

Aicyber's System for IALP 2016 Shared Task: Character-enhanced Word Vectors and Boosted Neural Networks

Steven Du, Zhang Xi

www.aicyber.com

No. 39 Nanjing Road, Capitaland International Trade Center Block B 1604

Tianjin China

{steven,zhangxi}@aicyber.com

Abstract

This paper introduces Aicyber's system for the Dimensional Sentiment Analysis of Chinese Words ¹ in IALP 2016. The system is an ensemble of several boosted one layer neural networks, each one is trained on a different type of Chinese word vector. Our best system mainly use position-based character-enhanced word embedding and FastText as word vectors and achieve mean absolute error (MAE) 0.577(1st) in valence prediction and Pearson Correlation Coefficient (PCC) 0.671(1st) in arousal prediction.

1 Introduction

The purpose of this shared task is to predict a given traditional Chinese word's affective states in continuous numerical values (from 1 to 9) on valence-arousal space (Yu et al., 2016). It is a supervised learning task, 1653 human labeled training dataset is given, participants are asked to predict valence and arousal value of 1149 words in the test set. System performance is measured by MAE and PCC.

This paper will present Aicyber's system, from feature engineering, model building, evaluation towards final submission.

2 Feature Engineering

The training and testing data are all traditional Chinese words. A better representation is required to cover the semantic aspects of them. Several methods could learn word representation from text corpus.

For example the Latent Semantic Analysis, Latent Dirichlet Allocation, Vector Space Model, Explicit Semantic Analysis and the distributed word representation, which is also known as word2vec embedding (Mikolov et al., 2013). In this paper, we use two set of word embeddings derived from the word2vec approach.

2.1 Character-enhanced Word Embedding

The first set of word embedding is character-enhanced word embedding (Chen et al., 2015b) (CWE). Their work shows semantic meaning of a word is also related its composing characters. The effectiveness of new word vector is confirmed by evaluating on word relatedness computation and analogical reasoning. Two type of embeddings in CWE, the position-based character embeddings (CWE+P) and cluster-based character embeddings (CWE+L) are evaluated, but only CWE+P is included in submitted system due to a bug in script.

These two CWE models are trained with windows size of 5, 5 iterations, 5 negative examples, minimum word count of 5, Skip-Gram (Mikolov et al., 2013) with starting learning rate of 0.025, the output word vectors are of 300 dimensions.

2.2 FastText

Similar to CWE, (Bojanowski et al., 2016) proposes enriching word vectors with sub-word information. Where a word vector is associated to each character n-grams. The open sourced toolkit for their model is known as FastText ². It's the second embedding

¹http://nlp.innobic.yzu.edu.tw/tasks/dsa_w/

²<https://github.com/facebookresearch/fastText>

model used in submitted system.

FastText word embedding is trained with same setting as CWE training. Please noticed that default character n-gram for English is 3, we set it to 1 for Chinese.

2.3 Data for Word Embedding Training

We don't have any traditional Chinese text corpus, so word embedding is trained from simplified Chinese.

Following public datasets are used:

1. Chinese Wikipedia Dumps (Time stamp: 2011-02-05T03:58:02Z), however use of the latest dumps³ is encouraged.
2. Douban movie review⁴.
3. Aicyber synthesized 200 sentences⁵. These are intended to cover the out of vocabulary words.

Embedding trained on these datasets contains 445662 word vectors. Word not frequently used in mainland China are mapped to its synonyms in simplified Chinese, for example 强暴犯(Rapists) is mapped to 强奸犯(Rapists). After mapping, most of words in this task are in vocabulary, except ㄟ ㄟ. (Which indeed means hehe, an onomatopoeia, we noticed this after made submission.)

3 Prediction Model

In submitted system, boosted neural networks is adopted as regression method. It is a boosting system applied on neural networks. Both neural network and boosting method are from Scikit-learn (Pedregosa et al., 2011), an awesome machine learning package.

3.1 Neural Network

The neural network used in this work is very simple. It has only one hidden layer, and its size is 100, with *relu* (Glorot et al., 2010) activation function, *adam* (Kingma and Ba, 2014) as its training algorithm and a constant learning rate of 0.001.

3.2 Boosted Neural Networks

Boosting is a machine learning ensemble strategy to build a committee of learners that may be superior to a single learner (Freund, 1996; Drucker, 1997). In this experiment we propose use neural network as the fundamental building blocks, so named as Boosted Neural Networks (BNN). The selected implementation of boosting algorithm is known as AdaBoost.R2 (Drucker, 1997).

4 Evaluation

Evaluation is conducted locally on valence estimation only, the metric used is MAE, it is measured by mean value of 3 rounds of 10 folds cross-validation, each round has constant and unique random seed.

First, we evaluate the performance of different features then we fix the feature type and evaluate different regression methods.

4.1 Evaluation of Features

Besides three embeddings (CWP+P, CWP+L and FastText) mentioned in Section 2, we extend evaluation to include word2vec as the baseline feature. Two regression methods, the Linear Support Vector Regression (Fan et al., 2008; Ho and Lin, 2012)(LSVR) and proposed BNN are used to evaluate the goodness of word embedding features.

This work explores the size of word vector and two training schema, the continuous bag-of-words model (CBOW) (Mikolov et al., 2013) and Skip-Gram of different type of word embeddings.

In Table 1, MAE is measured on valence estimation based on LSVR. Regression target value is in range of 1 to 9. Table shows that in all the cases, Skip-Gram has better performance than CBOW. Higher dimension in Skip-Gram training leads to reduction in prediction error, but it doesn't guarantee better MAE for CBOW training. This is confirmed with Table 2, where the proposed boosted neural network is chosen as regression method. It is also clearly seen that new word vectors are better than baseline word2vec when BNN is applied, and the best feature are embeddings trained with Skip-Gram of 300 dimensions (S-300). Next we will exam different regression methods based on S-300.

³<https://dumps.wikimedia.org/zhwiki/>

⁴<http://www.datatang.com/data/45075>

⁵https://github.com/StevenLOL/ialp2016_Shared_Task/blob/master/data/extend.txt

Valence MAE by LSVR Baseline				
Embeddings	C-100	S-100	C-300	S-300
word2vec	0.936	0.839	0.950	0.819
CWE+P	0.940	0.773	0.940	0.765
CWE+L	0.953	0.827	0.923	0.769
FastText	1.093	0.796	1.343	0.765

Table 1: A LSVR is applied to different type of embedding features, grouped by training methods and size of feature vector. **C** denotes CBOW, **S** denotes Skip-Gram, **100/300** denotes embedding size is 100 or 300. Regression target value is range from 1 to 9. Reported MAE(the lower the better) is an average of 3 rounds of 10 folds cross-validation on training set. Overall, there was a general decrease in MAE when embeddings are trained with Skip-Gram on same size of vector. Increasing the size of embedding vector will improve MAE for Skip-Gram based embeddings.

Valence MAE by Boosted Neural Network				
Embeddings	C-100	S-100	C-300	S-300
word2vec	0.878	0.757	0.952	0.756
CWE+P	0.823	0.702	0.837	0.670
CWE+L	0.823	0.741	0.816	0.662
FastText	0.876	0.695	0.947	0.668

Table 2: Boosted neural network regression method applied to different type of features. This table can conclude that with same number of dimension, the Skip-Gram model is always better than CBOW. The best feature is embeddings trained with Skip-Gram with vector size of 300 dimensions (S-300).

4.2 Evaluation of Regression Methods

This section introduces the regression methods, training schema, evaluation results and issue with training.

Since we use the boosted neural network for this task. To evaluated boosting effect, the performance of ordinary neural work (NN) is required to be examined. Meanwhile two of regressors from boosting family are worth to be mentioned, gradient boosting machine (Friedman, 2000; Friedman, 2002) (GBM) and the eXtreme Gradient Boosting (Chen et al., 2015a) (XGB).

During 10 fold cross-validation, regression methods are not trained in the same way. LSVR and GB are trained on 90% of training data. XGB, NN, BNN further divided this 90% of training data into 90% training and 10% validation set, their training will stop once there is no improvement obtained on validation set, this strategy is know as early-stopping.

Table 3 presents the results of above regression methods along with LSVR baseline. It's quite obvious that BNN obtain better MAE than ordinary NN, both BNN and NN are better than baseline LSVR.

Embedding S-300	Valence MAE by Different Regression Methods				
	LSVR	GBM	XGB	NN	BNN
word2vec	0.819	0.829	0.881	0.801	0.756
CWE+P	0.765	0.757	0.809	0.729	0.670
CWE+L	0.769	0.795	0.860	0.730	0.662
FastText	0.765	0.791	0.847	0.711	0.668

Table 3: Evaluation of different regression methods applied to S-300 embedding. BNN and NN outperform LSVR baseline. Meanwhile BNN is not only better than ordinary NN, but also better than two boosting approaches, the GBM and XGB.

Embedding S-300-PCA100	Valence MAE		Valence PCC	
	BNN	BNN_Norm	BNN	BNN_Norm
word2vec	0.702	0.686	0.858	0.867
CWE+P	0.678	0.623	0.874	0.891
CWE+L	0.671	0.627	0.873	0.891
FastText	0.662	0.639	0.879	0.889

Table 4: This table presents MAE and PCC (the higher the better) of valence estimation by BNN and normalized BNN (BNN_Norm) approach. Features are S-300 after PCA projected to 100 dimensions. The BNN_Norm achieved better MAE and PCC than BNN.

GBM and XGB didn't yield superior results to baseline method. We will further investigate the performance of XGB as it is a well known reliable large-scale tree boosting system (Chen and Guestrin, 2016), ruling leader-board of many machine learning competitions(Holloway et al., 2016).

One problem with BNN is that it has the longest running time, to make it run faster, a Principal Component Analysis (Jolliffe, 1986) (PCA) is applied. This reduces the dimensions of word embedding from 300 to 100, with a trade-off in MAE. As demonstrated in Table 4, after PCA, word2vec and FastText have better MAE, but CWE+P and CWE+L are degraded. We found that normalize target value by removing mean and scaling to unit variance (BNN_Norm), will lead to notable improvement for both MAE and PCC.

Through evaluation, the best feature and regression methods are found. The final submission Run1 is an average of BNN_Norm applied on CWE+P and FastText, if CWE+L is included in the ensemble, better result could achieve. Run2 is generated by BNN_Norm on CWE+P feature alone.

4.3 Discussion

For Run1 submission the MAE obtained on valence testing set is 0.577. It's even better than 3 out of 5 human annotators mentioned in (Yu et al., 2016). But the PCC value is only 0.848, it is not as good

as expected. We are still investigating this problem. Anyone interested in solving this problem or try to reproduce our results could refer to our public repository⁶.

5 Conclusion

This system paper had covered Aicyber’s system for Dimensional Sentiment Analysis of Chinese words. During feature engineering CWE+P, CWE+L and FastText word-embedding are identified as the best features. Then local cross-validation demonstrated the effectiveness of BNN. To tackle training speed problem, PCA is applied on embeddings. The target normalization further improved system performance. Submitted system achieved MAE 0.577(1st) in valence estimation and PCC 0.671(1st) in arousal estimation.

References

- Liang Chih Yu, Lung Hao Lee, Shuai Hao, Jin Wang, Yunchao He, Jun Hu, K Robert Lai, and Xuejie Zhang. 2016. Building chinese affective resources in valence-arousal dimensions. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Xinxiong Chen, Lei Xu, Zhiyuan Liu, Maosong Sun, and Huanbo Luan. 2015b. Joint learning of character and word embeddings. In *International Conference on Artificial Intelligence*.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching word vectors with subword information.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in Neural Information Processing Systems*, 26:3111–3119.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12(Oct):2825–2830.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2010. Deep sparse rectifier neural networks. *Journal of Machine Learning Research*, 15.
- Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *Computer Science*.
- Harris Drucker. 1997. Improving regressors using boosting techniques.
- Rong En Fan, Kai Wei Chang, Cho Jui Hsieh, Xiang Rui Wang, and Chih Jen Lin. 2008. Liblinear: A library for large linear classification. *Journal of Machine Learning Research*, 9(9):1871–1874.
- Y. Freund. 1996. Experiments with a new boosting algorithm. In *Thirteenth International Conference on Machine Learning*, pages 148–156.
- Jerome H. Friedman. 2000. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5):1189–1232.
- Jerome H. Friedman. 2002. Stochastic gradient boosting. *Computational Statistics & Data Analysis*, 38(4):367–378.
- Chia Hua Ho and Chih Jen Lin. 2012. Large-scale linear support vector regression. *Journal of Machine Learning Research*, 13(1):3323–3348.
- Tianqi Chen, Tong He, and Michael Benesty. 2015a. xgboost: Extreme gradient boosting.
- Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. *arXiv:1603.02754*.
- Eric Holloway, II Marks, et al. 2016. High dimensional human guided machine learning. *arXiv:1609.00904*.
- I. T. Jolliffe. 1986. Principal component analysis. *Springer Berlin*, 87(100):41–64.

⁶https://github.com/StevenLOL/ialp2016_Shared_Task