

УДК 004.89

**КЛАССИФИКАЦИЯ ТЕКСТОВ С
ПОМОЩЬЮ СВЕРТОЧНЫХ
НЕЙРОННЫХ СЕТЕЙ***Н. В. Воробьев, Е. В. Пучков*

Донской государственной технической
университет, Ростов-на-Дону, Российская
Федерация

niknegeralmail@gmail.compuchkoff@i-intellect.ru

Исследованы методы применения сверточных нейронных сетей для классификации коротких текстов — подзадачи разработки чат-бота. Проведено сравнение качества работы построенного классификатора с результатами, полученными с помощью метода опорных векторов и рекуррентных нейронных сетей. Показано, что реализованный классификатор демонстрирует результаты лучше, чем другие методы, а также может быть рекомендован для применения в чат-боте.

Ключевые слова: классификация, обработка естественного языка, сверточные нейронные сети, глубокое обучение.

Введение. Сверточные нейронные сети (*CNN*) успешно зарекомендовали себя в задачах распознавания и классификации изображений. Это привело к множественным попыткам использования этой модели в других областях с целью улучшить результаты в сравнении с традиционными методами. На сегодняшний день классификация текста, т. е. определение принадлежности текста к какой-либо категории (классу) в условиях постоянно возрастающего объема информации является актуальной задачей. Классификация текста применяется в решении многих практических задач, таких как фильтрация документов, распознавание спама, классификация новостей и т. д. Очень часто задача классификации текста является отдельным элементом сложной проблемы, например, создание интеллектуального онлайн-консультанта или чат-бота. В этом случае имеется множество сообщений, разделённых некоторым образом на классы.

Цель данной работы — построить модель (классификатор) с помощью сверточных нейронных сетей для определения типа новых сообщений, которые пользователь посылает чат-боту. Проведение исследований при анализе данных и классификации текста осуществлялось в соответствии с методологией *CRISP-DM* [1]. В рамках поставленной задачи можно выделить следующие этапы:

1. Сбор и подготовка данных для обучения;
2. Построение классификатора;
3. Реализация классификатора в виде программного кода;

UDC 004.89

**CLASSIFICATION OF TEXTS USING
CONVOLUTIONAL NEURAL NETWORKS***N. V. Vorobev, E. V. Puchkov*

Don State Technical University, Rostov-on-Don,
Russian Federation

niknegeralmail@gmail.compuchkoff@i-intellect.ru

The article investigates the convolutional neural networks for the classification of short texts as subtasks for the chat-bot development. The constructed classifier performance has been compared with the support vector machine and recurrent neural networks. The implemented classifier has showed better results than other methods and can be recommended for using in chat-bot.

Keywords: classification, natural language processing, convolution neural networks, deep learning.

4. Обучение и тестирование созданной модели;
5. Верификация модели.

Сбор и подготовка данных. Данные представляют собой короткие сообщения — диалоги, собранные из разных источников: социальные сети; субтитры к фильмам; сообщения, написанные самостоятельно. Количество сообщений приблизительно равно 1700. Было определено 4 класса вышеупомянутых сообщений:

1. беседа;
2. запрос пользователя;
3. контактная информация;
4. информация по доставке.

Сообщения помещены в базу данных. С помощью веб-интерфейса корпус был размечен, т. е. для каждого сообщения был определен класс. Для удобства работы весь корпус переведен в *json*-формат. Проведена предподготовка входных данных. Использовано два подхода — токенизация и технология *Word2Vec*.

Построение классификатора. Сверточные нейронные сети [2], первоначально не предназначались для работы с текстом, они использовались в «компьютерном зрении» и распознавании образов. Сверточная нейронная сеть — это особый вид нейронных сетей прямого распространения (рис. 1). Под прямым распространением понимается то, что распространение сигналов по нейронам идет по порядку, от первого слоя до последнего. Скрытых слоев в сети может быть достаточно много, всё зависит от количества данных и сложности задачи.



Рис. 1. Модель сверточной нейронной сети

Основной особенностью таких сетей является наличие чередующихся слоев типа «свертка — субдискретизация», которых может быть множество. Операция свертки (рис. 2) подразумевает, что каждый фрагмент входа поэлементно умножается на небольшую матрицу весов (ядро), а результат суммируется. Эта сумма является элементом выхода, который называется картой признаков. Взвешенная сумма входов пропускается через функцию активации.

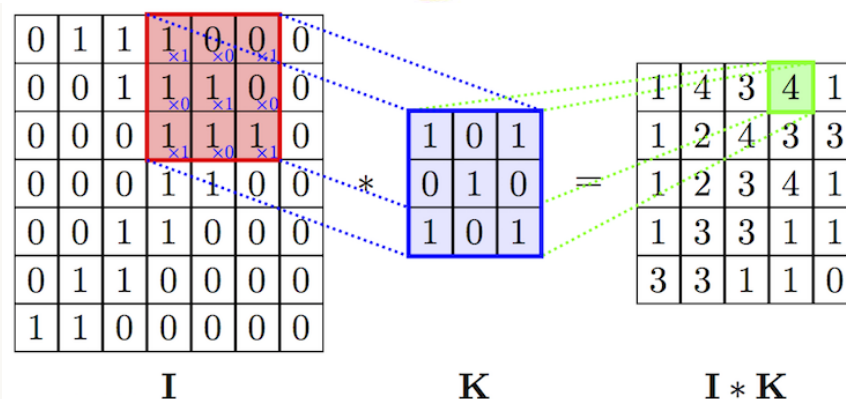


Рис. 2. Слой свертки

Слой пулинга (подвыборки, субдискретизации) представляет собой нелинейное уплотнение карты признаков, проходя нелинейное преобразование. Пулинг интерпретируется как разбиение карты признаков на более мелкие матрицы, нахождение их максимальных элементов, т. е. происходит увеличение «глубины» значений.

Для решения задачи была выбрана следующая архитектура сети:

1. входной слой;
2. слой свертки (*convolution layer*);
3. слой субдискретизации (*max pooling*);
4. полносвязный скрытый слой;
5. выходной слой.

Между основными слоями будут добавлены дополнительные слои для защиты от переобучения сети (*dropout*). Принцип их работы заключается в том, что сеть «забывает» некий процент весов (передается параметром в *dropout*). В качестве активации в слое свертки и скрытом слое будет использована выпрямленная линейная функция активации «*Relu*». Для выходного слоя выбрана функция активации «*Softmax*» — один из частных случаев функции сигмоиды, который используется для многоклассовой классификации.

Программная реализация классификатора. Для реализации данной задачи было выбрано следующее программное обеспечение:

1. язык программирования *Python*, на котором создано множество библиотек для работы с данными и нейронными сетями;
2. библиотека «*Keras*» [3], которая позволяет имплементировать *CNN* с помощью «*TensorFlow*» [4] на более высоком программном уровне;
3. среда разработки *Spyder IDE*.

Отрывок кода реализации сверточной нейронной сети представлен на рис. 3.

```
119 model = Sequential()
120
121 #input layer
122 model.add(Embedding(max_features,
123                     embedding_dims,
124                     input_length=maxlen))
125 model.add(Dropout(0.2))
126
127 #convlution layer 1D
128 model.add(Conv1D(filters,
129                 kernel_size,
130                 padding='valid',
131                 activation='relu',
132                 strides=1))
133 # we use max pooling:
134 model.add(GlobalMaxPooling1D())
135
136 # hidden layer:
137 model.add(Dense(hidden_dims))
138 model.add(Dropout(0.2))
139 model.add(Activation('relu'))
140
141 #output layer:
142 model.add(Dense(3))
143 model.add(Activation('softmax'))
144
```

Рис. 3. Программная реализация модели в *Keras*

Обучение и тестирование модели. Для обучения модели были выбраны следующие параметры:

1. Метрика для ошибки: «*categorical_crossentropy*». Используется, когда желательна вероятностная интерпретация оценок. Данная метрика измеряет сходство между истинным классом и прогнозируемым.

2. Оптимизатор: «*adam*». Алгоритм сочетает в себе идею накопления движения градиента и идею более слабого обновления весов для типичных признаков.

3. Метрика для оценки точности классификатора: «*F-мера*». Для многоклассовой классификации выделяют два подхода расчета *F-меры*: *micro*-усреднение и *macro*-усреднение. При *micro*-усреднении значения матрицы несоответствий усредняются по всем классам, а затем вычисляется *F-мера*. При *macro*-усреднении сначала вычисляется *F-мера* для каждого класса, а затем результаты усредняются по всем классам.

Все исходные данные были разделены на обучающую (80 %) и тестовую выборку (20 %). Обучение прекращалось, когда ошибка на тестовой выборке начинала расти, а на обучающей продолжала падать (метод ранней остановки). Обучение модели проходило на ноутбуке *Asus K551L* с видеокартой *NVIDIA 840M*. Графический процессор этой видеокарты поддерживает параллельные вычисления (*CUDA*), что сокращает время обучения в несколько раз по сравнению с вычислениями на *CPU* [5]. Одна эпоха длилась 2–3 секунды (рис. 4).

```

1556/1556 [=====] - 11s - loss: 0.3557 - fmeasure: 0.8933 - acc: 0.9332
Epoch 2/10
1556/1556 [=====] - 3s - loss: 0.2158 - fmeasure: 0.9408 - acc: 0.9389
Epoch 3/10
1556/1556 [=====] - 3s - loss: 0.1227 - fmeasure: 0.9647 - acc: 0.9653
Epoch 4/10
1556/1556 [=====] - 3s - loss: 0.0702 - fmeasure: 0.9785 - acc: 0.9756
Epoch 5/10
1556/1556 [=====] - 3s - loss: 0.0425 - fmeasure: 0.9849 - acc: 0.9833
Epoch 6/10
1556/1556 [=====] - 3s - loss: 0.0272 - fmeasure: 0.9873 - acc: 0.9897
Epoch 7/10
1556/1556 [=====] - 3s - loss: 0.0160 - fmeasure: 0.9940 - acc: 0.9955
Epoch 8/10
1556/1556 [=====] - 3s - loss: 0.0102 - fmeasure: 0.9958 - acc: 0.9961
Epoch 9/10
1556/1556 [=====] - 3s - loss: 0.0056 - fmeasure: 0.9994 - acc: 0.9994
Epoch 10/10
1556/1556 [=====] - 3s - loss: 0.0029 - fmeasure: 0.9987 - acc: 0.9987
    
```

Рис. 4. Обучение модели

Для проведения сравнительно анализа были построены классификаторы с помощью метода опорных векторов [6] и рекуррентных нейронных сетей (*LSTM*) [7]. Результаты тестирования моделей представлены в таблице 1.

Таблица 1

Результаты тестирования

	<i>CNN</i>	<i>LSTM</i>	<i>SVM</i>
<i>Word2Vec</i>	<i>Macro</i> =0,54 <i>Micro</i> =0,6	<i>Macro</i> =0,48 <i>Micro</i> =0,56	<i>Macro</i> =0,43 <i>Micro</i> =0,51
<i>Tokenization</i>	<i>Macro</i> =0,6 <i>Micro</i> =0,62	<i>Macro</i> =0,5 <i>Micro</i> =0,59	<i>Macro</i> =0,47 <i>Micro</i> =0,55

На основе проведенного тестирования можно сделать следующие выводы: *SVM* и *LSTM* показали результаты хуже, чем *CNN*; подготовка данных с помощью токенизации оказалась эффективней, чем *Word2Vec* из-за небольшой по размеру обучающей выборки. Отметим, что время построения классификатора на основе *LSTM* было больше, чем с помощью *CNN*. Таким образом, для большой выборки данных использование *CNN* будет предпочтительнее.

Верификация модели. Модель *CNN* была проверена на контрольных примерах (размер выборки составил 60 сообщений), неиспользованных при обучении. Для демонстрации результатов создан консольный интерфейс (рис. 5).

```

>>>>Привет, как дела?
[ 9.99715626e-01 2.84020789e-04 7.65799939e-08 1.81221552e-07]
Беседа
>>>>как мне к вам обращаться?
[ 9.99962807e-01 3.57042081e-05 1.30191063e-06 9.45389047e-08]
Беседа
>>>>Имеется ли у вас доставка?
[ 7.46589940e-05 5.42177120e-03 2.61384831e-03 9.91889775e-01]
Информация по доставке
>>>>Как мне связаться с вами?
[ 9.80289042e-01 1.92561131e-02 5.67155112e-05 3.98124714e-04]
Беседа
>>>>Где находится ваш офис?
[ 1.15036564e-02 9.61918473e-01 3.20720195e-04 2.62571145e-02]
Контактная информация
    
```

Рис. 5. Результаты проверки модели

Заключение. В данной работе проводилось исследование сверточных нейронных сетей для классификации текстов с целью использования результатов в реализации чат-бота. Проведено

сравнение качества работы построенного классификатора с результатами, полученными с помощью *SVM* и *LSTM*. Доказано, что реализованный классификатор демонстрирует результаты лучше, чем другие методы. Так как создание чат-бота связано с обработкой большого количества текстов, применение сверточных нейронных сетей будет эффективнее с точки зрения вычислительной сложности, чем *LSTM*. Важно отметить, что в данной работе реализована простая архитектура сверточной нейронной сети. Улучшение параметров модели, посимвольная подготовка данных [8], а также увеличение размера исходного корпуса текстов позволит значительно повысить характеристики классификатора.

Исследование выполнено при финансовой поддержке РФФИ в рамках научного проекта № 17-01-00888 а.

Библиографический список.

1. IBM SPSS Modeler CRISP-DM Guide [Электронный ресурс] / IBM Support. — Режим доступа : [ftp:// public.dhe.ibm.com /software /analytics /spss /documentation/modeler/14.2/en/CRISP_DM.pdf](ftp://public.dhe.ibm.com/software/analytics/spss/documentation/modeler/14.2/en/CRISP_DM.pdf) (дата обращения : 23.04.17).
2. Yann Le Cun, Léon Bottou, Yoshua Bengio and Patrick Haffner: Gradient Based Learning Applied to Document Recognition, Proceedings of IEEE, 86(11):p.2278-2324, 1998.
3. Keras: The Python Deep Learning library [Электронный ресурс] / Официальный сайт библиотеки «Keras». — Режим доступа : <https://www.keras.io/> (дата обращения : 23.04.17).
4. About TensorFlow [Электронный ресурс] / Официальный сайт фреймворка «Tensorflow». — Режим доступа : <https://www.tensorflow.org/> (дата обращения : 23.04.17).
5. Параллельные вычисления от NVIDIA [Электронный ресурс] / NVIDIA. — Режим доступа : <https://www.nvidia.ru/object/cuda-parallel-computing-ru.html> (дата обращения : 23.04.17).
6. Corinna Cortes and Vladimir N. Vapnik. Support vector networks. Machine Learning, 20, 1995, pp. 1–25.
7. S. Hochreiter, J. Schmidhuber Long short-term memory // Neural Computation, 1997, 9 (8), pp. 1735–1780.
8. Y. LeCun, X. Zhang. Text understanding from scratch // Computer Science Department, arXiv:1509.01626, 2016.