

РАСПОЗНАВАНИЕ ИНТЕНТОВ С ПОМОЩЬЮ НЕЙРОСЕТЕЙ

Баймурзина Диляра Римовна

Аспирант

Факультет управления и прикладной математики МФТИ, Москва, Россия

E-mail: dilyara.rimovna@gmail.com

Распознавание интенгов (намерений) пользователя является задачей классификации текстов, одной из составляющих модуля понимания естественного языка в диалоговых системах. Для решения данной задачи могут быть использованы логистическая регрессия, машины опорных векторов, наивные байесовские классификаторы, случайные леса, а также нейронные сети.

В данном исследовании представлено решение на основе нейросетевых методов, показаны способы улучшения архитектуры сети, а также произведено сравнение результатов с готовыми системами распознавания интенгов. Исследование проведено с использованием набора данных SNIPS [5]. Он содержит около 2400 примеров на каждый из семи интенгов. На текстах комментариев с сайта Reddit [6] была обучена модель векторных представлений слов (word embeddings) с использованием библиотеки fastText [2].

В качестве базовой модели использовалась конфигурация shallow-and-wide свёрточной нейронной из [4]. В векторной форме данные служат входом для нейронной сети, которая представляет из себя входной слой, состоящий из трёх свёрток с разными размерами ядра, за каждой из которых следует слой глобальной субдискретизации (global max pooling). Далее полученные вектора конкатенируются и передаются на вход полносвязным слоям.

Основная свёрточная модель может быть улучшена путем добавления дополнительных признаков, например, векторных представлений запроса в целом (*sentence embeddings*), для получения которых используется предобученная модель InferSent [3]. Единственное отличие модели состоит в том, что полученное векторное представление запросов конкатенируется с выходом слоев global max pooling.

В следующей модели используется возможность взаимодействия обучения распознавателей интенгов и именованных сущностей. На вход отдельно подаются также полученные векторные представления именованных сущностей, встретившихся в тексте, полученные с помощью BiLSTM-CRF модели по распознаванию именованных сущностей [1], конкатенируемые с выходом слоев global max pooling.

Таблица 1: Результаты экспериментов

Модель	Среднее значение F-меры для каждого интента							
	AddTo Playlist	Book Restaurant	Get Weather	Play Music	Rate Book	Search Creative Work	Search Screening Event	
api.ai	0.9931	0.9949	0.9935	0.9811	0.9992	0.9659	0.9801	
ibm.watson	0.9931	0.9950	0.9950	0.9822	0.9996	0.9643	0.9750	
microsoft.luis	0.9943	0.9935	0.9925	0.9815	0.9988	0.9620	0.9749	
wit.ai	0.9877	0.9913	0.9921	0.9766	0.9977	0.9458	0.9673	
snips.ai	0.9873	0.9921	0.9939	0.9729	0.9985	0.9455	0.9613	
recast.ai	0.9894	0.9943	0.9910	0.9660	0.9981	0.9424	0.9539	
amazon.lex	0.9930	0.9862	0.9825	0.9709	0.9981	0.9427	0.9581	
CNN	0.9956	0.9973	0.9968	0.9871	0.9998	0.9752	0.9854	
CNN & NER	0.9964	0.9958	0.9920	0.9865	0.9970	0.9652	0.9768	
CNN & InferSent	0.9956	0.9971	0.9969	0.9879	0.9994	0.9753	0.9845	
BiLSTM	0.9612	0.9488	0.9514	0.9351	0.9688	0.8769	0.8979	
CNN & NER-truth	0.9996	0.9987	0.9986	0.9997	1.0000	1.0000	1.0000	

Рекуррентные нейронные сети часто превосходят свёрточные сети в задачах обработки текста, поэтому была исследована модификация модели в которой сверточные слои были заменены двунаправленной долговременной краткосрочной памяти (Bidirectional Long-Short Term Memory).

Результаты обучения моделей, представленные в Таблице 1, наглядно демонстрируют преимущества shallow-and-wide свёрточной сети над рекуррентной BiLSTM моделью для данной задачи. Базовая свёрточная модель с подобранным параметрами выдаёт более высокий результат, чем все рассмотренные в презентации Intento «NLU. Intent Detection. Benchmark» модели, включая *api.ai* и *ibm.watson*. Разумеется, модели в презентации обучались в формате *black box*, то есть подбор параметров, если и осуществлялся, то только внутри самой системы автоматически.

Исследования и разработки выполнены при поддержке Фонда поддержки проектов Национальной технологической инициативы и ПАО «Сбербанк». Идентификатор проекта 0000000007417F630002.

Литература

1. Anh L. T., Arkhipov M. Y., Burtsev M. S. Application of a Hybrid Bi-LSTM-CRF model to the task of Russian Named Entity Recognition // arXiv preprint arXiv:1709.09686. 2017.
2. Bojanowski P. et al. Enriching word vectors with subword information // arXiv preprint arXiv:1607.04606. 2016.
3. Conneau A. et al. Supervised learning of universal sentence representations from natural language inference data // arXiv preprint arXiv:1705.02364. 2017.
4. Le H. T., Cerisara C., Denis A. Do Convolutional Networks need to be Deep for Text Classification? // arXiv preprint arXiv:1707.04108. 2017.
5. Набор данных SNIPS:
<https://github.com/snipsco/nlu-benchmark/tree/master/2017-06-custom-intent-engines>
6. Набор данных с сайта Reddit «RC_2011-01»: <http://files.pushshift.io/reddit/comments/>