

Who Am I? Personality Detection based on Deep Learning for Texts

Xiangguo Sun, Bo Liu, Jiuxin Cao, Junzhou Luo

School of Computer Science and Engineering

Southeast University, Nanjing, China 211189

Email: { sunxiangguo, bliu, jx.cao, jluo }@seu.edu.cn

Xiaojun Shen

School of Computing and Engineering

University of Missouri-Kansas City, Kansas city, USA 64110

Email: shenx@umkc.edu

Abstract—Recently, personality detection based on texts from online social networks has attracted more and more attentions. However, most related models are based on letter, word or phrase, which is not sufficient to get good results. In this paper, we present our preliminary but interesting and useful research results to show that the structure of texts can be also an important feature in the study of personality detection from texts. We propose a model named 2CLSTM, which is a bidirectional LSTMs (Long Short Term Memory networks) concatenated with CNN (Convolutional Neural Network), to detect users personality using structures of texts. Besides, a concept, Latent Sentence Group (LSG), is put forward to express the abstract feature combination based on closely connected sentences and we use our model to capture it. To the best of our knowledge, most related works only conducted their experiments on one data set, which may not well explain the versatility of their models. We implement our evaluations on two different kinds of datasets, containing long texts and short texts. Evaluations on both datasets have achieved better results, which demonstrate that our model can efficiently learn valid text structure features to accomplish the task.

Index Terms—computational personality, Big Five, deep learning, online social networks, NLP.

I. INTRODUCTION

In the field of psychology, a generally accepted and mostly influential model to characterize and measure the degree of one's personality is the Big Five Model[1], which consists of five traits: Openness, Conscientiousness, Extroversion, Agreeableness and Neuroticism.

With the development of online social networks such as Twitter, Sina microblog¹ and so on, it is easy to obtain online social texts created by users. A number of articles[2], [3], [4] have suggested that there are strong connections between one's personality traits and their behavior observable from web text records. Therefore, personality detection based on texts from online social networks has a significant meaning for plenty of applications such as individual RecSys, mental diagnosis, human resources management and so on.

However, capturing truly useful features from texts that have close connection to user's personality is still a challenging task to explore. Although some researchers have made comprehensive research and found a number of useful features such as LIWC (Linguistic Inquiry and Word Count)[5], Mairesse[6], letter's level features[7], [8], responsive patterns[9] and so on,

unfortunately, most their applications mainly rely on features of only letters, words, or phrases and could only get limited results, which implies the necessity to explore more useful features.

We have noticed, some researchers [10], [11] have started to use ensemble learning to integrate multiple kinds of features instead of just relying on one feature to improve performance. This is a correct direction we will follow in the future. Our research experience shows that it is extremely essential to explore more relevant features which collectively would more completely and accurately reflect the relations between one's personality traits and his original text data. As a matter of fact the structure of online texts created by users, has not been studied as thoroughly as other features obtained by traditional ways. Obviously, in order to develop an effective model for multiple features, we need first to learn, to analyze, and to model each single feature well.

In this paper, we attempt to push forward the study on structure of online texts and its relation to personality. Specifically, we present our research on three questions: (1) Whether a discovered structure from the texts of online social networks can effectively reflect one's personality? (2) How to make an appropriate definition for the structure of texts? (3) How can we capture this structural feature from online texts?

Recently deep learning methods have been introduced into personality prediction and have obtained good performance. Intuitively, the convolutional neural network (CNN) attempts to reconstitute the course of writing articles, while the recurrent neural network (RNN) tries to understand texts by imitating the reading process akin to humans. By integrating these two kinds of neural networks, we propose the 2CLSIM Model, a bidirectional LSTMs concatenated with CNN, to capture the structural features from texts. In addition, a concept of latent sentence group (LSG), is introduced to describe the abstract feature combination extracted from closely related sentences. Furthermore, to the best of our knowledge, most related works only conducted their experiments on one dataset, which may not well explain the versatility of their models. We did our evaluations on two datasets that are relatively heterogeneous. In both experiments our model has shown better results, which demonstrate that our model can capture the more advanced structural feature correctly and is effective to detect users personality traits. In summary, contributions of

¹ <https://weibo.com/login.php>

this work are as follows:

- We have introduced the concept of latent sentence group (LSG) to model the structural features of texts at the level of sentences and use CNN to learn the features.
- Based on CNN and RNN, we have proposed a combined neural network model to implement our task. The results of the experiments show that our model outperforms those that only use single RNN or CNN.

The rest of this paper is organized as follows. We present related works in Section II. Then our model will be introduced in Section III. The experiments and evaluations will be presented in Section IV. After that, we come to conclusion and future work in Section V.

II. RELATED WORKS

Broadly, the data type used for this research includes texts, avatars[9], likes[12] and so on. For texts data, most researches treat the task as a peculiar kind of texts classification. Although great progress has been made in the field of text classification in recent years, personality detection from texts is still in its early stage. As previously mentioned, capturing useful features suitable for personality detection remains a challenging problem. Most previous works have mainly paid attentions to the features at character or word level. Authors of [7] used RNN for character level and word level to build hierarchical, vector representations of words and sentences for trait inferences. Qiu et al.[13] have analyzed the relationship between participants' words in Twitter and their personalities. They have found there are connections between personality characteristics and specific words used in tweets. Beyond some specific words, many researchers found that the psychological meaning of words can also reflect ones characteristics. One of the methods used to express this kind of relations is called LIWC (Linguistic Inquiry and Word Count)[14]. Xiaoqian Liu and Tingshao Zhu[11] have used LIWC to represent each tweet and extracted the main part by DFT (Discrete Fourier Transform). Then they have used stacked auto-encoders to implement unsupervised feature learning. In addition, some other researchers also found the effectiveness of using structures extracted from texts in personality detection. Honghao Wei et al.[9] have used CNN with 1,2,3-grams kernels to capture the structure. However, using fixed size for kernels would impose a limitation on long texts. Since their work is based on stacked generalization and texts data is only a part of the whole, the final output cannot demonstrate whether the simple CNN is effective. Navonil Majumder et al.[10] have used 3-dimensional convolution to learn the structure of an article. They encoded the essay from word level to sentence level. Then they tried to continue using CNN to construct an essay vector based on previous work. However, when aggregating sentence vectors into the document vector, they could not continue using convolution operation because of non-convergence. To cope with this problem, they have used the max pooling instead.

Inspired by the work of Siwei Lai et al.[15], we combine both RNN and CNN in this study to capture structural features

but with the following improvements. First, different from their work, we use LSTM instead of simple RNN so that the model can get better performance in both short and long texts. Second, we have observed that most related works only consider features on character, word or sentence. We extend the structures to the LSG, which is derived from the group of related sentences.

III. 2CLSTM MODEL

In this section, we introduce our 2CLSTM model in detail. Subsection III-A overviews the model and present the process of information flowing. Following this process, some important components of the model will be presented in the subsections from III-B to III-E.

A. Overview of the Model

The architecture of our model is depicted in Figure 1, which includes word embedding, bidirectional LSTMs, CNN layers with 1,2,3-grams kernels and classification. The whole model mainly consists of two parts: 2LSTMs and CNNLSG. Basically, the first part focuses on assembling context and extracting word semantic features, while the second part attempts to learn features from sentences structure. Five traits are trained separately. As the output from CNNLSG, document vectors will be sent into the softmax layer and get the final category.

The process of a text takes the following steps. Firstly each word is embedded into word vector, which will be explained in section B, and then the context information is encoded by the first part of our model. We name it as 2LSTMs, which will be introduced in section C. In the second part of our model (CNNLSG), we use CNN to learn the structural feature LSG produced from the first part, which will be explained in section D. The final feature vectors produced by the second part (CNNLSG) will be sent into softmax to produce the final personality traits as section E presents.

B. Word Embedding

For word embeddings, we have used pre-trained word vectors from GloVe[16], which is an unsupervised learning algorithm for obtaining word vector representations. Each word was embedded into 100 dimensions. If a word doesn't exist in the GloVe word list, a random number in $[-0.25, 0.25]$ will be assigned across the whole coordinates. However, we strongly recommend to train your own word vectors from your datasets. We didn't do this because of our limited hardware and datasets available. You can also try word2vect model instead of GloVe, there is no too much difference for our task.²

C. 2LSTMs

Original LSTM (Long Short Term Memory networks) is shown in Figure 2. Instead of a unit that only executes simple activation function, LSTM has internal self loop as well as

² Training process of word2vect with negative sampling is usually much longer than GloVe.

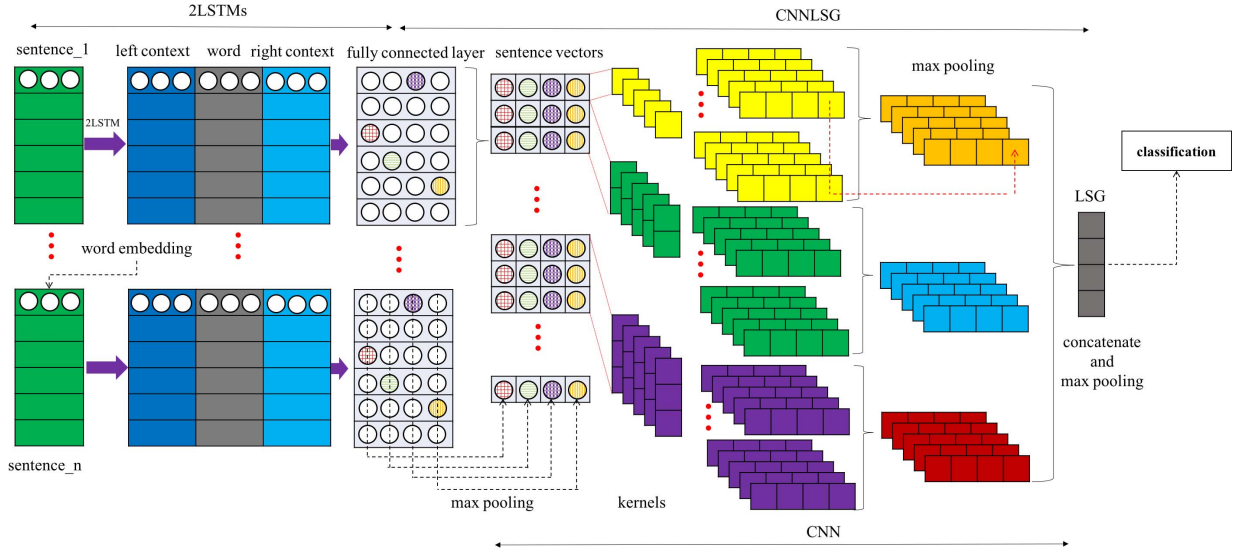


Fig. 1: The architecture of 2CLSTM

outer recurrence of the RNN. Each cell has more parameters and gates to control information flowing.

Specifically, LSTM can tolerate the unit forgetting left-context knowledge to some extent, which is realized by forget gate. The percentage of information leak is set between 0 and 1 via sigmoid function as follows:

$$f_i^{(t)} = \sigma \left(b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} \right) \quad (1)$$

where $x^{(t)}$ is the input while $h^{(t)}$ is the hidden layer vector. Meanwhile, LSTM also controls input knowledge absorption by input gate as follows:

$$g_i^{(t)} = \sigma \left(b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right) \quad (2)$$

which is similar to forget gate except using its own parameters. By filtering out the input and context information to some extent, the internal state is updated as follows:

$$s_i^{(t)} = f_i^{(t)} s_i^{(t-1)} + g_i^{(t)} \sigma \left(b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)} \right) \quad (3)$$

The output $h^{(t)}$ can also be shut off by the output gate as follows:

$$h_i^{(t)} = \tanh(s_i^{(t)}) q_i^{(t)} \quad (4)$$

where $q^{(t)}$ is calculated similarly to the forget and input gates:

$$q_i^{(t)} = \sigma \left(b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)} \right) \quad (5)$$

In order to retain more messages from context, we use bidirectional LSTMs concatenated with current word. One of the advantages is that more structural features from context (both left-context and right-context) are concentrated within

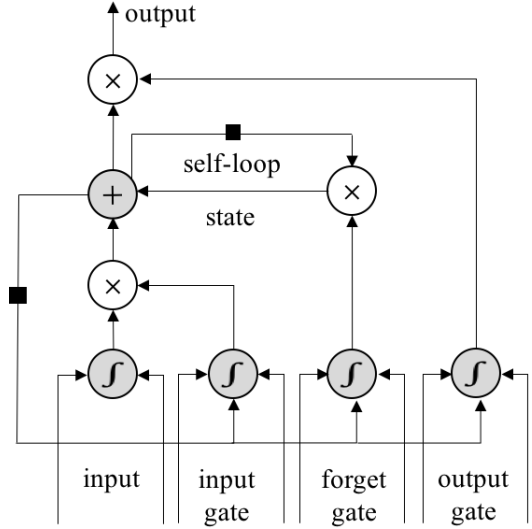


Fig. 2: LSTM cell

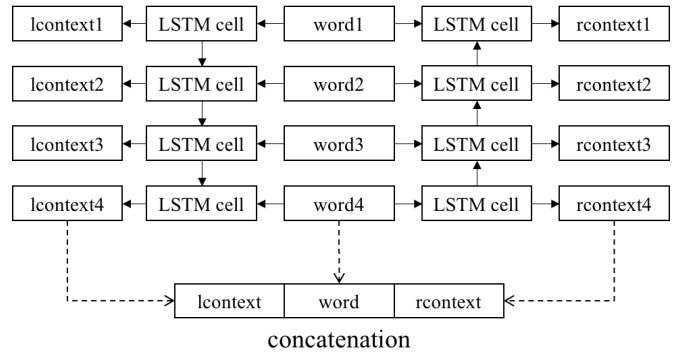


Fig. 3: The architecture of 2LSTMs

the current word, which lays the foundation of the following further feature extraction. We named this architecture as 2LSTMs, which is depicted as Figure 3.

D. Latent Sentence Group

The concept of sentence group derives from the linguistics. Basically, the sentence group means several successive sentences which are closely connected in logic and semantic structure such as coordinate relation, preference relation, causal relation and so on. However, to detect these concrete relations is impractical for most texts tasks. As a matter of fact, we often focus on the relationships between sentence vectors in some dimensions. These sentence relations are based on some specific vector dimensions, which means that is a "latent" relationship. Therefore, we use the phrase, latent sentences group (LSG), to denote the abstract combination of sentence features. The LSG means several sentences in close relationship from the view of dimension. These sentences in an LSG don't need to be strictly successive in spatial position. The definition of LSG is as follows:

Definition 1 (latent sentence group): Latent Sentence Group (LSG) is defined as a synthesis that consists of a number of sentence vectors which are closely connected in some coordinates.

We use CNN to learn LSG features. The details of this part of our model are shown in Figure 1. By means of the max pooling of each sentence, we get sentence vectors. Then in each dimension, we use 1,2,3-gram kernels to learn LSG features in each coordinate. After that, a dense layer and max pooling layer follow immediately, where the final vector will be produced and sent into the classifier.

E. Classification

Various classifiers can be used in classification part such as softmax, SVM, KNN or even Adaboost. In most deep learning models, softmax is one of the most widely used classifiers. Softmax receives feature vectors as input, and calculates the probability in each category as follows:

$$p(y = i|x; \theta) = \frac{e^{\theta_i^T x}}{\sum_j e^{\theta_j^T x}}, \quad (6)$$

where y is the classes, and x is the input vector. The class which gets higher score is the output of softmax. In our model, we use softmax as the final classification.

IV. EXPERIMENTS

In this section, we present the experiments to evaluate the effectiveness of our model for personality detection.

A. Dataset

We use the following data sets in our experiments:

1) *Stream-of-consciousness essays*: We used stream-of-consciousness essay dataset from James Pennebaker and Laura King[17]. It contains 2,467 valid anonymous essays tagged with the authors personality traits: EXT (Extroversion), NEU (Neuroticism), AGR (Agreeableness), CON (Conscientiousness), and OPN (Openness).

2) *YouTube*: The YouTube personality dataset³ comes from about 400 YouTube vloggers' webcam videos. It contains speech transcriptions, genders and behavioral features translated manually from the videos. A main difference from the first kind of dataset is that most texts are shorter in this dataset. Another difference is that the labels here (personality impressions) are collected from the annotators rating impressions by watching each vlog, not from the authors own.

Remark: We choose these two data sets because we wish to test if our model can handle different cases. Table I shows a number of statistics we summarized from these two kinds of datasets. The mainly reasons why we choose these two datasets are as follows:

- Firstly, as depicted in Figure 4, most documents from YouTube are shorter than stream-of-consciousness essays, which can verify whether the 2CLSTM model is effective both for short and long texts.
- Secondly, the labels of stream-of-consciousness essays data come from the author's own questionnaire, which can be explained as autognosis while the other dataset (YouTube dataset) can be treated as outer-perception because the personality labels of this dataset come from the volunteers watching the vlogger's videos. Therefore, these two datasets can demonstrate that our model is valid in both cases, no matter whether the labels are generated by authors or by other people.

TABLE I: Macro statistics for datasets

Property	stream-of-consciousness	youtube
average word count	648	526
max word count	2488	1972
min word count	33	36
average sentence count	46	41
max sentence count	327	147
min sentence count	1	2

B. Contrast Models

In this subsection, we choose five well-known models as our contrast models, which are listed in table III.

The first one is Bayes classification based on TF-IDF feature, which is one of the most basic and common methods for texts classification. TF-IDF is the product of term frequency (TF) and inverse document frequency (IDF). The value expresses the importance of some specific words distinguishing documents. However, in practice, words with too small (nearly zero for example) or too high TF-IDF values are not necessary to consider. In our experiment, we set the valid range of TF-IDF as [0.2, 0.5].

³ <https://www.idiap.ch/dataset/youtube-personality>

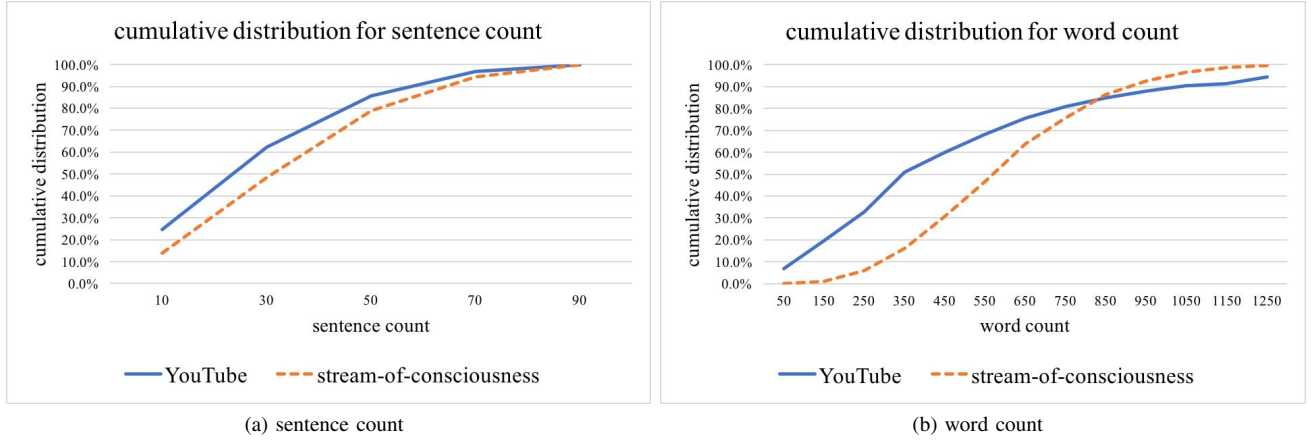


Fig. 4: Cumulative distributions of sentences and words for YouTube and stream-of-consciousness datasets

The second contrast model we use is 2CNN, 2-dimension convolution networks coming from Honghao Wei et al[9]. Since text data is only part of their whole dataset and the task is only one of the classifiers in their heterogeneous ensembles, we removed LIWC features and reserved document vectors from their CNN. We choose it as one of our contrast models because this model is the latest method which focuses on the relationship between words.

The third model is 3-dimensional CNN, which is proposed by Navonil Majumder et al[10]. They use convolution kernels in 3-dimensions (word vector, phrase, and sentence). Since the computation is huge, we joined sentences into a single one, and then use 2-dimension convolution. All maps from 1,2,3-gram kernels were padded with zeros to keep the same shape. At last, we concatenated the maps again and use max pooling layer to get the document vector. The reason why we choose it is that this model is the latest method which attempts to learn the feature based on sentences.

As for the last two contrast models, they all belong to RNN method. One is single RNN (1LSTM), and the other is bi-directional LSTMs without current word concatenated. They are all typical models in text classification.

C. Experimental Setting

Our experiments are based on ubuntu 16.04 LTS with 8GB memory and Intel Core i7-4790 CPU. The maximal word count in one document is set 800 (pad zero if less than it), and the word vector is 100 dimensions. According to these, each layer’s output shape is shown in Table II. Besides, we opened our source code for research purposes⁴.

As ordinary operations, we split our dataset as test data and train data (about 9 : 1). In each training epoch, about 20% percent of the train data is treated as validation data. The maximum epoch times set here is 100, but in the practical experiment, most tasks can early stop at 20-30 epochs. In order to overcome the overfitting, we added some dropout layers

TABLE II: Architecture shape of main layers

Layer	Output Shape
Bidirectional LSTM	(800, 600)
Concatenate with current word	(800, 700)
Conv2D with 1-gram	(25, 300, 10)
Conv2D with 2-gram	(24, 300, 10)
Conv2D with 3-gram	(23, 300, 10)
MaxPooling2D	(3, 300, 10)
MaxPooling2D	(1, 300, 1)
Softmax	(2)

* In order to overcome the overfitting, we added some dropout layers with drop rate from 0.2 to 0.3, which are not listed in the table.

with drop rate from 0.2 to 0.3, which depends on the situation.

D. Results

The macro precision of these models for both stream-of-consciousness dataset and YouTube dataset is listed in Table III. We use bold font to mark the top two best results in each category. As a whole, we can find that the 2CLSTM model performs pretty well in all categories. Although in some classes such as CON or OPN for stream-of-consciousness dataset and AGR or OPN for YouTube dataset, 2CLSTM doesn’t come to the top, it still keeps relatively good precision. In addition, for the stream-of-consciousness dataset, 2CLSTM get better in categories such as AGR and NEU, while for YouTube dataset, it gets better in EXT and NEU.

The macro precision of YouTube is generally superior to that of stream-of-consciousness, which reflects, to some extent, that models listed here (including 2CLSTM) seem to be better at detecting the personality traits in the eyes of other people than the author himself.

⁴ Source code is available for research purposes. You can download from this URL: <https://github.com/sunxiangguo/2CLSTM>

TABLE III: Macro Precision

	stream-of-consciousness					YouTube				
	AGR	EXT	CON	NEU	OPN	AGR	EXT	CON	NEU	OPN
TF-IDF+Bayes	0.3702	0.5134	0.5109	0.5363	0.3994	0.6112	0.4780	0.5670	0.5840	0.4987
2CNN	0.5029	0.5321	0.4923	0.4859	0.4773	0.5083	0.4963	0.5071	0.5796	0.4954
3CNN	0.5055	0.5020	0.4810	0.5047	0.5060	0.5046	0.5208	0.5185	0.5994	0.5133
1LSTM	0.5842	0.4972	0.5181	0.4926	0.5363	0.5538	0.6185	0.6117	0.5325	0.5791
2LSTMs	0.5743	0.5379	0.5449	0.5317	0.5654	0.5547	0.5220	0.5567	0.5613	0.5183
2CLSTM	0.5887	0.5564	0.5352	0.5677	0.5419	0.5802	0.6769	0.6117	0.6128	0.5546

* We use the bold item to represent the top two results.

V. CONCLUSION AND FUTURE WORK

In this paper, we have proposed a deep learning based model, called 2CLSTM, to detect users personality traits from online texts written by themselves. The results from experiments demonstrate that compared with the state-of-art methods in this field, our 2CLSTM performs better in detecting personality traits both from self-cognition or outer-perception point of view, no matter long texts or short texts.

Our future work can be summarized as follows:

- Besides text data, we will extend to include other data types into our work such as photos, voice, videos, relationships with friends, and responsive patterns.
- We will design an application to be performed on WeChat⁵ platform to allow users to know themselves better.

VI. ACKNOWLEDGMENT

This work is supported by National Natural Science Foundation of China under Grants, No. 61370208, No. 61772133, No. 61472081, No. 61402104, No. 61370207, No. 61300024, No. 61320106007, No. 61272531, No. 61202449, No. 61272054. Collaborative Innovation Center of Wireless Communications Technology Collaborative Innovation Center of Social Safety Science and Technology Jiangsu Provincial Key Laboratory of Network and Information Security under Grants No. BM2003201, and Key Laboratory of Computer Network and Information Integration of Ministry of Education of China under Grants No. 93K-9.

REFERENCES

- [1] Paul T Costa and Robert R McCrae. The revised neo personality inventory (neo-pi-r). *The SAGE handbook of personality theory and assessment*, 2:179–198, 2008.
- [2] H Andrew Schwartz, Johannes C Eichstaedt, Margaret L Kern, Lukasz Dziurzynski, Stephanie M Ramones, Megha Agrawal, Achal Shah, Michal Kosinski, David Stillwell, Martin EP Seligman, et al. Personality, gender, and age in the language of social media: The open-vocabulary approach. *PloS one*, 8(9):e73791, 2013.

- [3] Lin Qiu, Jiahui Lu, Jonathan Ramsay, Shanshan Yang, Weina Qu, and Tingshao Zhu. Personality expression in chinese language use. *International Journal of Psychology*, 2016.
- [4] Yilin Li, Tingshao Zhu, Ang Li, Fan Zhang, and Xinguo Xu. Web behavior and personality: A review. In *Web Society (SWS), 2011 3rd Symposium on*, pages 81–87. IEEE, 2011.
- [5] James W Pennebaker Martha E Francis and Roger J Booth. Linguistic inquiry and word count. Technical report, Technical Report, Dallas, TX: Southern Methodist University, 1993.
- [6] François Mairesse, Marilyn A Walker, Matthias R Mehl, and Roger K Moore. Using linguistic cues for the automatic recognition of personality in conversation and text. *Journal of artificial intelligence research*, 30:457–500, 2007.
- [7] Fei Liu, Julien Perez, and Scott Nowson. A language-independent and compositional model for personality trait recognition from short texts. *arXiv preprint arXiv:1610.04345*, 2016.
- [8] Cicero Nogueira Dos Santos and Maira Gatti. Deep convolutional neural networks for sentiment analysis of short texts. In *COLING*, pages 69–78, 2014.
- [9] Honghao Wei, Fuzheng Zhang, Nicholas Jing Yuan, Chuan Cao, Hao Fu, Xing Xie, Yong Rui, and Wei-Ying Ma. Beyond the words: Predicting user personality from heterogeneous information. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pages 305–314. ACM, 2017.
- [10] Navonil Majumder, Soujanya Poria, Alexander Gelbukh, and Erik Cambria. Deep learning-based document modeling for personality detection from text. *IEEE Intelligent Systems*, 32(2):74–79, 2017.
- [11] Xiaoqian Liu and Tingshao Zhu. Deep learning for constructing microblog behavior representation to identify social media users personality. *PeerJ Computer Science*, 2:e81, 2016.
- [12] Michal Kosinski, David Stillwell, and Thore Graepel. Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the National Academy of Sciences*, 110(15):5802–5805, 2013.
- [13] Lin Qiu, Han Lin, Jonathan Ramsay, and Fang Yang. You are what you tweet: Personality expression and perception on twitter. *Journal of Research in Personality*, 46(6):710–718, 2012.
- [14] Yla R Tausczik and James W Pennebaker. The psychological meaning of words: Liwc and computerized text analysis methods. *Journal of language and social psychology*, 29(1):24–54, 2010.
- [15] Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. Recurrent convolutional neural networks for text classification. In *AAAI*, volume 333, pages 2267–2273, 2015.
- [16] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, volume 14, pages 1532–1543, 2014.
- [17] James W Pennebaker and Laura A King. Linguistic styles: language use as an individual difference. *Journal of personality and social psychology*, 77(6):1296, 1999.

⁵ A very famous online social platform and communication tool in China. Website: <http://www.wechat.com/en/>