКИ И ИНФОРМАЦИОННЫХ
ТЕХНОЛОГИЙ

КАФЕДРА СИСТЕМНОГО АНАЛИЗА И ИНФОРМАЦИОННЫХ ТЕХНОЛОГИЙ

Направление: 02.04.02 — Фундаментальная информатика и
информационные технологии
Профиль: Математические основы и программное обеспечение
информационной безопасности и защиты информации

МАГИСТЕРСКАЯ ДИССЕРТАЦИЯ

**Разработка метода получения векторного представления слов,
отсутствующих в обучающей выборке**

Студент 2 курса
группы 09-635

"__" _____ 20__ г.

(Боробов П.А.)

Научный руководитель
к.ф.-м.н., ассистент КСАИТ

"__" _____ 20__ г.

(Разинков Е.В.)

Заведующий кафедрой
д.т.н., профессор

"__" _____ 20__ г.

(Латыпов Р.Х.)

Казань – 2018

Оглавление

Введение	4
1. Литература	7
1.1. Определение тональности текста	7
1.1.1. Лингвистическая структура предложения	7
1.1.2. Рекурсивные нейронные модели	8
1.1.3. Сверточные нейронные сети в классификации пред- ложений	9
1.1.4. LSTM в виде дерева	12
1.2. Семантическое векторное представление	13
1.2.1. NLP и задача семантического векторного представ- ления слов	13
1.2.2. Описание задачи для слов отсутствующих в обу- чающей выборке	14
1.2.3. Модель Mimick	14
1.2.4. Модель Enriching	16
1.2.5. Модель Glove	18
2. Разработки	21
2.1. Определение тональности	21
2.1.1. Распределение эмоционального окраса в предло- жениях	21
2.1.2. Увеличение выборки: замена синонимов	22
2.1.3. Увеличение выборки: замена слова на усредненное описание	23
2.1.4. Углубление архитектуры нейронной сети	24
2.2. Векторное представление слов	24
2.2.1. Модель Mimick и наша модификация	24
2.2.2. Наша модификация GloVe + SISG	25
3. Способы валидации	26
3.1. Валидация на моделировании английского языка	26

4. Детали реализации	28
4.1. Детали реализации методов векторного представления .	28
4.1.1. Mimick с контекстным вектором	28
4.1.2. GloVe + SISG	29
5. Результаты экспериментов	31
5.1. Результаты экспериментов разработанных методов опре- деления тональности	31
5.1.1. Оценка методик на существующих реализациях .	31
5.1.2. Эксперимент: замена слова на его описание	33
5.2. Результаты экспериментов разработанных методов век- торного представления слов	33
5.2.1. Результаты модификации модели Mimick	33
5.2.2. Результаты GloVe + SISG	35
Заключение	36
Список литературы	37

Введение

Область обработки естественных языков (англ. Natural Language Processing, далее NLP) активно развивается несколько десятков лет. Задачи решаемые этой областью имеют широкое применение в реальном мире, например, после появления различных интернет-сервисов, которые предоставляют возможность динамически создавать или менять контент интернет-страниц, текстовой информации стало очень много и ее ручная обработка стала занимать невероятное количество времени. NLP позволяет автоматизировать обработку текста с целью выявления некоторой структурированной информации, которая может быть использована в каких-либо бизнес-целях. В данной работе была рассмотрена задача определения тональности текста. Эта задача не нова и актуальна по сей день. Об этом говорят проводимые каждый год конференции (Sentiment Analysis Symposium [19], SemEval [18]). Что же такого интересного в этой задаче, почему она привлекает такой интерес, и почему решение этой задачи является важным для области анализа текста на естественном языке? Для начала стоит понять в чем заключается эта задача и уже потом посмотреть на ее применение. Смысл задачи заключается в том, чтобы по одной фразе, предложению или абзацу, предсказать содержащуюся в них эмоцию, сказать какой характер носит текст: позитивный или негативный, или усложняя эту задачу производить оценку по шкале (например, от 1 до 5 звезд). Получается, данный метод может быть применен по отношению к любому тексту. Существует много интернет-сервисов, где данный метод применяется для оценки пользовательского мнения о товарах. Хорошим примером является Google Product Search, Bing Shopping [12]. В этих сервисах совместно с определением эмоционального окраса текста используются техники извлечения атрибутов объекта и определяется мнение не всего объекта, а его атрибутах или характеристиках. Понятное дело, что задача не была бы такой актуальной, если бы ее применение было возможно только в интернет-сервисах. Так в Gollup Poll [12] исследовали корреляцию эмоционального окраса сообщений в твиттере и показате-

лем доверия потребителя, в результате было выявлено, что изменение некоторых экономических параметров может влиять на мнение людей в социальных сетях. Ключевым применением, конечно, является предсказание трендов на рынках биржевых товаров [12], которое дало бы преимущество в игре на таком рынке. Данную задачу можно решать как методами основанными на словарях, так и с использованием методов машинного обучения. В работе будут рассмотрены алгоритмы, основанные на машинном обучении. Задача может быть сформулирована, как задача двухклассовой или многоклассовой классификации. Для ее постановки необходимы некоторые текстовые признаки и предсказываемая метка этого текста, на их основе может быть построена выборка, которая будет использована при обучении моделей. Модели описанные в этой работе используют некоторое машинное представление текста, такой способ представления текста подразумевает использование чисел или вектора из чисел при представлении конкретного слова. Данный вектор может быть представлен различными способами. Например, наиболее интуитивный способ, представить слово в виде вектора состоящего из нулей и единицы в той позиции, где это слово расположено в словаре. В данной работе были рассмотрены различные способы представления слов. Семантическое векторное представление слов - это отдельная область в NLP. Она занимается разработкой алгоритмов, получающих векторные представления слов способных распознать семантику, например, $sim(v(good), v(bad)) \approx 1$. Подробнее о свойствах такого представления слов будет описано в разделе 1.2. Важно отметить, что решаемая задача, позволяет представлять вектор более компактно и, как показывают результаты из статей [4] [10] [15] [5] позволяют улучшать результат решения целевых задач, например, задачи определения тональности текстов. В области векторного представления остается много нераскрытых вопросов. Одним из таких вопросов является векторное представление слов, отсутствующих в словаре. Решаемая в этой диссертации задача только недавно получила широкое распространение, но уже сейчас можно наблюдать непреодолимый интерес исследователей в области обработки естественного языка. Об

этом могут говорить последние конференции ориентированные на эту область, такие как: EMNLP, CoNLL, ACL и т.п. В этих конференциях всегда специальным образом выделялись исследования связанные с возможностью получения векторных представлений по тем словам, которые могли не участвовать в процессе обучения. Данная проблема может быть поставлена как расширение основной задачи по векторному представлению слов, поэтому многие исследователи, занимающиеся этим направлением, являются также авторами статей, решающих и эту проблему. На текущий момент существует несколько работ решающих эту задачу. Одной из таких работ является [5], где авторы статьи моделируют векторное представление слова, как усредненную поэлементную сумму составных частей слова, представленных в векторной форме. Также существуют и более тривиальные подходы для решения этой задачи, например, предложенная в статье модель [16], где проверяется гипотеза, что при помощи информации о символах можно восстановить векторное представление слова по уже составленному словарю. Оба подхода будут рассмотрены в данной работе и будут приведены результаты валидации по каждому алгоритму. В данной работе были использованы различные подходы для решения этой задачи. Например, описанный в разделе 2.2.1 алгоритм, восстанавливает векторные представление лучше, чем тот алгоритм, на котором были основаны разработки описанные в главе 2, но при этом такой результат все равно сильно уступал текущим state-of-the-art работам в области векторного представления слов. Также в этой работе возникли проблемы по оценке качества полученных векторных представлений подробнее об этом будет рассказано в разделе 5.2.1. Был разработан и реализован еще один метод решающих данную задачу, который был основан на двух работах [5] [15]. Такой алгоритм получил state-of-the-art результат на некоторых валидационных выборках и был получен state-of-the-art результат при решении задачи языкового моделирования английского языка.

1. Литература

1.1. Определение тональности текста

Современные алгоритмы определения эмоционального окраса текста содержат в себе различные подзадачи из разных областей. Сюда входят семантическое векторное представление слов, определение частей речи и логическое представление текста.

1.1.1. Лингвистическая структура предложения

Важно понимать, что любое предложение на естественном языке представляет собой некоторую структуру, которая может быть по-разному интерпретирована. Существует много работ в данной области, которые пытаются организовать последовательности слов в такие структуры. Логическое представление – это такое представление, при котором выражается некоторая логическая взаимосвязь слов, определяется их порядок по смыслу текста. Это некоторое представление текста в виде дерева, где выражается нисходящая взаимосвязь слов. Ключевым в этом представлении является ее структура в виде дерева, которая может быть использована в рекурсивных моделях. Такая форма представления текста позволяет подсказать алгоритмам информацию о структуре предложения и о его смысле, что непременно дает улучшение результата в решаемых задачах. Так, например, была улучшена точность алгоритма, использующего рекуррентные нейронные сети [17], где впервые была использована структура предложения в виде дерева и архитектура сети, и где слова были представлены векторами. В этом разделе будут описаны state-of-the-art работы последних лет [10][17]. Также в этих алгоритмах были раскрыты варианты использования различных представлений текста, которые позволяют использовать последние наработки в области глубинного обучения. Описанные в секциях (1.1.2, 1.1.3, 1.1.4) алгоритмы решают эту задачу с совершенно различными моделями, каждая из которых интересна излагаемой

идеей, каждый из представленных алгоритмов раскрывает некоторый аспект задачи и ее подводные камни, эти модели интересны еще и тем, что здесь происходит качественный переход в области NLP от настраиваемых вручную признаков в сторону автоматически определяемых моделью, то есть являющимися обучаемыми признаками.

1.1.2. Рекурсивные нейронные модели

Рекурсивные нейронные модели (Recursive Neural Models) – модели, использующие векторное представление фраз, которые потом используются для классификации. [17]

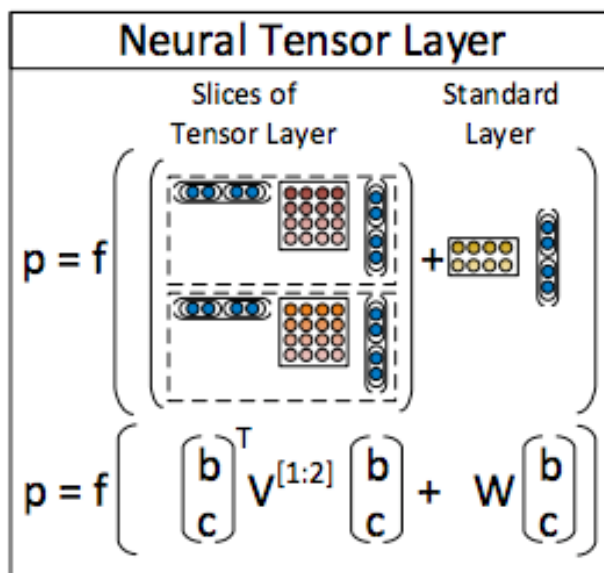


Рис. 1: Представление одного слоя в RNTN [17]

Входящий текст представляется в логической форме в виде дерева, где каждый узел является композицией двух дочерних узлов. Слова представлены векторами размерности $d = 300$, здесь же определена функция $g(w_1, w_2) = p_k$, где $|w_i| = d$, $i = 1, 2$ и $|p_k| = d$. На вход модели дается некоторое заранее заданное представление слов во всем словаре, которая обучается совместно с параметрами модели. Функция активации по вектору в конечном узле по заданной метке класса определена

следующим образом:

$$y^a = \text{softmax}(W_s a), \quad (1)$$

где $W_s \in \mathbb{R}^{c \times d}$ матрица весов, a - метка класса. Основные отличия таких моделей лежат в вычислении скрытых векторов p_k или, другими словами, в определении функции g [14] [17] [6]. Используемая в статье [17] модель, достигла наилучшего результата на момент выпуска статьи на выборках SST-1, SST-2, MR. Результаты представлены в таблице 1.

MR	SST1	SST2
80.5	45.7	85.4

Таблица 1: RNN результаты на различных выборках

Минусом используемой модели является некоторый предварительный этап преобразования сырого текста в его логическую форму. Это накладывает дополнительные ограничения на использование данного алгоритма в прикладных задачах.

1.1.3. Сверточные нейронные сети в классификации предложений

Одна из последних моделей в значительной мере отличающихся от рекурсивных моделей. Здесь заложена идея векторного представления слов и матричного представления документов. Между документом и графическим изображением проводится некоторая параллель, где указывается, что между отдельными элементами этой матрицы могут быть связи, то есть фразы в начале и конце предложения связаны между собой не через цепочку слов, а напрямую, чего нельзя получить в древовидном представлении текста. Подобными идеями задавались при анализе изображений и реализации модели классификации изображений. В работе [8] использовались сверточные нейронные сети (Convolutional Neural Network, далее CNN) для решения этой задачи.

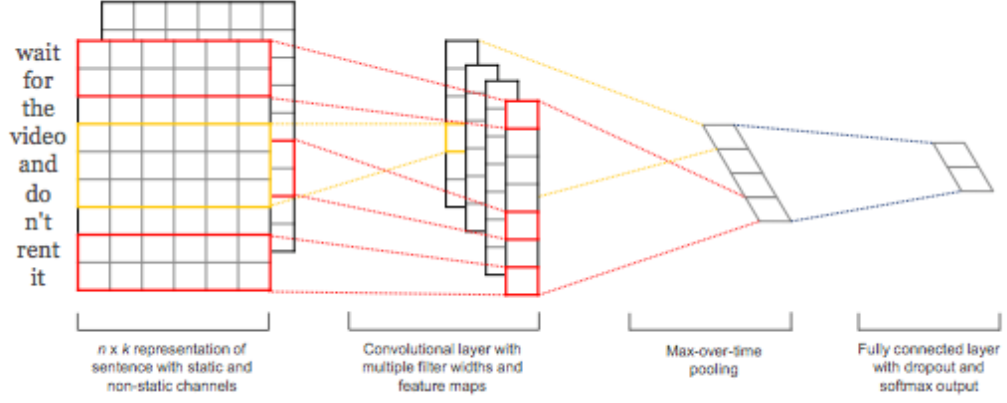


Рис. 2: Архитектура нейронной сети [10]

Все дальнейшее описание модели взято из статьи [10] Пусть у нас есть некоторое векторное представление слов $x_i \in \mathbb{R}^d$, $i = 1..s$, а s – количество слов в предложении. Представим документ в виде вектора, состоящего из слов, в векторной форме фиксированной размерности, в том порядке, в котором слова находились в исходном предложении. К полученной матрице применим операцию свертки с нелинейной функцией активации (ReLU). Размерность фильтра будет определяться размерностью вектора d , а высота вектора h будет варьируемой и является гиперпараметром нейронной сети.

$$X_{1:s} = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_s^T \end{pmatrix} \quad (2)$$

Пусть (2) матрица, где x_i – векторное представление слова на i -ой, тогда функция свертки, для фильтра размерностью $w \in \mathbb{R}^{hk}$, где k определяет количество таких фильтров, определена следующим образом:

$$o_i = f(w \cdot x_{i:i+h-1} + b), \quad (3)$$

где $b \in \mathbb{R}$ – смещение, а f – нелинейная функция (ReLU). После этого

получается некоторая карта признаков по всем регионам:

$$c = [c_1, c_2, \dots, c_{s-h+1}]. \quad (4)$$

После получения карты признаков c выбирается максимальный признак $\bar{c} = \max(c)$, такое поведение обусловлено тем, что длина документов может быть разной, и количество регионов в зависимости от длины документов может изменяться, но архитектура нейронной сети не должна меняться от этого. После проведения этих действий следует полносвязный слой с регуляризацией dropout, который на момент обучения выбрасывает часть переменных. Выход полносвязного слоя подается на вход функции *Softmax*. *Softmax* – функция принимающая на вход n -мерный вектор и с n -мерным выходом, где сумма компонент равна 1 и каждая компонента неотрицательна. Авторы приводят несколько вариантов такой модели, в которой меняется количество входных каналов, меняется способ их инициализации: предобученные векторы статичны, пред обученные векторы нестатичны, случайные векторы. В зависимости от такой инициализации результаты значительно отличаются. Такая модель показала высокий результат (таблица 2).

Модель	SST-1	SST-2	MR
CNN-static	42.2	83.5	80.5
CNN-random	44.8	85.6	75.9
CNN-non-static	46.7	87.0	81.3
CNN-multichannel	44.6	87.1	80.8

Таблица 2: Результаты разных моделей CNN

Здесь приведена архитектура нейронной сети без детального описания всех тонкостей и всех гиперпараметров модели. В статье [22] проводится валидация данной модели. Производится поиск и валидация лучших значений гиперпараметров модели, при которых достигается высокая точность классификации.

1.1.4. LSTM в виде дерева

В статье [20] приводится архитектура рекурсивных нейронных сетей с памятью (RNN-LSTM) решающая задачу определения эмоционального окраса предложения. Данный вид нейронной сети является развитием идеи из алгоритма описанного в разделе 1.1.2. В данном алгоритме также анализируются структура предложения в виде дерева, но ключевым отличием является тот факт, что в такой сети закономерности запоминаются на всей последовательности слов и не отображены только в каких-то конкретных узлах. Одним из преимуществ этой работы по сравнению с работой [17] является количество обучаемых параметров, авторами статьи удалось сократить количество параметров и улучшить результат.

Данная работа не представляет из себя типовое использование LSTM-сетей. Здесь произведена модификация LSTM-ячеек. Обычно у такой ячейки с индексом j имеется входящий и выходящий фильтр i_j , o_j , запоминающий фильтр c_j и скрытое состояние h_j . Отличие Tree-LSTM от обычного LSTM в том, что у Tree-LSTM обновление фильтров зависит от всех дочерних узлов этой ячейки. Также данный тип ячейки содержит забывающие фильтры для каждого дочернего узла.

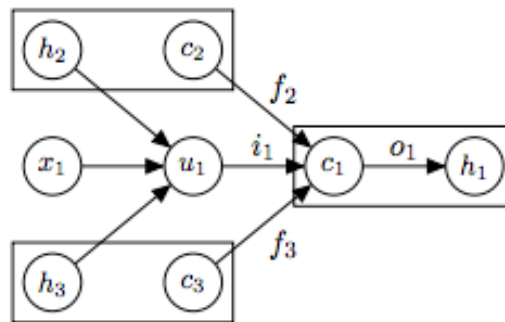


Рис. 3: Tree-LSTM с двумя дочерними узлами [20]

На вход алгоритму подаются заранее подготовленная выборка, состоящая из предложений с представлением в виде дерева. Сначала каждое дочернее слово подается на вход дочерним узлам Tree-LSTM, далее их скрытые состояния h_j , подается на вход родительскому узлу вме-

сте с родительским словом в структуре предложения. Таким образом производятся вычисления по всем ячейкам этой сети. Количество ячеек является гиперпараметром модели и требует валидации при обучении алгоритма. Данная работа показала на момент написания статьи state-of-the-art. Результаты обучения на различных выборках представлены в таблице 3

Модель	SST-1	SST-2	MR
Constituency Tree-LSTM, tuned	51.0	88.0	81.7
Constituency Tree-LSTM, fixed	49.7	87.5	-

Таблица 3: Результаты разных моделей Tree-LSTM [20]

1.2. Семантическое векторное представление

1.2.1. NLP и задача семантического векторного представления слов

Обычная форма векторного представления слов представляет собой one-hot вектор его размерность определена размером словаря и он состоит полностью из нулей кроме одной единственной позиции, при таком представлении можно точно определить слово и последовательность слов, представляя таким образом документ в виде матрицы, можно заметить, что такое представление сильно избыточно и не является вектором в математическом смысле. То есть математические операции над такими векторами не представляют из себя никакой семантики. Проблему векторного представления слов, отображающих семантику, решают множества алгоритмов: word2vec [4], GloVe [15], fastText [5] и т.д. Это некоторое векторное представление слова, размерность которого может быть задана при обучении модели. Такое представление более компактно и отвечает свойству семантической схожести, то есть два близких по смыслу слова в векторной форме обладают большим скалярным произведением. Применение математических операций в действительности может дать какой-то осмысленный результат, например,

$v(King) - v(Man) + v(Woman)$, где $v(w)$ - функция преобразования слова в векторное представление, результат этого выражения будет похож, согласно определенной метрике, на слово $v(Queen)$. Такое векторное представление в области обработки естественного языка (англ. Natural Language Processing) достаточно сильно распространено и показывает значительное улучшение результатов в решении не только задачи классификации текста, но и во многих других.

1.2.2. Описание задачи для слов отсутствующих в обучающей выборке

В данной работе для разработки алгоритма необходимо поставить задачу и определить требования по валидации получаемых векторных представлений слов. Задачу можно сформулировать следующим образом: для любого произвольного текста, необходимо получить векторные представления слов в этом тексте, даже для тех слов, которые отсутствовали при обучении. Полученные векторные представления слов, должны соответствовать определенным характеристикам схожести. Также полученные векторные представления должны улучшать результат какой-либо целевой задачи, использующие векторное представление слов, в зависимости от выбранного способа инициализации embedding слоев в искусственных нейронных сетях и с расширяемым словарем. Можно сказать, что для оценки качества получаемых векторных представлений должна проводиться валидация на выборках, оценивающих схожесть двух векторных представлений слов, и валидацию на целевой задаче. Способы валидации схожести и целевая задача подробнее описана в главе 3.

1.2.3. Модель Mimick

Авторы статьи [16] предлагают решать проблему отсутствующих векторов, как проблему генерации вектора, независимо от того, как они были получены. Они предполагают, что все вектора для слов могут быть вычислены согласно некоторому алгоритму. Обучая модель

на существующем словаре таких векторов, они предполагают, что можно получить вектор для любого произвольного слова по его символам. Для некоторого языка L , словаря $V \subseteq L$ и пред обученных векторов $W \in \mathbb{R}^{|V| \times d}$, где d - размерность векторов и каждый вектор $\{w_i\}_{i=1}^{|V|}$ размечен вектором e_i . Функция $f : L \rightarrow \mathbb{R}^{|V| \times d}$, приближающая значение функции к существующему вектору: $f(w_i) \approx e_i$. Таким образом появляется возможность получать вектора для слов $w_k^* \in L \setminus V$ следующим образом: $f(w_k^*) = e_k^*$. В статье предполагается использовать нейронную сеть глубокого обучения, как функцию f . Архитектура нейронной сети представляет из себя двунаправленную LSTM-сеть (рис. 4).

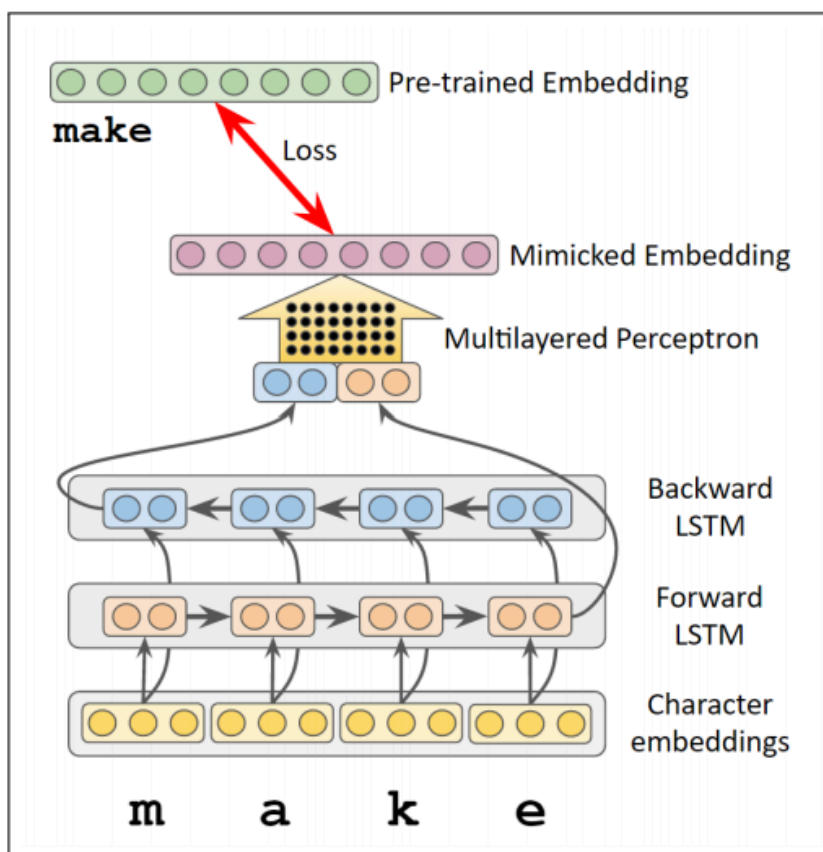


Рис. 4: Архитектура модели [16]

Где входом будут являються символы, генерируемого слова. Последовательность символов $w = \{c_i\}_1^n$ подаются на вход к слоям BiLSTM (forward, backward), как векторные представления символов $\{e_i^{(c)}\}_1^n$. Скры-

тые состояния для направленных вперед и назад слоев h_f^n и h_b^0 заключительные скрытые состояния соединяются и векторное представление слова будет представлено следующей формулой:

$$f(w) = O_T \cdot g(T_h \cdot (h_f^n \oplus h_b^0) + b_h) + b_T, \quad (5)$$

где T_h, b_h, O_T, b_T - параметры полносвязного слоя, а функция g - нелинейная функция активации, \oplus - конкатенация двух векторов. Функция ошибки представляет из себя средне-квадратичную ошибку:

$$L = \| f(w_k) - e_k \|^2$$

1.2.4. Модель Enriching

В данном разделе будет описан метод представленный в статье [5]. Метод описанный в статье предлагает решение проблемы векторного представления слов, отсутствующих в обучающей выборке. Векторные представления слов, получаемые этим алгоритмом, показывают state-of-the-art результат на валидационных задачах, таких как: схожесть слов, схожесть редких слов (англ. WordSim, RareWords). Основная идея метода - получать векторное представление слова, как линейную комбинацию векторных представлений символьных n-грамм. Таким образом пытаются выяснить, как отдельные части слова влияют на общее представление вектора, чтобы в последствии для слов, отсутствующих в выборке, была возможность получить его векторное представление, которое могло отвечать некоторым семантическим свойствам, как, например, слова получаемые алгоритмами word2vec, GloVe и т.п. Основой этого алгоритма стала модель Skip-gram предложенная в статье [4]. В этой статье было проверено предположение о распределении слов в зависимости, от его контекстных слов [9]. Целевая функция имеет вид:

$$\sum_{t=1}^T \sum_{c \in C_t} p(w_c | w_t), \quad (6)$$

где $p(w_c|w_t)$ - вероятность появления слова t в контексте c слова. В случае алгоритма word2vec используется функция softmax:

$$p(w_c|w_t) = \frac{e^{s(w_t, w_c)}}{\sum_{j=1}^{|W|} e^{s(w_t, w_j)}}. \quad (7)$$

Из-за такого способа представления функции вероятности, получается, что существует только одно предсказываемое контекстное слово. Поэтому решают несколько независимых задач, на каждое контекстное слово. Используя технику negative sampling, в нашу целевую функцию добавляются еще слова, которые не принадлежат контексту. В конечном виде целевая функция выглядит следующим образом:

$$\sum_{t=1}^T \left[\sum_{c \in C_t} \log(1 + e^{-s(w_t, w_c)}) + \sum_{n \in N_{t,c}} \log(1 + e^{s(w_t, n)}) \right]. \quad (8)$$

Такая функция параметризована относительно оценивающей функцией s , где в статье [4] было скалярное произведение: $s(w, v) = w^T v$. В описываемом методе использовалась другая оценивающая функция $s(\cdot, \cdot)$. С целью получить информацию о внутренней структуре слова, была предложена функция оценки на основе мешка с символьными n -граммами. Например, слово "какой" при $n = 3$ было бы представлено как: <ка, как, ако, кой, ой>. Символы <, > были добавлены с целью обозначить начало и конец слова. Допустим имеется словарь n -грамм размера G . Для некоторого слова w можно определить множество $\zeta_w \subset \{1, \dots, G\}$ - множество символьных n -грамм для слова w . Пусть z_g - векторные представления для всех n -грамм размера g . Таким образом наше слово теперь представимо в виде множества таких n -грамм:

$$s(w_t, w_c) = \sum_{g \in \zeta_{w_t}} z_g^T w_c. \quad (9)$$

Авторы также приводят решение ряда проблем с требованиями к памяти, но к основной модели они относятся несущественно и носят исключительно оптимизационный характер при реализации алгоритма. Таким образом целевая функция статьи [5], является модифицирован-

ной в форме оценочной функции формулы (8) с оценочной функцией из формулы (19).

1.2.5. Модель Glove

В данном разделе будет описан алгоритм метод получения векторного представления слов, для фиксированного словаря из статьи [15]. Авторы статьи предлагают более детально рассмотреть каким именно образом используется статистика совместной встречаемости в современных алгоритмах, решающих эту задачу. Авторы предлагают новый метод для генерации векторного представления на основе статистики совместной встречаемости собранной со всей обучающей выборки. Дальнейшее описание алгоритма сохранило нотацию авторов статьи [15]. Пусть матрица совместной встречаемости пар слов будет обозначена X , X_{ij} - количество раз слово j встретилось в контексте слова i . Пусть $X_i = \sum_k X_{ik}$ - количество слов, которые были в контексте слова i . Теперь необходимо обозначить вероятность появления слова j в контексте слова i : $P_{ij} = P(j|i) = \frac{X_{ij}}{X_i}$. Для иллюстрации моделируемой вероятности, авторы статьи приводят пример с термодинамическими фазами. Пусть $i = ice$, $j = steam$, $k = solid$ можно ожидать, что отношение $\frac{P_{ik}}{P_{jk}}$ будет высоким. И таким же образом, если слово k будет иметь большее отношение к слову j , то отношение будет меньше. На основе текущего предположения авторы статьи предлагают использовать вероятность основанную на статистике совместной встречаемости. Авторы предлагают обобщенную модель:

$$F(w_i, w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}, \quad (10)$$

где $w \in \mathbb{R}^d$ - векторное представление слова i , а $\tilde{w} \in \mathbb{R}^d$ - векторное представление контекстного слова. При текущем определении функции $F(\cdot, \cdot, \cdot)$ невозможно определить какую либо конкретную функцию, которая решала бы тождество 10. Для решения необходимо ввести ограничения эту функцию. Для этого необходимо ввести ограничение на $\frac{P_{ik}}{P_{jk}}$. Так как векторное пространство линейно, то векторная разница -

это более интуитивный способ отобразить это свойство свойство.

$$F(w_i - w_j, \widetilde{w}_k) = \frac{P_{ik}}{P_{jk}}. \quad (11)$$

Так как F может быть абсолютно всем чем угодно, например, искусственной нейронной сетью, то при такой формулировке можно потерять линейную структуру. Чтобы предотвратить эту проблему, необходимо взять скалярное произведение от аргументов функции:

$$F((w_i - w_j)^T \widetilde{w}_k) = \frac{P_{ik}}{P_{jk}}, \quad (12)$$

что предотвращает смешивание элементов вектора в произвольной модели. Чтобы отобразить симметрию относительно слова и контекстного слова авторы определяют следующее свойство:

$$F((w_i - w_j)^T \widetilde{w}_k) = \frac{F(w_i^T \widetilde{w}_k)}{F(w_j^T \widetilde{w}_k)}, \quad (13)$$

где $F(w_i^T \widetilde{w}_k) = P_{ik} = \frac{X_{ik}}{X_i}$. Таким образом можно выявить, что $F = \exp$. Тогда верно и следующее:

$$w_i^T \widetilde{w}_k = \log P_{ik} = \log X_{ik} - \log X_i. \quad (14)$$

Так как $\log X_i$ не зависит от слова k можно заменить его на дополнительное смещение b_i, \widetilde{b}_k . Основной проблемой такого подхода является тот факт, что при прямом использовании статистика совместной встречаемости оценивается одинаково для всех пар слов. То есть более редкие пары достаточно зашумлены и несут меньше информации. Авторы предлагают следующую функцию ошибки:

$$L = \sum_{i,j=1}^V f(X_{ij}) \cdot (w_i^T \widetilde{w}_j + b_i + \widetilde{b}_j - \log(X_{ij})), \quad (15)$$

где X_{ij} - показатель встречаемости слова j в контексте слова i , f - взвешивающая статистику встречаемости функции, V - размер словаря. Функция f должна отвечать некоторым свойствам:

1. $f(0) = 0$. Если f - непрерывная функция, то лучше рассматривать тождество, как сходящийся предел $\lim_{x \rightarrow 0} f(x) \log x^2$,
2. $f(x)$ неубывающая, чтобы не происходило перевзвешивания совместной встречаемости редких пар слов,
3. $f(x)$ должно быть достаточно низким на больших значениях x .

Существует множество функций, которые отвечают этим свойствам, но авторы статьи привели следующую:

$$f(x) = \begin{cases} \left(\frac{x}{x_{max}}\right)^\alpha & \text{if } x < x_{max}, \\ 1 & \text{otherwise.} \end{cases} \quad (16)$$

2. Разработки

2.1. Определение тональности

Реализация модели из статьи [10] производилась на языке Python и фреймворке по созданию искусственных нейронных сетей tensorflow [21]. Часть кода по предобработке текста: разбиение на слова и замена специальных символов была позаимствована из исходного кода предоставленная авторами статьи для сохранения равных условий. Нами была реализована архитектура нейронной сети из этой статьи и эта реализация была взята за некоторый основной алгоритм для дальнейших улучшений.

В процессе проведения экспериментов не удалось добиться точных результатов из статьи [22], но полученный результат обучения было принято взять за базовый, то есть все дальнейшие результаты оценивались относительно его, для этого было проведено обучение с использованием кросс-валидации, чтобы удостовериться в воспроизводимости результатов нашей реализации. В рамках работы не было цели реализовать алгоритм, который в точности воспроизводит описанную в статье [10] модель. Также для проведения дополнительной проверки были обучены существующие реализации на фреймворке Torch [3], результаты обучения и проверки описаны в главе 5.

2.1.1. Распределение эмоционального окраса в предложениях

В процессе создания нейронной сети были выдвинуты теории о словах, после которых мнение, содержащееся в левой части предложения не имеет смысла из-за второй части. Например, предложения по структуре похожие "X, но Y". В предложениях подобного типа могли встречаться абсолютно разные слова между X и Y, но их свойства сохранялись. Такой тип предложений легко воспринимался в модели RNTN [17], и было важным понять, а верно ли это для CNN, для этого был проведен эксперимент, где по каждому документу из выборки скользящее окно определенного размера, каждое окно было подано на предска-

ние метки класса в обученную модель. После этого строились графики в зависимости от расположения контекста (рис. 5, рис. 6).

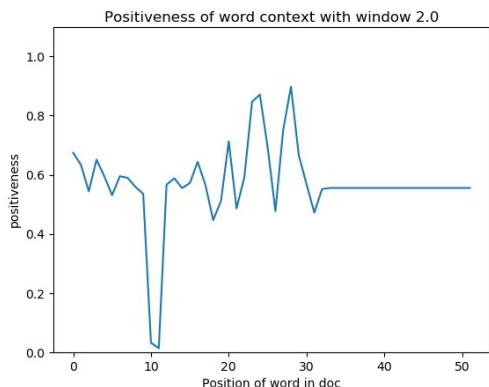


Рис. 5: Полярность документа с контекстом длины 2

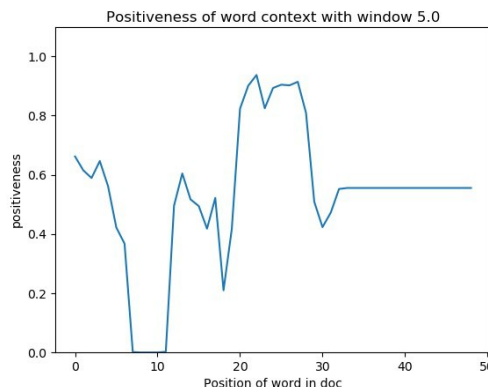


Рис. 6: Полярность документа с контекстом длины 5

Выявление подобных слов позволяло увеличить размер выборки в два раза за счет переворачивания частей предложения, чтобы из "X, но Y", предложение менялось на "Y, но X" с заменой метки класса такого документа. Такой вид модификации мог быть использован для искусственного увеличения объема выборки. Об экспериментах и других способах синтетического увеличения датасета будет сказано далее в разделе 2.1.2 и главе 5

2.1.2. Увеличение выборки: замена синонимов

В процессе реализации и последующего анализа результатов было выявлено, что структура нейронной сети уже достаточно сложна и ее углубление будет приводить лишь к переобучению. Одна из возможных причин такого поведения нейросетей – недостаточный объем выборки. В процессе мозгового штурма были выдвинуты аналогии из области компьютерного зрения, решающие подобную проблему. Одной из таких аналогий стало применение техник искусственного увеличения объемов данных (англ. data augmentation). Эти техники направлены на искусственное увеличение датасета за счет элементарных преобразований

над изображениями, такими как: зеркальное отображение, поворот на 90 градусов и т.п. В случае же текста было принято решение на первом этапе заменять прилагательные и наречия на их синонимы. Синонимы были получены при помощи WordNet [13]. Это некоторый граф слов и их описаний, он заранее обучен и имеется в открытом доступе. С использованием этой техники можно было получить выборку достаточно большого размера, но по некоторым экземплярам обучающей выборки не удавалось найти синонимы для используемых частей речи. Чтобы сеть могла достаточно хорошо обучиться и не переобучалась на конкретных экземплярах, присутствующих в увеличенной выборке, было введено ограничение на количество генерируемых документов по одному экземпляру. В этой работе были рассмотрены несколько вариантов такой генерации с рассмотрением процесса обучения. Такая модификация не затрагивала метку класса. И поэтому такое искусственное увеличение объема выборки не влияло на первоначальную сбалансированность выборки, то есть соотношение документов разных классов было с тем же отношением, что и до искусственного увеличения.

2.1.3. Увеличение выборки: замена слова на усредненное описание

В результате поиска современных статей о векторном представлении слов и текста, был найден один алгоритм, который мог представить весь документ в виде одного вектора [1]. Этот алгоритм был использован в качестве замены какого-то слова в предложении на описание этого слова представленное алгоритмом [1]. Данная модификация не изменяла метки класса и также нейтральна к изменению сбалансированности выборки, то есть свойства выборки должны сохраняться. Также данный алгоритм идет в двух модификациях, по которым были проведены эксперименты. Результаты экспериментов приведены в разделе 5.1.2.

2.1.4. Углубление архитектуры нейронной сети

Все действия, описанные в разделе 2.1.2, были направлены на то, чтобы попытаться искусственно увеличить выборку, для поиска более сложных взаимосвязей, которые могут быть распознаны более глубокими сетями. На примере работ из области компьютерного зрения, можно выявить, что более глубокие сети могут находить более сложные образы. Как следствие такие сети получают более высокий показатель точности. Но совместно с глубиной у нас увеличивается и количество обучаемых параметров, что требует большой обучающей выборки. Также важным знанием в этой области является и то, что каждый сверточный уровень представляет из себя некоторое абстрактное представление и добавление дополнительных слоев носит не случайный характер. Делается предположение, что каждый следующий слой является некоторой комбинацией простых признаков предыдущего слоя, что позволяет в конечном итоге распознавать очень сложные объекты. Подобные рассуждения были использованы в текущей работе, первый слой находил сочетания слов, а следующий сверточный слой находил отношения между сочетаниями слов.

2.2. Векторное представление слов

2.2.1. Модель Mimick и наша модификация

Суть разработки алгоритма лежит в предположении, что для поставленной задачи необходима дополнительная информация, например, о векторном представлении контекстного слова. Было выдвинуто предположение, что добавление контекстного слова при обучении модели *mimick*, позволит получить более общую модель. Была модифицирована функция f из формулы (5) следующим образом:

$$f(w, e_j) = O_T \cdot g(T_h \cdot (h_f^n \oplus h_b^0 \oplus e_j) + b_h) + b_T, \quad (17)$$

где e_j - векторное представление j -го слова, контекст слова w .

2.2.2. Наша модификация GloVe + SISG

Основой разработанного алгоритма является алгоритм GloVe. Его описание может быть найдено в главе 1. В алгоритме из раздела 1.2.4 было предложено обобщенное представление отношения между словом и его контекстом. Так и в алгоритме GloVe существует возможность замены части целевой функции, отвечающее за произведение векторов $w_i^T \tilde{w}_j$. Необходимо модифицировать функцию ошибки из формулы (15) схожим образом:

$$L = \sum_{i,j=1}^V f(X_{ij}) \cdot (s(w_i, \tilde{w}_i) + b_i + \tilde{b}_j - \log(X_{ij})). \quad (18)$$

Функция $s(\cdot, \cdot)$ была адаптирована из формулы 19, но с некоторыми изменениями относительно количества используемых размеров n -грамм. Фактически в функции оценки используется только один фактический размер символьных n -грамм. Пусть G_w - все векторные представления n -грамм слова w для заданного n , тогда оценивающую функцию можно записать в следующем виде:

$$s(w_t, w_c) = \sum_{g \in G_{w_t}} g^T w_c. \quad (19)$$

Предполагается, что за счет статистики совместной встречаемости контекста слова, можно получить качественные представления для символьных n -грамм, с помощью которых могут быть получены слова, отсутствующие в обучающей выборке.

Язык	Токены	V
CS	1m	46k
DE	1m	36k
EN	1m	17k
ES	1m	27k
FR	1m	25k
RU	1m	62k

Таблица 4: Статистика набора данных DATA-1M [2]

3. Способы валидации

3.1. Валидация на моделировании английского языка

В области векторного представления слов использование способа валидации на целевой задаче достаточно распространено. Такой способ валидации также демонстрирует некоторый способ применения техники transfer learning - использование результатов алгоритмов из другой задачи в решении целевой задачи. Сгенерированные вектора являются результатом решения нашей задачи, поэтому ключевым критерием оценки качества является оценка на задаче по обработке естественного языка. Авторы статьи [5] проводят такую валидацию на задаче языкового моделирования английского языка. Авторы статьи предлагают обучить искусственную нейронную LSTM сеть, решающую эту задачу. Качество обученной сети оценивается по метрике perplexity. Для проведения качественной и честной валидации, используется заранее подготовленная и обработанная выборка, где есть фиксированная тестовая и обучающая часть. Авторы статьи [5] позаимствовали выборку из работы [2]. Выборка представляет из себя набор текста разделенный по языкам. Базовые статистики по этому набору данных можно посмотреть в таблице 3.1. В каждом языке были отобраны около миллиона токенов, которые были заранее получены описанным в статье скриптом. Таким образом мы можем исключить отклонения из-за какой-либо обработки данных и другому разбиению текста на токены. Авторы статьи [5] не

приводили результаты на английском языке, поэтому нам необходимо было обучить эту сеть на этой выборке и также обучить сеть с различными инициализациями, подробнее об этом описано в разделе 5.2.2. Многие исследователи в области генерации векторных представлений слов используют и другие целевые задачи при валидации, но в рамках текущей работы они не были рассмотрены.

4. Детали реализации

4.1. Детали реализации методов векторного представления

В данном разделе содержится описание используемых языков программирования, библиотек и алгоритмов, которые позволили реализовать эту работу. Здесь же будут приведены базовые гиперпараметры, которые использовались при обучении искусственных нейронных сетей. Все разработанные алгоритмы и системы реализовывались на языке Python. Для обработки больших объемов текстов использовались скрипты на языке Perl и Bash. Предварительная обработка текста была написана на языке программирования Bash, часть обработки текстов для валидации была реализована на языке Python с использованием библиотек: nltk [11], Word Embeddings Benchmarks. Для реализации искусственных нейронных сетей использовался фреймворк глубинного обучения Tensorflow. Также важно отметить использование библиотеки tflearn, позволяющая описывать архитектуру нейронной сети в декларативном стиле. Все обучение производилось на CPU на процессоре Intel Core i7 2.8 Ghz и 16 Gb RAM. Размерность векторов была выбрана на основе максимально доступных данных, на которых была возможность проверить достоверность полученных результатов авторами статей, поэтому размерность вектора равнялась $d = 300$.

4.1.1. Mimick с контекстным вектором

При реализации данной сети возникли вопросы относительно того, каким образом обучать на большом объеме данных, ведь в оригинальной статье сложность обучения равнялась размеру словаря умноженному на сложность сети. И в их программной системе не было использовано алгоритмов оптимизирующих работу с большим объемом данных. Для решения этой задачи были использованы эффективные способы представления матриц с использованием библиотек hdf5, dask. Эти библиотеки позволяли хранить большие данные в памяти только тот ку-

сок матрицы, который использовался только в текущий момент. Данное улучшение программной реализации позволило сократить время обучение и в принципе сделало обучение сети возможным. Гиперпараметры архитектуры искусственной нейронной сети описанной в разделе (5.2.1) были от части адаптированы из работ [5] и [16]. Таким образом были получены гиперпараметры: количество эпох обучения (4), количество ячеек в LSTM-слое (50), а также это дало основу для различных предположений относительно гиперпараметров размера выборки для одной итерации алгоритма и параметр learning rate (0.001). Векторные представления слов для валидации были сэмплированы из обучающей выборки. Из обучающей выборки по составленному словарю случайным образом выбиралось контекстное слово и подавалось на вход этой модели, полученные векторные представления записывались в словарь и таким образом был получен финальный словарь, на котором проводилась вся валидация. Также важно отметить, что для слов отсутствующих в обучающей выборке был составлен дополнительный словарь, где контекстное слово было заменено на единичный вектор. Также важно отметить, что это было сделано только для тех выборок, где у анализируемого слова не было контекстного слова. Эта проблема была описана в разделе (5.2.1).

4.1.2. GloVe + SISG

При реализации данной системы основной точкой интереса являлась модификация целевой функции, а не эффективное написание подсчета статистики встречаемости слов. Поэтому была использована библиотека с открытым исходным кодом (<https://github.com/maciejkula/glove-python/tree/master/glove>), где подсчет данной статистики был реализован на языке C++. Архитектуру сети реализовали на чистой библиотеке Tensorflow с использованием off-heap хранилища для таблицы n-грамм, и разреженной матрицы совместной встречаемости. В целях дополнительной оптимизации программного алгоритма были использованы библиотеки по работе с off-heap хранилищем dask. Также возникали проблемы при обучении самой сети, ведь приходилось об-

новлять векторные представления в достаточно большой таблице векторных представлений символьных n-грамм. Для этого использовалось только частичное обновление параметров только участвующих при вычислении целевой функции и также было использован дополнительный сервис для распределения всей таблицы на несколько вычислительных машин. При обучении такой сети использовались аналогичные гиперпараметры из статьи [15]. Количество эпох при обучении равнялось 50, а размер контекста 5. Все остальные параметры подробно описаны в самой статье. Размер символьных n-грамм также выбирался исходя из лучших результатов из статьи [5].

5. Результаты экспериментов

5.1. Результаты экспериментов разработанных методов определения тональности

В данной главе описаны эксперименты по искусственному увеличению выборки, влияние этого на обучение модели и результаты проводимых экспериментов. Обучение нейронных сетей производилось без использования GPU и обучалось на стационарном компьютере. Обучение модели занимало от 30 минут до 3 часов, в случае использования CNN, и около одной недели при обучении модели Tree-LSTM.

5.1.1. Оценка методик на существующих реализациях

Всего было проведено два типа экспериментов с искусственно увеличенной выборкой:

1. обучение разных моделей;
2. углубление сети.

Весь эксперимент заключался в том, чтобы искусственно увеличить объем обучающей выборки и потом использовать его, для обучения той же самой модели с абсолютно идентичными гипер параметрами. Лучшие результаты можно посмотреть в таблице 5.

Модель	без увеличения	с увеличением	тип проверки
Tree-LSTM	81.7	82.58, 82.46	тестовая выборка
CNN (реализация torch)	81.3	81.8	кросс- валидация
CNN наша реализация	80.5	81	кросс- валидация
CNN 2conv 1x2	-	79	кросс- валидация
CNN 2conv 2x2	-	80	кросс- валидация

Таблица 5: Результаты проведенных экспериментов: последние две строчки с таблице, эксперимент по добавлению еще одного сверточного слоя; реализация harvard <https://github.com/harvardnlp/sent-conv-torch>

Полученная точность больше базового на 0.5, результат хоть и не значительный, но как показали эксперименты, можно брать любой алгоритм и производить его обучение с использованием этой методики, то есть берется абсолютно любой алгоритм, берется выборка, на котором он обучался, производится его искусственное увеличение, и получается результат лучше на 0.5%. Чтобы удостовериться в корректности нашего предположения был проведен еще ряд экспериментов, как на существующих имплементациях алгоритма, так и других алгоритмов, которые вышли после начала текущей работы. Этот эксперимент был повторен на реализации <https://github.com/harvardnlp/sent-conv-torch>, на алгоритме Tree-LSTM [20].

Второй эксперимент был по усложнению архитектуры сети с добавлением еще одного слоя свертки, для получения более сложных шаблонов (описано в разделе 2.1.4). В таблице 5 представлены результаты экспериментов двух сетей, на кросс валидации получен результат

на двух версиях сети, где отличается гиперпараметр размера свертки, один был размерностью 1×2 , что предположительно, должно было дать шаблоны находящие сочетания соседних словосочетаний одной длины. И был рассмотрен другой вариант со сверткой 2×2 , где находились сочетания не только сочетаний одинаковой длины, но и сочетаний разных длины. Как видно из таблицы, данный эксперимент не принес успеха и базовый результат модели.

5.1.2. Эксперимент: замена слова на его описание

Здесь приведено описание эксперимента по еще одному способу искусственного увеличения выборки, описанного в разделе 2.1.3. Эксперимент производился для одной модели из статьи [10]. На нашей реализации этого алгоритма. Представление документа из статьи [1] идет в двух вариантах: простое усреднение всех векторов-слов и усреднение с вычитание главных компонент из алгоритма SVD [7]. Результаты описаны в таблице 6

Модель	no-SVD	SVD
CNN наша реализация	75	76

Таблица 6: Результат экспериментов по замене слова на его усредненное описание

5.2. Результаты экспериментов разработанных методов векторного представления слов

5.2.1. Результаты модификации модели Mimick

В данном разделе будут описаны проведенные эксперименты по обучению и валидации алгоритма описанного в разделе 2.2.1. При обучении этой сети возникали различные проблемы с валидацией реализации основного алгоритма из статьи [16]. Реализация этого алгоритма была несовместимой с машиной на которой проводились эксперименты и было необходимо использовать виртуализацию для решения этой

проблемы. Результаты валидации на нескольких выборках по задаче векторного представления представлены в таблице 7.

Модель	Rare Words	WordSim353
fastText	0.43	0.73
Mimick	0.27	0.17
Модификация Mimick	0.29	0.21

Таблица 7: Результат экспериментов по получению векторных представлений из различных моделей

При проведении алгоритма выявилась проблема по составлению словаря используемого при валидации. Так как в описанной нами модификации при составлении вектора должно участвовать контекстное векторное представление слова, было необходимо каким-то образом инициализировать это слово из обучающей выборки или же использовать некоторый случайный вектор. Был выбран подход по составлению словаря по обучающей выборке, если же для какого-то слова не существовало контекстного слова или же само слово не присутствовало в обучающей выборке, то производилась случайная инициализация этого вектора. По этим результатам можно сказать, что как и работа [16], так и наша разработка не являются достаточным результатом, чтобы говорить о state-of-the-art результате, но важно отметить, что обучение этих моделей было проведено на основе какого-то словаря векторных представлений и логично предположить, что абсолютное улучшение результата здесь невозможно и важно оценивать не способность получать векторные представления лучше тех, на которых проходило обучение, а смотреть на способность алгоритма восстановить семантику этих векторных представлений. Таким образом можно сказать, что алгоритм [16] и наша реализация, могут достаточно хорошо распознавать семантику по существующим словарям. Также важно отметить, что результат полученный в данной работе превосходит результат [16].

5.2.2. Результаты GloVe + SISG

В данном разделе описаны проведенные эксперименты по обучению моделей из разделов (1.2.4, 2.2.2). При проведении экспериментов возникали некоторые проблемы по обучении сети. Эти проблемы были описаны в деталях реализации. По проведенным экспериментам по обучении искусственных нейронных сетей была проведена валидация на целевой задаче описанной в главе 3. Результаты различных векторных представлений на целевой задаче и на задачах по семантической схожести векторных представлений представлены в таблице 8.

Модель	LM (PPL)	Rare Words	WordSim353
random init	241	-	-
fastText	233	0.47	0.73
GloVe	229	0.478	0.759
наша модификация	225	0.48	0.757

Таблица 8: Результат экспериментов по получению векторных представлений из различных моделей на целевой задаче

Анализируя полученные результаты можно сделать вывод, что разработанный нами алгоритм способен получать получать state-of-the-art результат на целевой задаче и на задаче семантической схожести слов. Также важно отметить, что помимо этого, такая модификация алгоритма позволяет получать векторные представления для слов, отсутствующих в словаре.

Заключение

В данной магистерской диссертации были изучены и реализованы работы последних лет в области задачи определения тональности текста и векторного представления слов. Были разработаны алгоритмы машинного обучения и реализованы в виде программных систем способных определять тональность текста и способных получать векторные представления слов. В рамках диссертационной работы был разработан и реализован алгоритм векторного представления слов для слов, отсутствующих в словаре, были проведены эксперименты по обучению искусственных нейронных сетей для решения этой задачи и проведена валидация полученных векторных представлений.

Список литературы

- [1] Arora Sanjeev, Liang Yingyu, Ma Tengyu. A simple but tough-to-beat baseline for sentence embeddings. — 2016.
- [2] Botha Jan, Blunsom Phil. Compositional morphology for word representations and language modelling // International Conference on Machine Learning. — 2014. — P. 1899–1907.
- [3] Collobert R., Kavukcuoglu K., Farabet C. Torch7: A Matlab-like Environment for Machine Learning // BigLearn, NIPS Workshop. — 2011.
- [4] Distributed representations of words and phrases and their compositionality / Tomas Mikolov, Ilya Sutskever, Kai Chen et al. // Advances in neural information processing systems. — 2013. — P. 3111–3119.
- [5] Enriching word vectors with subword information / Piotr Bojanowski, Edouard Grave, Armand Joulin, Tomas Mikolov // arXiv preprint arXiv:1607.04606. — 2016.
- [6] Goller Christoph, Kuchler Andreas. Learning task-dependent distributed representations by backpropagation through structure // Neural Networks, 1996., IEEE International Conference on / IEEE. — Vol. 1. — 1996. — P. 347–352.
- [7] Golub Gene H, Reinsch Christian. Singular value decomposition and least squares solutions // Numerische mathematik. — 1970. — Vol. 14, no. 5. — P. 403–420.
- [8] Gradient-based learning applied to document recognition / Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner // Proceedings of the IEEE. — 1998. — Vol. 86, no. 11. — P. 2278–2324.
- [9] Harris Zellig S. Distributional structure // Word. — 1954. — Vol. 10, no. 2-3. — P. 146–162.

- [10] Kim Yoon. Convolutional neural networks for sentence classification // arXiv preprint arXiv:1408.5882. — 2014.
- [11] Loper Edward, Bird Steven. NLTK: The natural language toolkit // Proceedings of the ACL-02 Workshop on Effective tools and methodologies for teaching natural language processing and computational linguistics-Volume 1 / Association for Computational Linguistics. — 2002. — P. 63–70.
- [12] McAuley Julian, Leskovec Jure, Jurafsky Dan. Learning attitudes and attributes from multi-aspect reviews // Data Mining (ICDM), 2012 IEEE 12th International Conference on / IEEE. — 2012. — P. 1020–1025.
- [13] Miller George A. WordNet: a lexical database for English // Communications of the ACM. — 1995. — Vol. 38, no. 11. — P. 39–41.
- [14] Parsing natural scenes and natural language with recursive neural networks / Richard Socher, Cliff C Lin, Chris Manning, Andrew Y Ng // Proceedings of the 28th international conference on machine learning (ICML-11). — 2011. — P. 129–136.
- [15] Pennington Jeffrey, Socher Richard, Manning Christopher. Glove: Global vectors for word representation // Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). — 2014. — P. 1532–1543.
- [16] Pinter Yuval, Guthrie Robert, Eisenstein Jacob. Mimicking word embeddings using subword RNNs // arXiv preprint arXiv:1707.06961. — 2017.
- [17] Recursive deep models for semantic compositionality over a sentiment treebank / Richard Socher, Alex Perelygin, Jean Y Wu et al. // Proceedings of the conference on empirical methods in natural language processing (EMNLP) / Citeseer. — Vol. 1631. — 2013. — P. 1642.

- [18] SemEval-2016 Task 5: Aspect Based Sentiment Analysis / Maria Pontiki, Dimitrios Galanis, Haris Papageorgiou et al. // Proceedings of the 10th International Workshop on Semantic Evaluation. — 2016.
- [19] Symposium Sentiment Analysis. Sentiment Analysis Symposium // Sentiment Analysis Symposium. — 2016.
- [20] Tai Kai Sheng, Socher Richard, Manning Christopher D. Improved semantic representations from tree-structured long short-term memory networks // arXiv preprint arXiv:1503.00075. — 2015.
- [21] Tensorflow: Large-scale machine learning on heterogeneous distributed systems / Martin Abadi, Ashish Agarwal, Paul Barham et al. // arXiv preprint arXiv:1603.04467. — 2016.
- [22] Zhang Ye, Wallace Byron. A Sensitivity Analysis of (and Practitioners' Guide to) Convolutional Neural Networks for Sentence Classification // arXiv preprint arXiv:1510.03820. — 2015.

А. Исходный код

А.1. Определение тональности

```
1 import gzip
2 import os
3 import pickle
4
5 import numpy as np
6 import tensorflow as tf
7 import tflearn
8 from gensim.models.keyedvectors import \
9     KeyedVectors
10 from nltk.corpus import wordnet as wn
11 from nltk.tag import pos_tag
12 from sklearn.decomposition.truncated_svd import \
13     TruncatedSVD
14 from sklearn.model_selection import KFold
15 from sklearn.utils import shuffle
16 from tflearn import conv_1d
17 from tflearn.datasets import imdb
18 from tflearn.layers.core import input_data, \
19     dropout, fully_connected, flatten
20 from tflearn.layers.estimator import regression
21 from tflearn.layers.merge_ops import merge
22
23 from data_augmentation import wordnet_pos_code
24 from data_helpers import raw2csv
25
26
27 def load_embeddings(w2v, vocab, new_emb=None):
28     """
29     :param w2v: Word2Vec gensim instance
30     :param new_emb:
31     :param vocab:
32     :return:
33     """
34     if new_emb is None:
35         new_emb = dict()
36         # initial matrix with random uniform TODO: replace 300 with arg
37         initW = np.random.uniform(-0.25, 0.25,
38                                   (len(vocab), 300))
39     for word, idx in vocab._mapping.iteritems():
40         if word in new_emb:
41             initW[idx] = new_emb[word]
42             break
```



```

43     if word in w2v:
44         initW[idx] = w2v[word]
45     return initW
46
47
48 def augment_with_avg(df, vocab, w2v):
49     """
50     :param w2v:
51     :param df:
52     :param vocab:
53     :return:
54     """
55
56     def definition_to_word_vector(definition,
57                                   with_pca=False):
58         result = []
59         for token in definition:
60             if token in w2v:
61                 result.append(w2v[token])
62         a = np.array(result)
63         avg = np.average(a, axis=0)
64         if with_pca:
65             svd = TruncatedSVD()
66             svd.fit(a)
67             pc = svd.components_
68             return avg - avg.dot(pc.T) * pc
69         else:
70             return avg
71
72     out_emb = dict()
73     out_data, out_labels = [], []
74     for _, row in df.iterrows():
75         document, label = row.drop(
76             ['positive', 'negative'], row[
77                 ['positive',
78                 'negative']]
79         document = document.as_matrix()
80         raw_doc = []
81         non_unk_words = []
82         for i, w_id in enumerate(document):
83             word = vocab.reverse(w_id)
84             if w_id != 0:
85                 non_unk_words.append((i, word))
86             raw_doc.append(word)
87         num_of_words = len(non_unk_words)
88         cnt = 0

```

```

89     res = [document, document]
90     while cnt < num_of_words ** 2:
91         cnt += 1
92         selected_word = np.random.randint(0,
93                                         len(
94                                             non_unk_words))
95         orig_idx, aug_word = non_unk_words[
96             selected_word]
97         aug_tag = None
98         for w, t in pos_tag(raw_doc):
99             if w == aug_word:
100                 aug_tag = wordnet_pos_code(t)
101                 break
102
103         new_word = aug_word + '_mavg'
104         syns = wn.synsets(aug_word, aug_tag)
105         if len(syns) > 0:
106             definition = syns[
107                 0].definition().split(' ')
108             wv, _ = definition_to_word_vector(
109                 definition)
110             out_emb[new_word] = wv
111             vocab.freeze(False)
112             vocab.add(new_word)
113             copy = document.copy()
114             copy[orig_idx] = vocab.get(
115                 new_word)
116             res[1] = copy
117             break
118         out_data += res
119         out_labels += [label.as_matrix(),
120                       label.as_matrix()]
121
122     return np.array(out_data), np.array(
123         out_labels), vocab, out_emb

```

```

1  from itertools import chain, product
2
3  import numpy as np
4  import pandas as pd
5  from nltk import pos_tag_sents
6  from nltk import word_tokenize
7  from nltk.corpus import wordnet as wn
8  from pip._vendor.distlib._backport import shutil
9  from tflearn.data_utils import VocabularyProcessor
10 from sklearn.utils import shuffle
11

```

```

12 from data_helpers import load_data_and_labels
13
14
15 def synonyms(word, tag):
16     try:
17         synonyms = wn.synsets(word, tag)
18         lemmas = set(
19             chain.from_iterable(
20                 [syn.lemma_names() for syn in
21                  synonyms]))
22         return list(lemmas)
23     except Exception:
24         return []
25
26
27 def wordnet_pos_code(tag):
28     if tag.startswith('NN'):
29         return wn.NOUN
30     elif tag.startswith('VB'):
31         return wn.VERB
32     elif tag.startswith('JJ'):
33         return wn.ADJ
34     elif tag.startswith('RB'):
35         return wn.ADV
36     else:
37         return ""
38
39
40 def my_tokenizer(raw_docs):
41     for raw_doc in raw_docs:
42         splitted = raw_doc.split(' ')
43         yield splitted
44
45
46 def load_docs_as_word_ids(posfile, negfile,
47                            vocab_pkl):
48     """
49     :param posfile:
50     :param negfile:
51     :param vocab_pkl:
52     :return: :type tuple(:type np.ndarray :type np.ndarray :type VocabularyProcessor)
53     """
54     import pandas as pd
55     pos = pd.read_csv(posfile,
56                       header=None, ).as_matrix()
57     neg = pd.read_csv(negfile,

```

```

58         header=None, ).as_matrix()
59 data = np.concatenate((pos, neg), axis=0)
60 pos_labels = [[0, 1] for _ in pos]
61 neg_labels = [[1, 0] for _ in neg]
62 labels = np.concatenate(
63     [pos_labels, neg_labels], axis=0)
64 data, labels = shuffle(data, labels)
65 return data, labels, VocabularyProcessor.restore(
66     vocab_pkl)
67
68
69 def load_train_mr_data():
70     return load_docs_as_word_ids(
71         'data/train/rt-polarity.pos',
72         'data/train/rt-polarity.neg',
73         'data/vocab_original.pkl')
74
75
76 def load_train_imdb_data():
77     return load_docs_as_word_ids(
78         'imdb_data/train/imdb.pos',
79         'imdb_data/train/imdb.neg',
80         'imdb_data/vocab_original.pkl')
81
82
83 def augment(load_data_fn, augmentation_path):
84     data, labels, vocab = load_data_fn()
85     MAX_DOCS = 2
86     shutil.rmtree(augmentation_path,
87                   ignore_errors=True)
88     import os
89     os.mkdir(augmentation_path)
90
91     def reverse(documents):
92         print "reverse doc"
93         for _doc in documents:
94             out = []
95             for class_id in _doc:
96                 if class_id == 0:
97                     continue
98                 out.append(
99                     vocab.vocabulary__.reverse(
100                         class_id))
101             yield out
102
103     tagged_sent = pos_tag_sents(reverse(data))

```

```

104 print "tagged sent"
105 doc_num = 0
106 total = len(tagged_sent)
107 for doc, label in zip(tagged_sent, labels):
108     print "process document {} of {} with doc len {}".format(
109         doc_num,
110         total,
111         len(doc))
112 doc_word_variants = []
113 total_varian_count = 1
114 for word, tag in doc:
115     if total_varian_count > 2000:
116         print "exit: too many variants"
117         break
118     variants = [word]
119     wn_tag = wordnet_pos_code(tag)
120     if wn_tag is wn.ADJ: # or wn_tag is wn.ADV:
121         syns = synonyms(word, wn_tag)[:6]
122         variants += syns
123         total_varian_count *= len(
124             variants)
125     doc_word_variants.append(
126         list(set(variants)))
127 docs = []
128 for new_doc in product(
129     *doc_word_variants):
130     try:
131         raw_doc = ' '.join(new_doc)
132         docs.append(raw_doc)
133     except Exception:
134         continue
135 docs = shuffle(docs)[:MAX_DOCS]
136 while len(docs) < MAX_DOCS:
137     docs *= 2
138 docs = docs[:MAX_DOCS]
139 # need to unfreeze vocab for dynamic fitting of it
140 vocab.vocabulary_.freeze(False)
141 print "generate {} number of documents from one".format(
142     len(docs) - 1)
143 docs = [doc_ for doc_ in
144         vocab.fit_transform(docs)]
145 df = pd.DataFrame(docs, dtype=np.uint32)
146 f = augmentation_path + '/rt-polarity.pos' if \
147     label[
148     1] else augmentation_path + '/rt-polarity.neg'
149 df.to_csv(f, index=False, header=False,

```

```

150         mode='a+')
151     doc_num += 1
152     vocab.save(augmentation_path + '/vocab.pkl')

1  import os
2  import pickle
3  import re
4
5  import numpy as np
6  import pandas as pd
7  from tflearn.data_utils import VocabularyProcessor
8
9
10 class DataSet(object):
11     def __init__(self, folder, filename,
12                 extensions=('pos', 'neg')):
13         self._folder = folder
14         self._filename = filename
15         self._extensions = extensions
16         data, labels, vocab = self.__preprocess_data()
17         self.data = data
18         self.labels = labels
19         self.vocab__ = vocab
20
21     def save(self, filepath):
22         with open(filepath, 'ab+') as f:
23             pickle.dump(self, f)
24
25     @staticmethod
26     def load(filepath):
27         """
28         :param filepath: path to pickle file of :type DataSet instance
29         :return: instance of :type DataSet
30         """
31         with open(filepath, 'r') as f:
32             return pickle.load(f)
33
34     def word_index(self, word):
35         return self.vocab_.vocabulary_.get(word)
36
37     def __preprocess_data(self):
38         positive_filename = self._filename + '?' + \
39             self._extensions[0]
40         negative_filename = self._filename + '?' + \
41             self._extensions[1]
42         positive_filepath = os.path.join(
43             self._folder, positive_filename)

```

```

44     negative_filepath = os.path.join(
45         self._folder, negative_filename)
46     return preprocess(positive_filepath,
47                       negative_filepath,
48                       save_vocab=False)
49
50
51 def __generate_dataset_loader(folder, filename):
52     def load_from_dump(dump_path):
53         if os.path.exists(dump_path):
54             dataset = DataSet.load(dump_path)
55         else:
56             dataset = DataSet(folder, filename)
57
58             dataset.save(dump_path)
59     return dataset
60
61     return load_from_dump
62
63
64 polarity_ds = __generate_dataset_loader('data',
65                                       'rt-polarity')
66
67
68 def clean_str(string):
69     """
70     Tokenization/string cleaning for all datasets except for SST.
71     Original taken from https://github.com/yoonkim/CNN_sentence/blob/master/process_data.py
72     """
73     string = re.sub(r"[^A-Za-z0-9()!?'\"'"]", " ",
74                   string)
75     string = re.sub(r"\s", " \s", string)
76     string = re.sub(r"\ve", " \ve", string)
77     string = re.sub(r"n\t", " n\t", string)
78     string = re.sub(r"\re", " \re", string)
79     string = re.sub(r"\d", " \d", string)
80     string = re.sub(r"\ll", " \ll", string)
81     string = re.sub(r",", " , ", string)
82     string = re.sub(r"!", " ! ", string)
83     string = re.sub(r"\(", " \(", string)
84     string = re.sub(r"\)", " \)", string)
85     string = re.sub(r"?", " \? ", string)
86     string = re.sub(r"s{2,}", " ", string)
87     return string.strip().lower()
88
89

```

```

90 def load_data_and_labels(positive_data_file,
91                          negative_data_file):
92     """
93     Loads MR polarity data from files, splits the data into words and generates labels.
94     Returns split sentences and labels.
95     """
96     # Load data from files
97     positive_examples = list(
98         open(positive_data_file, "r").readlines())
99     positive_examples = [s.strip() for s in
100                         positive_examples]
101     negative_examples = list(
102         open(negative_data_file, "r").readlines())
103     negative_examples = [s.strip() for s in
104                         negative_examples]
105     # Split by words
106     x_text = positive_examples + negative_examples
107     x_text = [clean_str(sent) for sent in x_text]
108     # Generate labels
109     positive_labels = [[0, 1] for _ in
110                       positive_examples]
111     negative_labels = [[1, 0] for _ in
112                      negative_examples]
113     y = np.concatenate(
114         [positive_labels, negative_labels], 0)
115     return x_text, y
116
117
118 def space_split_tokenizer(iterator):
119     for iterate in iterator:
120         yield iterate.split(' ')
121
122
123 def build_vocab(data):
124     max_doc_length = max(
125         map(len, space_split_tokenizer(data)))
126     vocab_processor = VocabularyProcessor(
127         max_doc_length,
128         tokenizer_fn=space_split_tokenizer)
129     return vocab_processor
130
131
132 def preprocess(positive_path, negative_path,
133               save_vocab=True,
134               vocab_path='dataset_vocob.pkl'):
135     data, labels = load_data_and_labels(

```



```

136     positive_path, negative_path)
137 vocab = build_vocab(data)
138 docs = np.array(
139     [w for w in vocab.fit_transform(data)])
140 if save_vocab:
141     vocab.save(vocab_path)
142 return docs, labels, vocab
143
144
145 def raw2csv(
146     file_path='data/rt-polarity_data_labels.csv'):
147     data, labels, v = preprocess(
148         'data/rt-polarity.pos',
149         'data/rt-polarity.neg',
150         vocab_path='data/vocab_original.pkl')
151     frame = pd.DataFrame(
152         np.concatenate([data, labels], axis=1))
153     columns = [str(i) for i in
154                xrange(len(frame.columns))]
155     columns[-1] = 'positive'
156     columns[-2] = 'negative'
157     frame.columns = columns
158     frame.to_csv(file_path, header=False,
159                 index=False)
160     return frame, v
161
162
163 def save_separate_files(tt, type):
164     tt[tt['positive'] == 1].drop(
165         ['positive', 'negative'], axis=1).to_csv(
166         "data/{}/rt-polarity.pos".format(type),
167         index=False, header=False, mode='w+')
168     tt[tt['positive'] == 0].drop(
169         ['positive', 'negative'], axis=1).to_csv(
170         "data/{}/rt-polarity.neg".format(type),
171         index=False, header=False, mode='w+')

```



```

1 import gzip
2 import math
3 import os
4 import pickle
5 import time
6
7 import numpy as np
8 import pandas as pd
9 import tensorflow as tf
10 import tflearn

```

```

11 from gensim.models.keyedvectors import \
12     KeyedVectors
13 from gensim.models.word2vec import Word2Vec, \
14     LineSentence
15 from sklearn.model_selection import KFold
16 from sklearn.utils import shuffle
17 from tflearn import conv_1d
18 from tflearn.datasets import imdb
19 from tflearn.layers.conv import global_max_pool, \
20     conv_2d
21 from tflearn.layers.core import input_data, \
22     dropout, fully_connected, flatten
23 from tflearn.layers.estimator import regression
24 from tflearn.layers.merge_ops import merge
25 from data_augmentation import \
26     load_docs_as_word_ids, my_tokenizer, augment
27 from data_helpers import raw2csv, \
28     save_separate_files
29
30
31 def load_embeddings(filename, vocab):
32     """
33     :param filename:
34     :type vocab: tensorflow.contrib.learn.python.learn.preprocessing.categorical_vocabulary.CategoricalVocabulary
35     :param vocab:
36     :return:
37     """
38     w2v = KeyedVectors.load_word2vec_format(
39         filename, binary=True)
40     # initial matrix with random uniform TODO: replace 300 with arg
41     initW = np.random.uniform(-0.25, 0.25,
42                               (len(vocab), 300))
43     for word, idx in vocab._mapping.iteritems():
44         if word in w2v:
45             initW[idx] = w2v[word]
46     return initW
47
48
49 def load_imdb_vocab(path="imdb.dict.pkl"):
50     path = imdb.get_dataset_file(
51         path, "imdb.dict.pkl",
52         "http://www.iro.umontreal.ca/~lisa/deep/data/imdb.dict.pkl.gz")
53     if path.endswith(".gz"):
54         f = gzip.open(path, 'rb')
55     else:
56         f = open(path, 'rb')

```

```

57
58 vocab = pickle.load(f)
59 return vocab
60
61
62 def get_network(
63     vocab_size,
64     max_document_length,
65     filter_sizes=(2, 3, 4, 5),
66     num_filters=100,
67     dropout_keep_prob=0.5,
68     n_classes=2,
69     embedding_dim=300
70 ):
71     pass
72
73
74 def train_embed(embedding_dim, vocab):
75     model = Word2Vec(
76         LineSentence('data/rt-data.txt'),
77         embedding_dim, window=10, min_count=1)
78     initW = np.random.uniform(-0.25, 0.25, (
79         len(vocab), embedding_dim))
80     for word, idx in vocab._mapping.iteritems():
81         if word in model:
82             initW[idx] = model[word]
83     return initW
84
85
86 def learn():
87     from data_augmentation import my_tokenizer
88     folds = 10
89     k_fold = KFold(folds, shuffle=True)
90     file_path = 'data/rt-polarity_data_labels.csv'
91     df, _ = raw2csv(file_path)
92     for train_indicies, test_indicies in k_fold.split(
93         df):
94         # save_separate_files(df.loc[train_indicies, :], 'train')
95         # save_separate_files(df.loc[test_indicies, :], 'test')
96         # augment()
97         X_train, y_train, vocab = load_docs_as_word_ids(
98             'data/train/rt-polarity.pos',
99             'data/train/rt-polarity.neg',
100             'data/augmentation/vocab.pkl')
101         X_test, y_test, _ = load_docs_as_word_ids(
102             'data/test/rt-polarity.pos',

```

```

103     'data/test/rt-polarity.neg',
104     'data/augmentation/vocab.pkl')
105 X_train, y_train = shuffle(X_train,
106                             y_train)
107 X_test, y_test = shuffle(X_test, y_test)
108 max_document_length = X_train.shape[1]
109 vocab = vocab.vocabulary_
110
111 print "loaded vocab, size: {}".format(
112     len(vocab))
113 vocab_size = len(vocab)
114
115 print "Loaded data"
116 print "Max document length {}".format(
117     max_document_length)
118
119 dropout_keep_prob = 0.5
120 batch_size = 64
121 n_classes = 2
122 embedding_dim = 300
123 filter_sizes = [2, 3, 4]
124 num_filters = 100
125 # snapshot_every = 200
126 folds = 3
127 # max_document_length = 280
128 multichannel = False
129 nonstatic = True
130
131 def get_run_name():
132     mult = "M"
133     stat = "S"
134     nonstat = "NS"
135     params = "dkp{}_filters{}_nb_filters{}_{}_{}" \
136         .format(dropout_keep_prob,
137                 "_".join(
138                     map(lambda b: str(b),
139                         filter_sizes)),
140                 num_filters, time.time())
141     name = ""
142     if multichannel:
143         name += mult
144     elif nonstatic:
145         name += nonstat
146     else:
147         name += stat
148     name += "_"

```

```

149     return name + params
150
151 print "loading embeddings"
152 # pretrained = train_embed(embedding_dim, vocab)
153 pretrained = load_embeddings(
154     'GoogleNews-vectors-negative300.bin',
155     vocab)
156 # pretrained = None
157 # for train_indices, test_indices in ten_fold.split(X, y):
158 tf.reset_default_graph()
159 network = input_data(
160     shape=[None, max_document_length],
161     name='input')
162 emb = tflearn.embedding(network,
163     input_dim=vocab_size,
164     output_dim=embedding_dim)
165 branches = []
166 for filter_size in filter_sizes:
167     conv = conv_1d(emb, num_filters,
168         filter_size,
169         padding='valid',
170         activation='relu',
171         regularizer="L2",
172         name="conv-{}".format(
173             filter_size))
174     branches.append(conv)
175 branches = map(
176     lambda b: tf.expand_dims(b, 2),
177     branches)
178 network = merge(branches, axis=1,
179     mode='concat')
180 # network = conv_2d(network, 32, [1, 2], padding='valid', activation='relu', name='2conv2x2')
181 network = tf.reduce_max(network, axis=1)
182 network = flatten(network)
183 network = dropout(network,
184     dropout_keep_prob)
185 network = fully_connected(network,
186     n_classes,
187     activation='softmax')
188 network = regression(network,
189     optimizer='adam',
190     learning_rate=0.0001,
191     loss='categorical_crossentropy',
192     name='target')
193 run_name = 'mr_std_conf'
194 checkpoint_path = 'results/' + run_name

```

```

195     os.makedirs(checkpoint_path)
196     model = tflearn.DNN(network,
197         tensorboard_verbose=3,
198         checkpoint_path=checkpoint_path,
199         tensorboard_dir='logs')
200     if nonstatic or multichannel:
201         model.set_weights(emb.W, pretrained)
202
203     model.fit(X_train, y_train,
204         n_epoch=20,
205         validation_set=(X_test, y_test),
206         show_metric=True,
207         batch_size=batch_size,
208         snapshot_step=100,
209         snapshot_epoch=False,
210         validation_batch_size=512,
211         run_id=run_name)

```

A.2. Векторное представление слов

```

1  from __future__ import division
2
3  import datetime
4  import pickle
5  import sys
6
7  import numpy as np
8  import tensorflow as tf
9  from nltk import ngrams
10
11 import config
12 import data_helper
13
14
15 class GloveModule:
16     def __init__(self, coocur_file, ngram_table_file):
17         print("Load vocab and cooc matrix...")
18         self.vocab, cooccurrence_matrix, self.ngram_vocab, self.max_ngram_cnt = \
19             data_helper.get_wiki_corpus_and_dump(coocur_file, ngram_table_file)
20
21         self.cooccurrence_matrix = list(cooccurrence_matrix.items())
22         vocab_size = len(self.vocab)
23         ngram_table_size = len(self.ngram_vocab)
24         print("Done loading vocab and cooc.")
25         print("Creating tensorflow graph...")
26         self.graph = tf.Graph()

```

27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72

```
with self.graph.as_default():
    with self.graph.device("/cpu:0"):
        count_max = tf.constant([config.COUNT_MAX], dtype=tf.float32) # tf -
        alpha = tf.constant([config.SCALING_FACTOR], dtype=tf.float32) # tf -
        # word_ids
        self.focal_input_w_id = tf.placeholder(tf.int32, shape=[config.BATCH_SIZE])
        self.context_input_w_id = tf.placeholder(tf.int32, shape=[config.BATCH_SIZE])
        # ngram_ids
        self.focal_input = tf.placeholder(tf.int32, shape=[config.BATCH_SIZE, self.max_ngram_cnt])
        self.context_input = tf.placeholder(tf.int32, shape=[config.BATCH_SIZE, self.max_ngram_cnt])

        self.cooccurrence_count = tf.placeholder(tf.float32, shape=[config.BATCH_SIZE])
        self.ngram_embeddings = tf.Variable(
            tf.random_uniform([ngram_table_size, config.EMBEDDING_SIZE], 1.0, -1.0)
        )

        focal_word_biases = tf.Variable(
            tf.random_uniform([vocab_size], 1.0, -1.0)
        )

        context_word_biases = tf.Variable(
            tf.random_uniform([vocab_size], 1.0, -1.0)
        )

        focal_embedding = tf.nn.embedding_lookup([self.ngram_embeddings], self.focal_input_w_id)
        context_embedding = tf.nn.embedding_lookup([self.ngram_embeddings], self.context_input_w_id)
        focal_bias = tf.nn.embedding_lookup([focal_word_biases], self.focal_input)
        context_bias = tf.nn.embedding_lookup([context_word_biases], self.context_input)

        weighting_factor = tf.minimum(
            1.0,
            tf.pow(
                tf.div(self.cooccurrence_count, count_max),
                alpha
            )
        )

        #  $s(w_i, w_j) = \sum_{z \in G_j} \{w_i^T z\}$ 
        context_word_embedding = tf.reduce_sum(context_embedding, axis=0)
        embedding_product = tf.reduce_sum(tf.multiply(focal_embedding, context_word_embedding), axis=0)

        log_cooccurrences = tf.log(tf.to_float(self.cooccurrence_count))

        distance_expr = tf.square(tf.add_n([
            embedding_product,
            focal_bias,
```

```

73         context_bias,
74         tf.negative(log_cooccurrences)
75     ]))
76
77     single_losses = tf.multiply(weighting_factor, distance_expr)
78     self.total_loss = tf.reduce_sum(single_losses)
79     # with self.graph.device("/cpu:0"):
80     tf.summary.scalar('loss', self.total_loss)
81     # self.optimizer = tf.train.AdamOptimizer(config.LEARNING_RATE).minimize(self.total_loss)
82     self.optimizer = tf.train.AdagradOptimizer(config.LEARNING_RATE).minimize(self.total_loss)
83     print("Done creating")
84
85     def get_ngram_ids(self, word_id):
86         result = np.zeros(self.max_ngram_cnt)
87         for i, subword in enumerate(ngrams('<' + self.vocab[word_id] + '>', config.NGRAM_SIZE)):
88             result[i] = self.ngram_vocab[subword]
89         return np.array(result)
90
91     def get_batch(self, idx):
92         ww, counts = zip(*self.cooccurrence_matrix[idx:idx + config.BATCH_SIZE])
93         i_s, j_s = zip(*ww)
94         nis = []
95         for w_id in i_s:
96             nis.append(self.get_ngram_ids(w_id))
97         njs = []
98         for w_id in j_s:
99             njs.append(self.get_ngram_ids(w_id))
100        return i_s, j_s, np.array(nis), np.array(njs), counts
101
102    def begin(self):
103        print("Prepare cooccurrence matrix...")
104
105        # i_indices, j_indices, counts = zip(*self.cooccurrence_matrix)
106        print("Done preparing")
107        print("Get batches...")
108        # self.batches = list(self.batchify(i_indices, j_indices, cooc))
109        print("Done get batches.")
110        print("Begin training: {}".format(datetime.datetime.now().time()))
111        print("=====")
112        sys.stdout.flush()
113        with tf.Session(graph=self.graph) as session:
114            merged = tf.summary.merge_all()
115            batch_writer = tf.summary.FileWriter('./logs/batch')
116            tf.initialize_all_variables().run()
117            for epoch in range(config.NUM_EPOCHS):
118                # shuffle(self.batches)

```



```

119     print("Batches shuffled")
120     print("-----")
121     sys.stdout.flush()
122     accumulated_loss = 0
123     total = len(self.cooccurrence_matrix)
124     num_batches = total / config.BATCH_SIZE
125     batch_index = 0
126     for i in range(0, total, config.BATCH_SIZE):
127         batch_index += 1
128         i_s, j_s, nis, njs, counts = self.get_batch(i)
129
130         if len(counts) != config.BATCH_SIZE:
131             continue
132         feed_dict = {
133             self.focal_input_w_id: i_s,
134             self.context_input_w_id: j_s,
135             self.focal_input: nis,
136             self.context_input: njs,
137             self.cooccurrence_count: counts
138         }
139         summaries, __, total_loss_, = session.run(
140             [merged, self.optimizer, self.total_loss], feed_dict=feed_dict)
141         accumulated_loss += total_loss_
142         if (batch_index + 1) % config.REPORT_BATCH_SIZE == 0:
143             print("Epoch: {0}/{1}".format(epoch + 1, config.NUM_EPOCHS))
144             print("Batch: {0}/{1}".format(batch_index + 1, num_batches))
145             print("Average loss: {}".format(accumulated_loss / config.REPORT_BATCH_SIZE))
146             print("-----")
147             batch_writer.add_summary(summaries, epoch * num_batches + batch_index)
148             sys.stdout.flush()
149             accumulated_loss = 0
150         print("Epoch finished: {}".format(datetime.datetime.now().time()))
151         print("=====")
152
153     sys.stdout.flush()
154     batch_writer.close()
155     final_embeddings = self.ngram_embeddings.eval()
156     print("End: {}".format(datetime.datetime.now().time()))
157     final_dict = {}
158     for word, idx in self.vocab.items():
159         w_emb = np.zeros((1, config.EMBEDDING_SIZE), dtype=np.float32)
160         for ngram_id in self.get_ngram_ids(idx):
161             w_emb += final_embeddings[ngram_id]
162         final_dict[word] = w_emb
163     final_table = {}
164     for ngram_id in range(final_embeddings.shape[0]):

```

```

165         final_table[ngram_table[ngram_id]] = final_embeddings[ngram_id]
166     return final_dict, final_table

1  import os
2  import pickle
3  import glove
4  from nltk import ngrams
5  import argparse
6  import config
7  from GloVe import GloveModule
8
9  def get_wiki_corpus_and_dump(coocur_file, ngram_table_file):
10     corpus = glove.Corpus.load(coocur_file)
11     vocab = corpus.dictionary
12     coocur = corpus.matrix
13     max_ngram_cnt = -1
14     if os.path.exists(ngram_table_file):
15         ngram_table, max_ngram_cnt = pickle.load(open(ngram_table_file))
16     else:
17         unique_ngrams = set()
18         for _, word in vocab.items():
19             cur_ngrams = ngrams('<' + word + '>', config.NGRAM_SIZE)
20             cur_ngram_cnt = len(cur_ngrams)
21             if cur_ngram_cnt > max_ngram_cnt:
22                 max_ngram_cnt = cur_ngram_cnt
23             unique_ngrams += set(cur_ngrams)
24         ngram_table = {ngram: n_id for n_id, ngram in enumerate(unique_ngrams)}
25         pickle.dump((ngram_table, max_ngram_cnt), open(ngram_table_file, 'wb+'))
26
27     return coocur, vocab, ngram_table, max_ngram_cnt
28
29
30
31
32
33
34 if __name__ == '__main__':
35
36
37     parser = argparse.ArgumentParser()
38     parser.add_argument('--coocur-file')
39     parser.add_argument('--nragm-table-file')
40     args = parser.parse_args()
41     module = GloveModule(args.coocur_file, args.ngram_table_file)
42     embeddings, final_table = module.begin()
43     with open('data/emb/glove{}.pkl'.format(config.EMBEDDING_SIZE), 'wb+') as f:
44         pickle.dump(embeddings, f, protocol=4)

```

```

45     with open('data/emb/glove_ngram{}.pkl', 'wb+') as f:
46         pickle.dump(final_table, f, protocol=4)

1  from __future__ import print_function
2  import argparse
3  import pprint
4  import gensim
5  import subprocess
6
7  from tqdm import tqdm
8
9  from glove import Glove
10 from glove import Corpus
11
12
13 def get_line_no(fp):
14     p = subprocess.Popen(['wc', '-l', fp], stdout=subprocess.PIPE, stderr=subprocess.PIPE)
15     result, err = p.communicate()
16     if p.returncode != 0:
17         raise IOError(err)
18     return int(result.strip().split()[0])
19
20
21 def read_corpus(filename):
22     delchars = [chr(c) for c in range(256)]
23     delchars = [x for x in delchars if not x.isalnum()]
24     delchars.remove(' ')
25     delchars = ''.join(delchars)
26
27     table = str.maketrans(dict.fromkeys(delchars))
28
29     with open(filename, 'r') as datafile:
30         for line in tqdm(datafile, total=get_line_no(filename)):
31             yield line.lower().translate(table).split(' ')
32
33
34 def read_wikipedia_corpus(filename):
35     # We don't want to do a dictionary construction pass.
36     corpus = gensim.corpora.WikiCorpus(filename, dictionary={})
37
38     for text in corpus.get_texts():
39         yield text
40
41
42 if __name__ == '__main__':
43
44     # Set up command line parameters.

```

```

45 parser = argparse.ArgumentParser(description='Fit a GloVe model.')
46
47 parser.add_argument('--create', '-c', action='store',
48                     default=None,
49                     help=('The filename of the corpus to pre-process. '
50                           'The pre-processed corpus will be saved '
51                           'and will be ready for training.))
52 parser.add_argument('-wiki', '-w', action='store_true',
53                     default=False,
54                     help=('Assume the corpus input file is in the '
55                           'Wikipedia dump format'))
56 parser.add_argument('--train', '-t', action='store',
57                     default=0,
58                     help=('Train the GloVe model with this number of epochs.'
59                           'If not supplied, '
60                           'We\'ll attempt to load a trained model'))
61 parser.add_argument('--parallelism', '-p', action='store',
62                     default=1,
63                     help=('Number of parallel threads to use for training'))
64 parser.add_argument('--query', '-q', action='store',
65                     default='',
66                     help='Get closes words to this word.')
67 args = parser.parse_args()
68
69 if args.create:
70     # Build the corpus dictionary and the cooccurrence matrix.
71     print('Pre-processing corpus')
72
73     if args.wiki:
74         print('Using wikipedia corpus')
75         get_data = read_wikipedia_corpus
76     else:
77         get_data = read_corpus
78
79     corpus_model = Corpus()
80     corpus_model.fit(get_data(args.create), window=10)
81     corpus_model.save('corpus.model')
82
83     print('Dict size: %s' % len(corpus_model.dictionary))
84     print('Collocations: %s' % corpus_model.matrix.nnz)
85
86 if args.train:
87     # Train the GloVe model and save it to disk.
88
89     if not args.create:
90         # Try to load a corpus from disk.

```

```

91     print('Reading corpus statistics')
92     corpus_model = Corpus.load('corpus.model')
93
94     print('Dict size: %s' % len(corpus_model.dictionary))
95     print('Collocations: %s' % corpus_model.matrix.nnz)
96
97     print('Training the GloVe model')
98
99     glove = Glove(no_components=100, learning_rate=0.05)
100    glove.fit(corpus_model.matrix, epochs=int(args.train),
101              no_threads=args.parallelism, verbose=True)
102    glove.add_dictionary(corpus_model.dictionary)
103
104    glove.save('glove.model')
105
106    if args.query:
107        # Finally, query the model for most similar words.
108        if not args.train:
109            print('Loading pre-trained GloVe model')
110            glove = Glove.load('glove.model')
111
112        print('Querying for %s' % args.query)
113        pprint.pprint(glove.most_similar(args.query, number=10))

```

```

1  import subprocess
2  from collections import Counter
3  from time import clock
4  from typing import List, Dict
5
6  import numpy as np
7  import pandas as pd
8  from tqdm import tqdm
9
10
11  def get_line_no(fp):
12      p = subprocess.Popen(['wc', '-l', fp], stdout=subprocess.PIPE, stderr=subprocess.PIPE)
13      result, err = p.communicate()
14      if p.returncode != 0:
15          raise IOError(err)
16      return int(result.strip().split()[0])
17
18
19  def load_data(path: str):
20      with open(path, 'r') as f:
21          for line in tqdm(f, total=get_line_no(path)):
22              yield tokenize(line)
23

```

```

24
25 def tokenize(sentence_str: str):
26     # TODO: find better solution for this
27     return sentence_str.split(" ")
28
29
30 def generate_contexts(sentence: List[str], context_size):
31     last_index = len(sentence)
32
33     def window(start, end):
34         return sentence[max(start, 0): min(end, last_index)]
35
36     for i, word in enumerate(sentence):
37         yield window(i - context_size, i), word, window(i + 1, i + context_size)
38
39
40 def charseq(word, c2i):
41     chars = []
42     for c in word:
43         if c not in c2i:
44             c2i[c] = len(c2i)
45             chars.append(c2i[c])
46     return chars
47
48
49 def prepare_corpus(
50     path: str,
51     context_size: int = 5,
52     min_count: int = 3,
53     w2v_file: str = './sa_cnn/GoogleNews-vectors-negative300.bin'
54 ) -> (Dict[str, int], np.ndarray):
55     """
56     :param w2v_file:
57     :param min_count:
58     :param context_size:
59     :param path:
60     """
61     print('prepare data, loading vectors')
62     from gensim.models import KeyedVectors
63     w2v = KeyedVectors.load_word2vec_format(w2v_file, binary=True)
64     # filter rare words and index every word so we get id -> word
65     st = clock()
66     print('constructing vocab')
67     word_counts = Counter()
68     for sentence in load_data(path):
69         word_counts.update(sentence)

```

```

70 df = pd.DataFrame(list(word_counts.items()), columns=['word', 'counts'])
71 df = df[df['counts'] > min_count]
72 df = df.reset_index(drop=True)
73 vocab = {idx: w for idx, w in df['word'].to_dict().items() if w in w2v}
74 word_to_id = {v: k for k, v in vocab.items()}
75 vocab_size = len(vocab)
76 print("vocab get time: ", clock() - st)
77 print("vocab size", vocab_size)
78
79 all_pairs = []
80 char_vocab = {}
81 for sentence in load_data(path):
82     for l_ctx, word, r_ctx in generate_contexts(sentence, context_size):
83         if word in word_to_id:
84             charseq(word, char_vocab)
85             contexts = l_ctx + r_ctx
86             for ctx in contexts:
87                 if ctx in word_to_id:
88                     all_pairs.append((word_to_id[word], word_to_id[ctx]))
89 return all_pairs, vocab, char_vocab, w2v
90
91
92 def load_dataset():
93     return prepare_corpus(
94         'wiki.txt', context_size=5, min_count=1
95     )

1 import tensorflow as tf
2 import tflearn
3 import numpy as np
4 from tflearn import input_data, embedding, bidirectional_rnn, BasicLSTMCell, regression
5 from tflearn.data_utils import pad_sequences
6 from tflearn.models.dnn import DNN
7
8 from data_helper import load_dataset
9
10 max_chars = 128
11 emb_size = 300
12 char_emb_size = 20
13 num_lstm_units = 50
14
15
16 # network graph
17 def get_network(char_vocab_size):
18     context_word_vec = input_data([None, 300], name="context_w", dtype=tf.float32)
19     char_input = input_data([None, max_chars], name="word_chars", dtype=tf.float32)
20     char_embeddings = embedding(

```

```

21     char_input,
22     input_dim=char_vocab_size,
23     output_dim=char_emb_size,
24     trainable=True,
25     name='char_embeddings'
26 )
27
28 net = bidirectional_rnn(char_embeddings,
29                         BasicLSTMCell(num_lstm_units),
30                         BasicLSTMCell(num_lstm_units))
31 bilstm_out_with_context = tflearn.merge([net, context_word_vec], mode='concat')
32 net = tflearn.fully_connected(bilstm_out_with_context, 300)
33 return regression(net, optimizer='adam', loss='mean_square', metric='R2')
34
35
36 def main():
37     # preparing data to fit
38     dataset, word_vocab, char_vocab, w2v = load_dataset()
39     context_vectors = []
40     char_data = []
41     targets = []
42     print(len(dataset))
43     print(char_vocab)
44     for (w_id, c_id) in dataset:
45         chars = []
46         for ch in word_vocab[w_id]:
47             chars.append(char_vocab[ch] if ch in char_vocab else None)
48         char_data.append(chars)
49         context_vectors.append(w2v[word_vocab[c_id]])
50         targets.append(w2v[word_vocab[w_id]])
51     context_vectors = np.array(context_vectors)
52     char_data = pad_sequences(char_data, max_chars, value=0.)
53     print(char_data.shape)
54     print(context_vectors.shape)
55     # network session
56     net = get_network(len(char_vocab))
57     model = DNN(net, tensorboard_dir='./logs/')
58
59     # fitting the model
60     model.fit({
61         "context_w": context_vectors,
62         "word_chars": char_data
63     }, np.array(targets),
64         n_epoch=5,
65         validation_set=0.1,
66         show_metric=True,

```



```

67     batch_size=64,
68     shuffle=True
69 )
70
71 model.save('model.bin')
72
73
74 if __name__ == '__main__':
75     main()

1  #!/usr/bin/env bash
2  #
3  # Copyright (c) 2016-present, Facebook, Inc.
4  # All rights reserved.
5  #
6  # This source code is licensed under the BSD-style license found in the
7  # LICENSE file in the root directory of this source tree. An additional grant
8  # of patent rights can be found in the PATENTS file in the same directory.
9  #
10
11 set -e
12
13 normalize_text() {
14     sed -e "s/'/'/g" -e "s/ /'/g" -e "s"/'/g" -e "s'/ /'/g" -e "s/"/'/g" -e "s/'/'/g" \
15         -e 's"/' /' /g' -e 's/\./ \. /g' -e 's/<br \/>/ /g' -e 's/, / , /g' -e 's/( ( /g' -e 's/)/ ) /g' -e 's/!/ \! /g' \
16         -e 's/\?/ \? /g' -e 's/\;/ /g' -e 's/\:/ /g' -e 's/-/ - /g' -e 's/=/ = /g' -e 's/=/ = /g' -e 's/*/ /g' -e 's/|/ /g' \
17         -e 's/«/ /g' | tr 0-9 " "
18 }
19
20 export LANGUAGE=en_US.UTF-8
21 export LC_ALL=en_US.UTF-8
22 export LANG=en_US.UTF-8
23
24 NOW=$(date +"%Y%m%d")
25
26 ROOT="data/wikimedia/${NOW}"
27 mkdir -p "${ROOT}"
28 echo "Saving data in ""${ROOT}"
29 read -r -p "Choose a language (e.g. en, bh, fr, etc.): " choice
30 LANG="$choice"
31 echo "Chosen language: ""${LANG}"
32 read -r -p "Continue to download (WARNING: This might be big and can take a long time!)(y/n)? " choice
33 case "$choice" in
34     y|Y ) echo "Starting download...";;
35     n|N ) echo "Exiting";exit 1;;
36     * ) echo "Invalid answer";exit 1;;
37 esac

```

```

38 wget -c "https://dumps.wikimedia.org/"$LANG"wiki/latest/"${LANG}"wiki-latest-pages-articles.xml.bz2"\
39 -P "${ROOT}"
40 echo "Processing "$ROOT"/"$LANG"wiki-latest-pages-articles.xml.bz2"
41 bzip2 -c -d "$ROOT"/"$LANG"wiki-latest-pages-articles.xml.bz2" | awk '{print tolower($0);}' | \
42 perl wikifil.pl | normalize_text | awk '{if (NF>1) print;}' | tr -s " " | gshuf > "${ROOT}"/wiki."${LANG} ".txt

```

```

1 $/=">"; # input record separator
2 while (<>) {
3   if (<text /) {$text=1;} # remove all but between <text> ... </text>
4   if (/#redirect/i) {$text=0;} # remove #REDIRECT
5   if ($text) {
6     # Remove any text not normally visible
7     if (<\text>/) {$text=0;}
8     s/<.*>/;/ # remove xml tags
9     s/&amp;/&/g; # decode URL encoded chars
10    s/&lt;/</g;
11    s/&gt;/>/g;
12    s/<ref[^<]*<\ref>/;/g; # remove references <ref...> ... </ref>
13    s/<[^>]*>/;/g; # remove xhtml tags
14    s/\[http:[^ ]*\]/;/g; # remove normal url, preserve visible text
15    s/\|thumb/;/g; # remove images links, preserve caption
16    s/\|left/;/g;
17    s/\|right/;/g;
18    s/\|d+px/;/g;
19    s/\|[image:[^\]]*\]/;/g;
20    s/\|[category:([^\]]*)[^ ]*\]/[1]/g; # show categories without markup
21    s/\|[a-z-]*:[^\]]*\]/;/g; # remove links to other languages
22    s/\|[^\]]*\]/;/g; # remove wiki url, preserve visible text
23    s/{\{^}\}/;/g;
24    s/{^}\}/;/g;
25    s[/]/;/g;
26    s\[]/;/g;
27    s/&[^;]*;/ /g; # remove URL encoded chars
28    $_=" $_ ";
29    chop;
30    print $_;
31  }
32 }

```