

TFIDF meets Deep Document Representation: A Re-Visit of Co-Training for Text Classification

Zhiwei Chen and Aixin Sun

ABSTRACT

Many text classification tasks face the challenge of lack of sufficient labelled data. Co-training algorithm is a candidate solution, which learns from both labeled and unlabelled data for better classification accuracy. However, two sufficient and redundant views of an instance are often not available to fully facilitate co-training in the past. With the recent development of deep learning, we now have both traditional TFIDF representation and deep representation for documents. In this paper, we conduct experiments to evaluate the effectiveness of co-training with different combinations of document representations (e.g., TFIDF, Doc2vec, ELMo, BERT) and classifiers (e.g., SVM, Random Forest, XGBoost, MLP, and CNN) on two benchmark datasets (20 Newsgroup and Ohsumed). Our results show that co-training with TFIDF and deep contextualised representation offers improvement to classification accuracy.

1 INTRODUCTION

We often face the situation of inadequate labeled data to learn accurate text classification models. It is often expensive to obtain more labeled data. Among many possible solutions, the co-training algorithm is a solution designed to utilize both labeled and unlabeled data. Co-training assumes that each instance is described by two “views” (or two feature sets) [1]. The two views are expected to be conditionally independent and each view is sufficient to predict an instance’s class label. With a small number of labeled examples, two weak classifiers are first trained, each trained with one view. Both weak classifiers are used to predict the labels of unlabeled instances and the instances with most confident predictions are added as labeled instances to learn better classifiers. The co-training algorithm shows effectiveness in web page classification where a web page can be described by its content and also the anchor words pointing to the web page.

In practice, it is however often hard to find two views of an instance which are conditionally independent and both are sufficient for classification. Traditionally, a document is typically represented by Bag-of-Words (BoW) model in a vector space. Each word (i.e., a feature) is weighted by TFIDF scheme (i.e., term frequency \times inverse document frequency). For easy presentation, we refer to this document representation as TFIDF. Recently, with the rapid development of representation learning, we are now offered multiple different views of a document in low dimensional embedding spaces.

Examples range from Word2Vec based representations to deep contextualized representations like Embeddings from Language Models (ELMo) [10] and Bidirectional Encoder Representations from Transformers (BERT) [3].

Both TFIDF and word embedding based representations are sufficient to represent a document and they are from different feature spaces. This motivate us to re-visit co-training algorithm, to explore whether the two sets of document representations facilitate us to utilize unlabeled documents in text classification. To this end, we conduct experiments on two benchmark datasets, namely 20News-group and Ohsumed, using five different document representations: TFIDF, Doc2Vec [7] which is based on pre-trained word2vec embeddings, and three contextualized representations, i.e., ELMo, BERT, and Universal Sentence Encoder (USE) [2]. The two weak classifiers for each co-training setting, are selected from Support Vector Machine (SVM), Random Forest (RF), XGBoost (XGB), Multi-Layer Perceptron (MLP), and Convolutional Neural Network (CNN).

Experimental results suggest that deep contextualised representations together with the right choice of classification models is likely to achieve the best results. Nevertheless, TFIDF+SVM remains a strong baseline. Combination of different document representations and classification models give very diverse results, particularly when there is a lacking of sufficient training data. The different document representations do facilitate co-training to achieve significantly better classification performance when both weak classifiers could give relatively good performance with limited labeled data. On the other hand, on challenging dataset, when the weak classifiers could not achieve reasonable accuracy with small amount of training data, co-training leads to degradation in final predictions.

2 RELATED WORK

Co-training has been applied to web page classification in the original paper. The algorithm has also been applied to classify email using header and body as the two views [6]. Clustering and co-training [11] use text content and features derived from clusters of labeled and unlabeled data as the two views. Ghani [4] proposed to combine error correcting output coding and co-training for multi-class text classification. Co-training has also been extended to tri-training [14].

More relevant to our work is multi-co-training [5]. The authors exploit three document representations (TFIDF, LDA, and Doc2Vec) for multi-class text classification. The proposed method is similar to tri-training, where three weak classifiers are trained each with one document representation, by the same classification model. The authors evaluated Naïve Bayes and Random Forest as base classifiers. However, the objective of multi-co-training is to maximise the benefits of different document representations. Our work is to systematically evaluate the effectiveness of co-training with combinations of 5 document representations and 5 classification models.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICTIR’20, Sep 2020, Stavanger, Norway

© 2020 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

3 METHODOLOGY

In the original paper [1], co-training was introduced for a binary classification task. In our experiments, both datasets are for multi-class text classification. Before reporting our results, we briefly describe the co-training setting in our experiments.

Let L denote the set of labeled data; let U denotes the set of unlabeled data. Let $U' \in U$ denote a subset of U with size $|U'|$. Recall that we evaluate co-training with 5 different document representations and 5 classifiers. In each co-training setting, we choose two document representations x_1 and x_2 , and two classification models. Then all labeled instances in L are used to train two weak classifiers h_1 and h_2 , using representations x_1 and x_2 respectively. The trained weak classifier h_1 then predicts the category labels of documents in U' , and the top p documents with the highest confidence scores from U' are added to L as labeled data. The same applies to h_2 . The instances that have been added to L will be removed from U' , and U' is replenished with remaining data from U . This is one iteration. The process repeats for k iterations.

After k iterations, h_1 and h_2 are trained with the updated L . Then both h_1 and h_2 make predictions on the test data. A combined classifier h_3 is then defined, which makes prediction basing the outputs of h_1 and h_2 . In our implementation h_3 estimates probability $P(c_j|x)$ of class c_j given the instance $x = (x_1, x_2)$ by multiplying the probabilities derived from outputs h_1 and h_2 .

3.1 Document Representation

In our evaluation, a document can have the following representations: $x \in \{TFIDF, Doc2Vec, ELMo_s, ELMo_p, USE, BERT_s, BERT_p\}$. Among them, TFIDF is traditional vector space representation, Doc2Vec is based on pre-trained word2vec, where the embedding for a word is fixed. ELMo, USE, and BERT are deep contextualize representations, where the embeddings of the same word could be different depending on the context of the word.

TFIDF represents a document as a feature vector where each word is one feature and is weighted by its TFIDF value. We use scikit-learn's `TFIDFVectorizer`¹ to derive TFIDF vectors for documents.

Doc2Vec is an extension of Word2Vec [7, 9]. Although Doc2Vec is designed to make use of some contextual information during training, the training is based on pre-trained word vectors. We adopt pre-trained Doc2Vec embeddings² where the word vectors are learned from Wikipedia. Document vector dimension is 300.

USE. Universal-sentence-encoder is an extension of Transformer [2, 12]. We adopt pre-trained USE model.³ Each document is represented in a 512 dimension vector.

BERT applies bi-directional training of transformer to language modelling [3, 12]. We used pre-trained BERT-base⁴ for 20 Newsgroup. BERT-base generates two kinds of embeddings: one is *sequence* output, a 128 x 768 matrix for each document (document length is capped at 128), the other is *pooled* output, where a document is represented by a 768 dimension vector. Because Ohsumed

dataset is in medical domain, we use pre-trained **BioBERT** model⁵. BioBERT is pre-trained on large-scale biomedical corpora [8].

ELMo is also a language model that gives contextualised word embeddings [10]. We use pre-trained ELMo model⁶. Similar to BERT, this model also generates two kinds of embeddings (sequence and pooled), except that the word dimension is 1028 instead of 768.

3.2 Classification Model and Setting

To learn weak classifiers (h_1 and h_2) in co-training, we evaluated the following classifiers: Support Vector Machines (SVM), Random Forest classifier (RF), XGBoost (XGB), Multi-Layer Perceptron (MLP), and Convolutional Neural Network (CNN). These models are commonly used as baselines in many classification tasks and we used the implementations available in mainstream packages.⁷ Note that, CNN is only used for sequential word embeddings generated by BERT/ELMo, so as to learn the dependencies between words.

In our experiments, we simply adopt the pre-trained embeddings as detailed earlier, without finetuning the embeddings for the classification task. For classification models in our experiments, we adopt either default or commonly adopted parameters. The reason is, in co-training setting, we only has a small number of labeled examples to start with, which makes parameter finetuning less effective. Further, during each iteration, two new classifiers are learned.

4 EXPERIMENTS AND RESULTS

Datasets. We conduct experiments on two benchmark datasets: 20 Newsgroup and Ohsumed. The 20 NG dataset contains 11,314 labeled documents for training and 7,532 documents for test, in 20 classes. Ohsumed dataset contains abstracts of medical research papers. Some of its documents has multiple class labels. We follow the same preprocessing as in [13] to exclude all documents belonging to multiple classes. As the results, we 3,357 training documents and 4,043 test documents, in 23 classes.

Co-Training Setting. To evaluate co-training, we randomly sample 10% of labeled documents from each class to form the labeled set L . The remaining 90% of training data becomes U . We set the size of U' to be 400. After each iteration, the top p most confident documents from U' are added to L . We set $p = 1$ for the first 10 iterations, and increase the value of p by 1 after every 10 iterations. After $k = 40$ iterations, we have in total 100 documents added to L . Then we stop co-training and evaluate the combined prediction of the two weak classifiers.⁸

4.1 Results with 100% and 10% Labeled Data

To provide reference performances for co-training, we first evaluate the different classifiers trained on different document representations, using labeled documents in L (i.e., 10% of training data in the original dataset) and documents in $L+U$ (i.e., 100% training data in the original dataset) respectively. We report the classification performance by macro-averaged Precision/Recall/ F_1 and Accuracy. The results on the two benchmark datasets are reported in Table 1. In our discussion, we mainly focus on F_1 score.

¹https://scikit-learn.org/Class:sklearn.feature_extraction.text.TfidfVectorizer

²<https://github.com/jhlau/Doc2Vec>

³<https://tfhub.dev/google/universal-sentence-encoder/4>

⁴https://tfhub.dev/google/bert_uncased_L-12_H-768_A-12/1

⁵<https://github.com/naver/biobert-pretrained>

⁶<https://tfhub.dev/google/elmo/3>

⁷Code will be released upon acceptance.

⁸We have also evaluated another setting $p = 1$ and $k = 100$, and the results are comparable.

Table 1: Performance of different combinations of document representation and classification model, on 20Newsgroup and Ohsumed. Two sets of results are presented one is trained with 10% of training data, i.e., documents in L , and the other is trained with all training data, i.e., documents in $L + U$ in the co-training setting. BERT refers to BioBERT on Ohsumed dataset.

Dataset		20NG 100% labeled data				20NG 10% labeled data				Ohsumed 100% labeled data				Ohsumed 10% labeled data			
DocRep.	Model	Pre	Rec	F_1	Acc	Pre_L	Rec_L	F_{1L}	Acc_L	Pre	Rec	F_1	Acc	Pre_L	Rec_L	F_{1L}	Acc_L
TFIDF	SVM	0.78	0.77	0.77	0.78	0.69	0.68	0.68	0.69	0.68	0.56	0.59	0.67	0.49	0.26	0.28	0.47
	MLP	0.78	0.77	0.77	0.77	0.69	0.67	0.67	0.68	0.69	0.52	0.57	0.66	0.43	0.20	0.21	0.42
	RF	0.62	0.59	0.59	0.60	0.57	0.44	0.43	0.45	0.40	0.26	0.28	0.49	0.23	0.10	0.08	0.29
	XGB	0.64	0.62	0.63	0.63	0.39	0.35	0.35	0.36	0.42	0.33	0.34	0.52	0.04	0.06	0.03	0.02
Doc2Vec	SVM	0.53	0.53	0.52	0.54	0.41	0.41	0.40	0.42	0.39	0.37	0.37	0.48	0.30	0.21	0.21	0.38
	MLP	0.54	0.53	0.53	0.54	0.49	0.48	0.46	0.49	0.51	0.38	0.38	0.54	0.28	0.20	0.20	0.40
	RF	0.44	0.44	0.42	0.45	0.35	0.33	0.31	0.34	0.29	0.13	0.11	0.34	0.15	0.10	0.07	0.29
	XGB	0.43	0.43	0.42	0.44	0.31	0.31	0.29	0.31	0.26	0.18	0.18	0.38	0.10	0.10	0.09	0.27
USE	SVM	0.70	0.70	0.69	0.71	0.66	0.67	0.66	0.68	0.47	0.35	0.36	0.50	0.30	0.21	0.20	0.38
	MLP	0.70	0.69	0.69	0.70	0.67	0.66	0.65	0.68	0.42	0.33	0.34	0.48	0.25	0.20	0.19	0.38
	RF	0.68	0.68	0.67	0.69	0.64	0.64	0.62	0.65	0.48	0.24	0.25	0.44	0.20	0.15	0.14	0.33
	XGB	0.43	0.43	0.42	0.44	0.58	0.58	0.57	0.59	0.34	0.28	0.29	0.45	0.16	0.15	0.14	0.33
BERT _s	CNN	0.80	0.80	0.80	0.81	0.68	0.67	0.67	0.68	0.67	0.55	0.59	0.68	0.25	0.20	0.18	0.40
BERT _p	SVM	0.66	0.65	0.65	0.66	0.54	0.54	0.53	0.55	0.55	0.52	0.52	0.61	0.40	0.34	0.35	0.49
	MLP	0.58	0.51	0.50	0.53	0.37	0.29	0.24	0.30	0.58	0.53	0.54	0.62	0.41	0.28	0.29	0.46
	RF	0.38	0.38	0.36	0.39	0.51	0.50	0.48	0.51	0.61	0.25	0.27	0.47	0.18	0.14	0.12	0.34
	XGB	0.42	0.43	0.42	0.44	0.29	0.29	0.27	0.29	0.40	0.30	0.31	0.49	0.15	0.14	0.12	0.32
ELMo _s	CNN	0.76	0.76	0.76	0.76	0.62	0.62	0.61	0.63	0.49	0.39	0.40	0.57	0.21	0.17	0.16	0.37
ELMo _p	SVM	0.64	0.64	0.64	0.65	0.57	0.57	0.57	0.58	0.36	0.33	0.34	0.47	0.20	0.20	0.20	0.36
	MLP	0.68	0.66	0.66	0.67	0.61	0.55	0.54	0.56	0.37	0.33	0.32	0.48	0.20	0.19	0.18	0.37
	RF	0.58	0.58	0.57	0.66	0.51	0.50	0.48	0.51	0.25	0.33	0.18	0.40	0.17	0.14	0.12	0.32
	XGB	0.58	0.58	0.58	0.59	0.47	0.47	0.46	0.48	0.27	0.23	0.22	0.41	0.13	0.13	0.12	0.32

Performance on both datasets share similar trend. All classifiers show significant drop of performance when only 10% of labeled data are used for training, compared to the version using all training documents. Overall, Ohsumed is much more challenging to get high classification results, compared to 20NG. In terms of classifiers' performance, when all training data are used, BERT_s+CNN achieves the highest F_1 on 20NG, followed by TFIDF+SVM and ELMo_s+CNN.⁹ With 10% of training data, TFIDF+SVM becomes the best performer followed by BERT_s+CNN. On Ohsumed, with all training data are available, both TFIDF+SVM and BERT_s+CNN achieve the best F_1 . Both outperform the rest by a large margin. When only 10% of data is used for training, BERT_p+SVM is the best performer with F_1 score 0.35. TFIDF+SVM is in the second position with F_1 score 0.28. We note that, the F_1 scores by the best performers are fairly low.

Our results suggest that deep contextualised representations together with the right choice of classification models likely to achieve the best results. On the other hand, TFIDF+SVM remains a strong baseline. There is also an observation that the combination of different document representations and classification models give very diverse results, particularly when there is a lacking of sufficient training data.

⁹The best performance reported in our experiments is not as good as SOTA, as we do not finetune parameters and document representations. Further, our focus is to study the performance of co-training and not to achieve the new SOTA. On the other hand, the best results obtained by using all training data are not too far from SOTA results.

4.2 Results of Co-Training

As results of different classifiers on the same dataset are very diverse (see Table 1). Some of them give very poor results, hence it is not necessary to conduct co-training on every possible combination. Instead, on each dataset, for each kind of document representation, we choose the best performing classification model. Then we conduct co-training on this set of classifiers. With 5 document representations, we have 10 combinations of h_1 and h_2 for co-training.

Table 2 reports the performance of 10 co-training settings on both datasets. In this table, we list the two classifiers (h_1 and h_2), and report macro-averaged Precision/Recall/ F_1 and accuracy of the combined classifier h_3 , after co-training. Again, we mainly focus on F_1 . For comparison, we compute the change of F_1 obtained by h_3 , compared to the lower bound, denoted by ΔF_{1L} . The lower bound is determined by the F_{1L} of the better classifiers among h_1 and h_2 , trained by using 10% of training data (i.e., set L). Accordingly ΔF_1 denotes the difference between h_3 's F_1 and the upper bound, which is the F_1 of the better classifier among h_1 and h_2 , training by using all training data ($L + U$). We also list the error rate of h_1 and h_2 , which is the rate of wrong predictions made by the weak classifier when adding documents to L . From the results, we make the following observations.

On 20Newsgroup datasets, five classifier combinations in co-training setting lead to significant improvement in F_1 score. In particular $\langle \text{TFIDF+SVM}, \text{USE+SVM} \rangle$ and $\langle \text{TFIDF+SVM}, \text{BERT}_s+\text{CNN} \rangle$ and both manage to improve the F_1 score to 0.72. And another two

Table 2: Co-training results of combined classifier h_3 . Results are in bold if the performance of h_3 is better than either h_1 or h_2 in the evaluated combination. Paired t -test is conducted on F_1 scores only (not on Acc) and * denotes statistically significant.

	h_1	ER_{h1}	h_2	ER_{h2}	Pre	Rec	F_1	$\Delta F1_L$	ΔF_1	Acc	ΔAcc_L	ΔAcc
20Newsgroup	TFIDF+SVM	0.05	Doc2Vec+MLP	0.09	0.60	0.56	0.55	-0.13	-0.22	0.57	-0.12	-0.21
		0.05	USE+SVM	0.01	0.73	0.72	0.72*	+0.04	-0.05	0.74	+0.05	-0.04
		0.03	BERT _s +CNN	0.02	0.72	0.72	0.72*	+0.04	-0.08	0.73	+0.04	-0.08
		0.07	ELMo _s +CNN	0.02	0.66	0.65	0.65	-0.03	-0.12	0.66	-0.03	-0.12
	Doc2Vec+MLP	0.07	USE+SVM	0.00	0.59	0.57	0.56	-0.12	-0.13	0.58	-0.10	-0.13
		0.11	BERT _s +CNN	0.01	0.71	0.69	0.69*	+0.01	-0.11	0.71	+0.03	-0.10
		0.15	ELMo _s +CNN	0.01	0.63	0.61	0.60	-0.06	-0.16	0.62	-0.06	-0.14
	USE+SVM	0.00	BERT _s +CNN	0.02	0.73	0.72	0.71*	+0.04	-0.09	0.73	+0.06	-0.08
		0.00	ELMo _s +CNN	0.01	0.67	0.66	0.65	-0.03	-0.11	0.67	-0.01	-0.09
	BERT _s +CNN	0.03	ELMo _s +CNN	0.03	0.72	0.71	0.71*	+0.03	-0.09	0.72	+0.04	-0.09
Ohsumed	TFIDF+SVM	0.12	Doc2Vec+MLP	0.33	0.27	0.23	0.21	-0.07	-0.38	0.44	-0.03	-0.23
		0.13	USE+SVM	0.26	0.39	0.22	0.21	-0.07	-0.38	0.42	-0.05	-0.25
		0.11	BERT _p +SVM	0.19	0.45	0.36	0.37	+0.09	-0.22	0.52	+0.03	-0.15
		0.10	ELMo _s +CNN	0.25	0.20	0.19	0.17	-0.11	-0.11	0.41	-0.06	-0.26
	Doc2Vec+MLP	0.24	USE+SVM	0.23	0.24	0.21	0.20	0.00	-0.18	0.42	+0.02	-0.12
		0.19	BERT _p +SVM	0.22	0.37	0.28	0.27	-0.08	-0.25	0.48	-0.01	-0.13
		0.29	ELMo _s +CNN	0.20	0.22	0.20	0.19	-0.01	-0.38	0.42	+0.02	-0.15
	USE+SVM	0.26	BERT _p +SVM	0.18	0.42	0.31	0.32	-0.03	-0.20	0.49	0.00	-0.12
		0.26	ELMo _s +CNN	0.22	0.25	0.21	0.18	-0.02	-0.39	0.41	+0.03	-0.16
	BERT _p +SVM	0.31	ELMo _s +CNN	0.32	0.24	0.24	0.22	-0.13	-0.30	0.44	-0.05	-0.17

combinations achieve F_1 score 0.71. Comparing to the F_1 score of a single classifier obtained by using all labeled documents, the improvement is promising. We note that the error rate is fairly low for both h_1 and h_2 on 20Newsgroup dataset.

On Ohsumed, only one combination (\langle TFIDF+SVM, BERT_p+SVM) manages to get a positive $\Delta F1_L$. And the improvement is not statistically significant by paired t -test with two-tails. In other words, co-training results on this dataset are in general worse than the weak classifiers h_1 and h_2 trained on limited labeled data L . We also note that the error rate for both classifiers are fairly high.

In summary, the evaluated co-training combinations show very different results on the two datasets. On a typical text classification task like 20Newsgroup, we observe improvements on 5 out of 10 combinations of classifiers after co-training. The classification models in each combination leading to improvement, could be the same or different. However, in all the 5 combinations leading to improvement, the weak classifiers h_1 and h_2 are both among the better ones, compared to the other classifiers in the comparison (see Table 1). On a challenging classification task like Ohsumed. We do not observe the benefit of co-training, probably due to the weak performance of weak classifiers h_1 and h_2 . With a relatively high error rate, the quality of documents added to L is not ensured, leading to the classifiers in the following iterations to learn from noisy labeled data.

5 CONCLUSION

In this study, we conduct a systematic evaluation of co-training by using different document representations and classification models. Our study show that traditional TFIDF representation and deep learning based representation both are sufficient to represent documents for classification task, making co-training applicable. Co-training achieves significant improvement if the classification task

is not very challenging and weak classifiers are able to get reasonable good predictions with limited training data. The traditional setting of TFIDF+SVM remains a strong baseline.

REFERENCES

- [1] Avrim Blum and Tom M. Mitchell. 1998. Combining Labeled and Unlabeled Data with Co-Training. In *COLT*. 92–100.
- [2] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Universal Sentence Encoder. *CoRR* abs/1803.11175 (2018).
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT*. 4171–4186.
- [4] Rayid Ghani. 2002. Combining Labeled and Unlabeled Data for MultiClass Text Categorization. In *ICML*. 187–194.
- [5] DongHwa Kim, Deokseong Seo, Suhyoung Cho, and Pilsung Kang. 2019. Multi-co-training for document classification using various document representations: TF-IDF, LDA, and Doc2Vec. *Inf. Sci.* 477 (2019), 15–29.
- [6] Svetlana Kiritchenko and Stan Matwin. 2011. Email classification with co-training. In *Proc CASCON*. 301–312.
- [7] Quoc V. Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents. In *ICML*, Vol. 32. 1188–1196.
- [8] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *CoRR* abs/1901.08746 (2019).
- [9] Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. In *NeurIPS*. 3111–3119.
- [10] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations. In *NAACL-HLT*. 2227–2237.
- [11] Bhavani Raskutti, Herman L. Ferrá, and Adam Kowalczyk. 2002. Combining clustering and co-training to enhance text classification using unlabelled data. In *KDD*. 620–625.
- [12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *NeurIPS*. 5998–6008.
- [13] Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Graph Convolutional Networks for Text Classification. In *AAAI*. 7370–7377.
- [14] Zhi-Hua Zhou and Ming Li. 2005. Tri-Training: Exploiting Unlabeled Data Using Three Classifiers. *IEEE TKDE* 17, 11 (2005), 1529–1541.