Convolutional Neural Network for Chinese Sentiment Analysis Considering Chinese Slang Lexicon and Emoticons

Da Li^{*1} Rafal Rzepka^{*1} Michal Ptaszynski^{*2} Kenji Araki^{*1}

*1 Graduate School of Information Science and Technology, Hokkaido University
*2 Department of Computer Science, Kitami Institute of Technology

Nowadays, social media have become the essential part of our lives. Internet slang is an informal language used in everyday online communication which quickly becomes adopted or discarded by new generations. Similarly, pictograms (emoticons/emojis) have been widely used in social media as a mean for graphical expression of emotions. People can convey delicate nuances through textual information when supported with emoticons, and the effectiveness of computer-mediated communication is also improved. Therefore, it is important to fully understand the influence of Internet slang and emoticons on social media. In this paper, we propose a convolutional neural network model considering Internet slang and emoticons for sentiment analysis of Weibo which is the most popular Chinese social media platform. Our experimental results show that the proposed method can significantly improve the performance for predicting sentiment polarity.

1. Introduction

Today, many people share their lives with their friends by posting status updates on Facebook, sharing their holiday photos on Instagram or tweeting their views via Twitter or Weibo - the biggest Chinese social media network that was launched in 2009. Social media data contain a vast amount of valuable sentiment information not only for the commercial use, but also for psychology, cognitive linguistics or political science [Li 18a].

Sentiment analysis of microblogs became an important area of research in the field of Natural Language Processing. Study of sentiment in microblogs in English language has undergone major developments in recent years [Peng 17]. Chinese sentiment analysis research, on the other hand, is still at early stage [Wang 13] especially in the domains of lexicons and emoticons.

Pictograms (emoticons/emojis) have been widely used in social media as a mean for graphical expression of emotions. For example, if ace with tears of joy", an emoji that means that somebody is in an extremely good mood, was regarded as the 2015 word of the year by The Oxford Dictionary [Moschini 16]. In our opinion ignoring emoticons in sentiment research is unjustifiable, because they convey a significant emotional information and play an important role in expressing emotions and opinions in social media [Novak 15, Guibon 16].

Internet slang is ubiquitous on the Internet. The emergence of new social contexts like micro-blogs, discussion groups or social networks has enabled slang and non-standard expressions to abound on the Web. Despite this, slang has been traditionally viewed as a form of non-standard language, a form of language that is not the focus of linguistic analysis and has largely been neglected [Kulkarni 17].

Furthermore, we also noticed that when people use new words and pictograms, they tend to express a kind of humorous emotion which is difficult to be easily classified as positive or negative. It seems that some emoticons are used just for fun, self-mockery or jocosity which expresses an implicit humor which might be characteristic to Chinese culture. Figure 1 shows an example of a Weibo microblog posted with emoticons and Internet slang. In the third line of the post, 柠檬人 (ning meng ren) is a new word appeared in early 2019 on Chinese social media which means "lemon

man". Accordingly, to meet this new popular phrase, was added to the pictogram repoitoare by social media companies in January 2019. This lemon with a sad face also called "lemon man" which expresses the same emotion as slang *ning meng ren* – "sour grapes" or "jealous of someone's success". This entry seems to express a humorous nuance of a pessimistic attitude. Emoticons and slang seem to play an important role in expressing this kind of emotions. There is a high possibility that this phenomenon can cause a significant difficulty in sentiment recognition task.

172	<mark>你喜欢我</mark> 2-3 23:49	什么我就 来自小米8	尤叫什么 周年旗舰手	心 机	十关注
我冷静	下来了…	明知道日	自己中不	了奖…	我就不
参与这	种热闹了	(…抽个/	小奖都抽	不中…	怎么能
指望这	种大奖呢	(2) 新的·	一年我要	到个村	宁檬人

Figure 1: Example of Weibo post with Internet slang and emoticons. The entry says "I have calmed down. Knowing that I won't win the prize, I don't participate in this excitement. I even can't win a small prize. How can I expect a big one? I want to be a lemon man in the new year".

To address this phenomenon, in this paper we focus on the Internet slang and emoticons used on Weibo in order to establish if both slang and emoticons improve sentiment analysis by recognizing humorous entries which are difficult to polarize. To perform experiments, we built a Chinese Internet slang lexicon and a Chinese emoticon lexicon. Because the emoticons probably play a more important role in expressing emotion than textual features, we also analyzed the characteristics of this particular set of emoticons, report on their evaluation while dividing them into three categories: positive, negative and humorous. We also noticed that

Emoticon	Humorous {%}	Negative {%}	Positive {%}
	41.7	25.0	33.3
	91.7	8.3	0.0
	83.3	0.0	16.7
	75.0	8.3	16.7

Table 1: Examples of Emoticons Conveying Humor Typical for Chinese Culture.

Table 2: Examples of Chinese Emoticon Lexicon.

Emoticon	Textual Feature	Emotion/Implication
63	[阴险]	"smirking"
	[挖鼻]	"nosepick"
	[污]	"filthy"
*	[舔屏]	"screen lick"

among the resources of Chinese social media sentiment analysis, the labelled Weibo data containing emoticons are extremely rare which makes considering them in machine learning approaches difficult. To resolve this problem, we applied the emoticons polarity and utilized both lexicons with convolutional neural network (CNN) in a way which allows sentiment analysis on smaller annotated data sets. Our experimental results show that the proposed method can significantly improve the performance for predicting sentiment polarity on Weibo.

2. Related Research

In 2017, Felbo and others [Felbo 17] proposed a powerful system utilizing emoji in Twitter sentiment analysis model called DeepMoji. They trained 1,246 million tweets containing one of 64 common emoticons by Bi-directional Long Short-Term Memory (Bi-LSTM) model and applied it to interpret the meaning behind the online messages. DeepMoji is also one of the most advanced sarcasm-detecting model, sarcasm reverses the emotion of the literal text, therefore sarcasm-detecting capability can play a significant role in sentiment analysis, especially in case of social media. Although sarcasm and irony tend to convey negative emotions in general, we found that in Chinese social media (Weibo in our example), in addition to the expression of positive and negative emotions, people tend to express a kind of humorous emotion that escapes the traditional bi-polarity.

In our early research [Li 18b], we analyzed the usage of 67 emoticons with facial expression used on Weibo. We asked 12 Chinese native speakers to label these emoticons by applying one of three following categories: positive, negative and humorous. We have confirmed that 23 emoticons can be considered more as humorous than positive or negative. On this basis, we applied the emoticons polarity (see Table 1) in a Long Short-Term Memory recurrent neural network for sentiment analysis of undersized labelled data.

Chinese Internet slang is defined as an informal language used to express ideas on the Chinese Internet in response to events, to mass media and foreign cultures. It also expresses a natural human desire to simplify and update language. In [Li 19], we collected 448 frequent Internet slang expressions and created a slang lexicon (examples are shown in Table 2), then we converted the 109 Weibo emoticons into textual features creating Chinese emoticon lexicon (examples are shown in Table 3). To test the influence of slang and emoticons on sentiment analysis task, we also utilized both lexicons with several machine learning-based classifiers for detecting humorous expressions on Chinese social media.

3. Convolutional Neural Network Approach

Inspired by the above mentioned works on Internet slang and emoticons, in this paper, we utilized both lexicons and emoticon polarity with convolutional neural network (CNN) for sentiment classification of Weibo undersized labelled data.

In the first step, we added the Chinese slang lexicon and Chinese emoticon lexicon to segmentation tool for matching new words and emoticons. Then we used the updated tool to segment the sentences of a large data set. Secondly, we applied the segmentation results into the word embedding tool for training word vectors. Then, we applied the word embedding model which considered Internet slang and emoticons to train a CNN model with training data to learn an output representation. Next, we input testing data into deep learning model, and we use a softmax classifier to obtain the predicted results and output their probability. Since we assume that emoticons have relatively greater impact than textual features on emotional expression, we set a hyperparameter for each text and emoticon, and calculate the summation of both features' polarities with the hyperparameter. Finally, we can obtain the sentiment probability of a Weibo post which considers the effect of emoticons and Internet slang.

In order to verify the validity of our proposed method, we performed series of experiments described below.

3.1 Preprocessing

Initializing word vectors with those obtained from an unsupervised neural language model is a popular method to improve performance in the absence of a large supervised training set. For our experiment we collected a large dataset (7.6 million posts) from Weibo API from May 2015 to July 2017 to be used in calculating word embeddings. Firstly, we deleted the images, and videos treating them as noise. Secondly, we applied Chinese Internet slang lexicon and Chinese emoticon lexicon [Li 19] into the dictionary of Python Chinese word segmentation module Jieba^{*1}. Next, we used Jieba to segment the sentences of the microblogs, and applied the segmentation results into the word2vec model [Mikolov 13] for training word vectors. The vectors have dimensionality of 300 and were trained using the continuous skip-gram model.

Next, we collected 4,000 Weibo posts containing ambiguous (\bigcirc , , , , , , , , ,) emoticons, ensuring each entry has only one emoticon of given type (cases with more emoticons of the same type were allowed). To use these posts as our training data, we asked three Chinese native speakers to annotate them into three categories: "positive", "negative", and "humorous". After one annotator labelled polarities of all posts, two other native speakers confirmed correctness of his annotations. Whenever there was a disagreement, all decided the final polarity through discussion.

*1 https://github.com/fxsjy/jieba

Туре	Examples (Origin)	English Translation
Numbers	233 (哈哈哈)	"laughter"
Chinese contractions	人艰不拆 (人生已经如此的艰难 有些事情就不要拆穿)	"Life is so hard that some lies are better not exposed."
Slang derived from foreign language	欧尼酱 (お兄ちゃん)	"O-niichan" ("Brother" in Japanese)

Table 3: Examples of our Chinese Internet Slang Lexicon.

3.2 Convolutional Neural Network

Convolutional neural networks (CNN) utilize layers with convolving filters that are applied to local features [LeCun 98]. Originally invented for computer vision, CNN models have subsequently been shown to be effective for NLP and have achieved excellent results in semantic parsing [Yih 14], search query retrieval [Shen 14], sentence modeling [Kalchbrenner 14], and other traditional NLP tasks.

The equations of the CNN are as follows: $x_i \in \mathbb{R}^k$ is the kdimensional word vector corresponding to the *i*-th word in the sentence. A sentence of length n (padded where necessary) is represented as described in [Kim 14]:

$$x_{1:n} = x_1 \oplus x_2 \oplus \ldots \oplus x_n \tag{1}$$

where \oplus is the concatenation operator. In general, let $x_{i:i+j}$ refer

to the concatenation of words $x_i, x_{i+1}, ..., x_{i+j}$. A convolution operation involves a filter $w \in \mathbb{R}^k$, which is applied to a window of h words to produce a new feature. For example, a feature c_i is generated from a window of words $x_{i:i+h-1}$ by

$$c_i = f(wx_{i:i+h-1} + b)$$
 (2)

Here $b \in \mathbb{R}$ is a bias term and f is a non-linear function such as the

hyperbolic tangent. This filter is applied to each possible window of words in the sentence $\{x_{i:h}, x_{2:i+1}, \ldots, x_{n-h+1:n}\}$ to produce a feature map:

$$c = [c_1, c_2, \dots, c_{n-h+1}]$$
 (3)

with $c \in \mathbb{R}^{n-h+1}$. We then apply a max-overtime pooling oper-

ation [Collobert 11] over the feature map and take the maximum value $\hat{c} = max\{c\}$ as the feature corresponding to this particular filter. The idea is to capture the most important feature, one with the highest value, for each feature map. This pooling scheme naturally deals with variable sentence lengths. We experimented with the CNN architecture and a model with 10 epochs and the performance achieved the highest value when the dropout rate was 0.5, the filter size was 32 and strides number was 2. The validity of the model was examined by holdout method (90%/10%, training/validation). The activation functions, we used *RELU* in general, and the network output activation function was softmax.

3.3 Emoticon polarity

In order to predict sentiment category of Weibo posts considering the influence of emoticons for Chinese social media sentiment analysis, we assign the probability of the deep learning model's

Table 4: Comparison F scores results of three CNN approaches.

	Humorous	Negative	Positive
CNN	70.48%	66.67%	56.00%
CNN+Lexicons	70.59%	70.45%	70.58%
CNN+Lexicons+Polarities	76.84%	74.69%	73.56%

softmax output $S(z_i)$ a hyperparameter λ_1 . At the same time, we apply the labelled emoticons [Li 18b] as polarity P_e , and assign a hyperparameter λ_2 . P becomes the final probability output of the classification:

$$P = \lambda_1 S(z_i) + \lambda_2 P_e \tag{4}$$

where the summation of λ_1 and λ_2 is equal to 1.

3.4 Performance Test

Using a trained word2vec model, we passed word vectors of training data into the three deep learning models to train the model. We collected and annotated 180 Weibo entries with the eight emoticons mentioned above as a testing set, deleting images and videos. Then we used the proposed method to calculate probability of each category and confirmed the precision, recall and F1-score. Because we assumed that in emotion expression emoticons might play a greater role than text, in our experiment, we set the hyperparameters λ_1 and λ_2 to 0.4 and 0.6 respectively. We compared the results of sentiment classification by CNN only, CNN considering Internet slang and emoticons lexicons only, and our proposed methods: CNN model considering Internet slang and emoticons for shear the results of F1-score with above methods.

The results proved that our proposed method is more effective than a) convolutional neural network only categorization and b) convolutional neural network approach considering just slang and emoticon lexicons. Limited to small annotated data, the precision of the sentiment classification was relatively low, but by considering Internet slang and emoticons, the F1-score of each classifier outperformed previous method. Our proposed approach has improved the performance showing that low-cost, small-scale data labeling is sufficient to outperform widely used state-of-the-art when emoticon and slang information is added to the learning process.

4. Discussion

In our proposed approach, we focused on emoticons and Internet slang in microblogs and investigated how adding these features separately and together influences the previously proposed method for recognizing humorous posts which are problematic when it comes to semantic analysis.

Error analysis showed that some posts were wrongly predicted due to ambiguous usage of emoticon which brought clearly negative impact on the results. In Figure 2 we show an example of such misclassification into "positive" category annotated as "humorous" by annotators. ⁶⁹ was considered as more positive than humorous by our annotators (67%/0%/33%, positive/negative/humorous). It seems that this particular user wrote a joke just for fun, however, our proposed method was misguided by this "smirking" emoticon. Therefore, we plan to increase the number of evaluators for annotating Weibo emoticons in fine-grained humorous emotion to enhance the reliability of the polarity of emoticons.

Post: 吃饱了就有力气减肥了 😆 😆 Pinyin: Chi bao le jiu you li qi jian fei le 😆 😂 Segmentation: 吃饱了/就/有/力气/减肥/了/[阴险]/[阴险] Translation: When your stomach is full, you get the strength to reduce weight 😂 😂

Figure 2: Example of wrong classification into "positive" category.

5. Conclusions and Future Work

In this paper, we applied the emoticons polarity, Chinese Internet slang lexicon and Chinese emoticon lexicon with a convolutional neural network model for sentiment analysis of undersized labelled data. Our experimental results show that the proposed method can significantly improve the F1-score for predicting sentiment polarity on Weibo.

For improving the performance of our proposed method, in the near future we are going to increase the size of both slang lexicon and training dataset to improve further results. We also plan to test other deep learning approaches to compare the classification results.

6. Acknowledgment

This work was supported by JSPS KAKENHI Grant Number 17K00295.

References

- [Li 18a] Li, D., Rzepka, R., Araki, K.: Preliminary Analysis of Weibo Emojis for Sentiment Analysis of Chinese Social Media, Proceedings The 32th Annual Conference of the Japanese Society for Artificial Intelligence, 1J3-04 (2018)
- [Peng 17] Peng, H., Cambria, E., Hussain, A.: A review of sentiment analysis research in Chinese language, Cognitive Computation, 2017, 9(4): 423-435 (2017)
- [Wang 13] Wang, X., Zhang, C., Ji, Y.: A depression detection model based on sentiment analysis in micro-blog social network, Pacific-Asia Conference on Knowledge Discovery and

Data Mining. Springer, Berlin, Heidelberg, 2013: 201-213 (2013)

- [Moschini 16] Moschini, I.: The 'Face with Tears of Joy' Emoji, A Socio-Semiotic and Multimodal Insight into a Japan-America Mash-Up, HERMES-Journal of Language and Communication in Business, 2016 (55): 11-25 (2016)
- [Novak 15] Novak, P.K., Smailović, J., Sluban, B. et al.: Sentiment of emojis, PLOS One, 2015, 10(12): e0144296 (2015)
- [Guibon 16] Guibon, G., Ochs, M., Bellot, P.: From Emojis to Sentiment Analysis, WACAI 2016 (2016)
- [Kulkarni 17] Kulkarni, V., Wang, W.Y.: TFW, DamnGina, Juvie, and Hotsie-Totsie: On the Linguistic and Social Aspects of Internet Slang, arXiv preprint arXiv:1712.08291, 2017 (2017)
- [Felbo 17] B. Felbo, A. Mislove, A. Søgaard, et al. "Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm." arXiv preprint arXiv:1708.00524 (2017)
- [Li 18b] Li, D., Rzepka, R., Ptaszynski, M., Araki K.: Emoticon-Aware Recurrent Neural Network Model for Chinese Sentiment Analysis, in The Ninth IEEE International Conference on Awareness Science and Technology, iCAST 2018 (2018)
- [Li 19] Li, D., Rzepka, R., Ptaszynski, M., Araki K.: A Novel Machine Learning-based Sentiment Analysis Method for Chinese Social Media Considering Chinese Slang Lexicon and Emoticons, The AAAI-19 Workshop on Affective Content Analysis, AffCon 2019 (2019)
- [Mikolov 13] T. Mikolov, K. Chen, G. Corrado, et al. "Efficient estimation of word representations in vector space." arXiv preprint arXiv:1301.3781 (2013)
- [LeCun 98] LeCun Y., Bottou L., Bengio Y., et al.: Gradientbased learning applied to document recognition, Proceedings of the IEEE, 1998, 86(11): 2278-2324 (1998)
- [Yih 14] Yih W., He X., Meek C.: Semantic parsing for singlerelation question answering, Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). 2014, 2: 643-648 (2014)
- [Shen 14] Shen Y., He X., Gao J., et al.: Learning semantic representations using convolutional neural networks for web search, Proceedings of the 23rd International Conference on World Wide Web. ACM, 2014: 373-374 (2014)
- [Kalchbrenner 14] Kalchbrenner N., Grefenstette E., Blunsom P.: A convolutional neural network for modelling sentences, arXiv preprint arXiv:1404.2188 (2014)
- [Kim 14] Kim Y.: Convolutional neural networks for sentence classification, arXiv preprint arXiv:1408.5882 (2014)
- [Collobert 11] Collobert R., Weston J., Bottou L., et al.: Natural language processing (almost) from scratch, Journal of Machine Learning Research, 2011, 12(Aug): 2493-2537 (2011)