

Semi-supervised Content-based Detection of Misinformation via Tensor Embeddings

Gisel Bastidas Guacho
UC Riverside
gbast001@ucr.edu

Sara Abdali
UC Riverside
sabda005@ucr.edu

Neil Shah
Snap Inc.
nshah@snap.com

Evangelos E. Papalexakis
UC Riverside
epapalex@cs.ucr.edu

Abstract—Fake news may be intentionally created to promote economic, political and social interests, and can lead to negative impacts on humans beliefs and decisions. Hence, detection of fake news is an emerging problem that has become extremely prevalent during the last few years. Most existing works on this topic focus on manual feature extraction and supervised classification models leveraging a large number of labeled (fake or real) articles. In contrast, we focus on content-based detection of fake news articles, while assuming that we have a *small* amount of labels, made available by manual fact-checkers or automated sources. We argue this is a more realistic setting in the presence of massive amounts of content, most of which cannot be easily fact-checked. So, we represent collections of news articles as multi-dimensional tensors, leverage tensor decomposition to derive concise article embeddings that capture spatial/contextual information about each news article, and use those embeddings to create an article-by-article graph on which we propagate limited labels. Results on real-world datasets show that our method performs on par or better than existing fully supervised models, in that we achieve better detection accuracy using fewer labels. In particular, our proposed method achieves 75.43% of accuracy using only 30% of labels of a public dataset while an SVM-based classifier achieved 67.43%. Furthermore, our method achieves 70.92% of accuracy in a large dataset using only 2% of labels.

Index Terms—Fake news, tensor decomposition, semi-supervised learning, belief propagation.

I. INTRODUCTION

Misinformation on the web is a problem that has been greatly amplified by the use of social media, and the problem of fake news in particular has become ever more prevalent during the last years. Social media is a common platform for consuming and sharing news, due to its ease-of-use in diffusing content and promoting exposure/discussion. In fact, two-thirds of Americans reported getting some of their news from social media in 2017¹. Even though social media has become a news source for its advantages, it is especially vulnerable to the propagation of fake news mostly coming from unverified publishers and crowd-based content creators because there is practically no control over the information that is shared. The well-documented spread of misinformation on Twitter during events such as Hurricane Sandy in 2012 [1], the Boston Marathon blasts in 2013 [2] and US Presidential Elections on Facebook in 2016 [3] are all such examples. Since

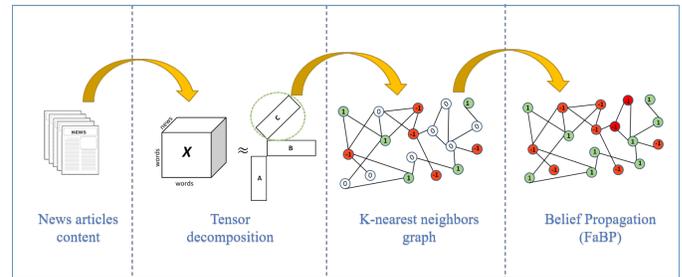


Fig. 1. Our proposed method discerns real from misinformative news articles via leveraging tensor representation and semi-supervised learning in graphs.

misinformation is intentionally created for malicious purposes such as obtain economic and political benefits or deceiving the public [4], it can clearly lead to negative user experience by either influencing their beliefs and impacting their decisions for the worse. Several approaches in recent literature have been proposed to automatically detect misinformation using supervised classification models.

Some works extract manually crafted features from news content such as the number of nouns, length of the article, fraction of positive/negative words, and more in order to discriminate fake news articles [5]–[7].

In addition to these works, several others proposed propagation-based models for evaluating news credibility [8]–[10]. Nonetheless, they initialized credibility values for the entire network using a supervised classifier. However, the reality is that such labels are often very limited and sparse. Fact-checking websites such as Snopes.com, PolitiFact.com, and FactCheck.org can be used to assess claims, but these websites require domain experts to assign credibility values to claims and are therefore, limited by human capacity. Moreover, fact-checking is a time-consuming process, often requiring surveying multiple articles and sources, evaluating reputation and likelihood of the claims before coming to a decision.

In this paper, we propose a new *semi-supervised* approach for fake news detection based on news content, which requires *limited* labels. On a high level, our approach exploits tensor representation and decomposition of news articles, careful construction of a k -nearest neighbor graph, and propagation of limited labeled article information to conduct inference on a larger set.

Our main contributions are:

¹<http://www.journalism.org/2017/09/07/news-use-across-social-media-platforms-2017/>

- We leverage tensor-based article embeddings, which are shown to produce concise representations of articles with respect to their spatial context, in order to derive a graph representation of news articles.
- We formulate fake news detection as a semi-supervised method that propagates known labels on a graph to determine unknown labels.
- We collect a large dataset of misinformation and real news articles publicly shared on social media.
- We evaluate our method on real datasets. Experiments on two previously used datasets demonstrate that our method outperforms prior works since it requires a fewer number of known labels and achieves comparable performance.

II. PROBLEM DEFINITION

We consider a misinformative, or fake, news article as one that is “*intentionally and verifiably false*”, following the definition used in [4]. With this definition in mind, we aim to discern fake news articles from real ones based on their content. Henceforth, by “content”, we refer to the text of the article. We reserve the investigation of other types of content (such as image and video) for future work.

Let $\mathcal{N} = \{n_1, n_2, n_3, \dots, n_M\}$ be a collection of news articles of size M where each news article is a set of words and $\mathcal{D} = \{w_1, w_2, w_3, \dots, w_I\}$ be a dictionary of words of size I . Note that articles can have varying length. Assuming that labels of some news articles are available. Let $l \in \{-1, 0, 1\}$ denote a vector containing the partially known labels, such that entries of 1 represent real articles, -1 represents fake articles and 0 denotes an unknown status. We address the problem as a binary classification problem; hence, a news article is classified either fake or real.

III. PROPOSED METHOD

Our proposed method consists of the following steps:

Step 1: Tensor Decomposition We build similar tensor-based article embeddings as proposed in [11]. Specifically, we propose the use of *binary-based tensor* construction method. That is, we build a three-mode tensor $\mathcal{X} \in \mathbb{R}^{I \times I \times M}$ (*words, words, news*) where for each news article, we create a co-occurrence matrix where all co-occurrence entries are boolean and indicate ($word_1, word_2$) appeared within a window parameter of w (5-10) words² at least once. We then use CP/PARAFAC tensor decomposition [12] to factorize the tensor. As [11] demonstrates, such tensor-based article embeddings captures spatial/contextual nuances of different types of news articles and result in homogeneous article groups. After decomposing the tensor, we obtain the factor matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$ whose columns correspond to different latent topics, clustering news articles and words in the latent topic space. More specifically, each row of \mathbf{C} is the representation of the corresponding article in the resulting embedding space.

²We experimented with small values of that window and results were qualitatively similar.

Step 2: k -NN graph of news articles The k -nearest-neighbors of a point in n -dimensional space are defined using a “closeness” relation where proximity is often defined in terms of a distance metric [13] such as Euclidean ℓ_2 distance. We use the factor matrix \mathbf{C} in order to construct a k -NN graph G of news articles. As we mentioned before, each column in \mathbf{C} is the representation of the corresponding news article in the latent topic space; thus, by constructing a k -NN graph on \mathbf{C} , we can find similar articles in that space. So, we consider each row in $\mathbf{C} \in \mathbb{R}^{M \times R}$ as a point in R -dimensional space. We then compute ℓ_2 distance among news and find the k -closest points for each point in \mathbf{C} .

Step 3: Belief Propagation Using the graphical representation of the news articles above, and considering that for a small set of those news articles we have ground truth labels, our problem becomes an instance of semi-supervised learning over graphs. We use a belief propagation algorithm which assumes homophily, because news articles that are connected in the k -NN graph are likely to be of the same type due to the construction method of the tensor embeddings. More specifically, we use the fast and linearized FaBP variant proposed in [14].

IV. EXPERIMENTAL EVALUATION

We implemented our method in MATLAB using Tensor Toolbox [15] and MATLAB FaBP implementation [14].

A. Dataset description

We use the following datasets:

Public datasets The two public datasets were used in previous studies. Specifically, *Dataset1* consists of 150 political news articles, balanced to have 75 articles of each class, and was provided by [7]. *Dataset2* contains 68 real and 69 fake news articles, and was provided by [5]. Our dataset contains 31,739 articles from different fake news categories such as Fake, Conspiracy, Rumor, Satire and Junk Science.

Our dataset In constructing our dataset, we collected news article URLs from Twitter tweets during a 3-month period from June-August 2017. These URLs were filtered based on website domain. We then crawled those URLs to get the news article content. To that end, we used web API boilerpipe³, Python library Newspaper3k⁴, and Diffbot⁵. All real news articles were featured on 367 domains obtained from Alexa⁶, and fake news articles belong to 367 domains from the BSDetector browser extension domain list [16].

B. Evaluation

For evaluation, we measured Accuracy, Precision, Recall and $F1_score$. In order to find the best-performing parameters for our method, we run an iterative process using cross-validation where we evaluated different settings with respect to R (i.e. decomposition rank) and k (i.e. the number of nearest neighbors, controlling the density of the k -NN graph). We

³<http://boilerpipe-web.appspot.com/>

⁴<http://newspaper.readthedocs.io/en/latest/>

⁵<https://www.diffbot.com/dev/docs/article/>

⁶<https://www.alexa.com/>

considered values of R from 1 to 20, since decomposition rank is often set to be low for time and space reasons in practice [17]. Likewise, we tested k with values from 1 to 100, trading off greater bias for less variance with increasing k . We found that the best accuracy is obtained when both parameters R and k are set to be 10. We find that for values of k and R greater than 10, performance is qualitatively similar as shown in Figure 2, and thus we fix the parameters as such in evaluation. Notice that using a small k value (e.g. 1 or 2), the accuracy is relatively poor; this is because building a k -NN graph with small k results in a highly sparse graph which offers limited propagation capacity. In all experiments, we tested accuracy over the test set of all articles whose labels were “unknown” or unspecified in the propagation step. We evaluated our method with different percentages p of

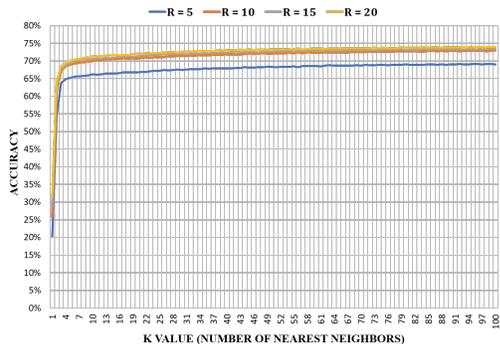


Fig. 2. Performance using different parameter settings for decomposition rank (R) and number of nearest neighbors (k).

known labels. Table I shows the performance of our method using $p \in \{5\%, 10\%, 20\%, 30\%\}$ of labeled news articles from our dataset. Our results demonstrate that we can achieve an accuracy of 70.76% only using 10% of labeled articles. We also evaluated the performance of our approach using extremely sparse known labels. That is, we evaluated our method using $p < 5\%$ and varying the number of nearest neighbors. Figure 3 shows that we can achieve an accuracy of 70.92% using 2% of known labels when the number of nearest neighbors is set to be 200. In fact, the performance of our approach degrades fairly gracefully with even smaller proportions of known labels.

TABLE I
PERFORMANCE OF THE PROPOSED METHOD USING OUR DATASET WITH DIFFERENT PERCENTAGES OF LABELED NEWS.

%Labels	Accuracy	Precision	Recall	F1
5%	69.12 ± 0.003	69.09 ± 0.004	69.24 ± 0.009	69.16 ± 0.004
10%	70.76 ± 0.003	70.59 ± 0.003	71.13 ± 0.010	70.85 ± 0.004
20%	72.39 ± 0.001	71.95 ± 0.002	73.32 ± 0.004	72.63 ± 0.002
30%	73.44 ± 0.001	73.13 ± 0.003	74.14 ± 0.003	73.63 ± 0.001

Additionally, to evaluate the quality of tensor embeddings over traditional vectorial representations, we compared performance between our approach and a variant in which between

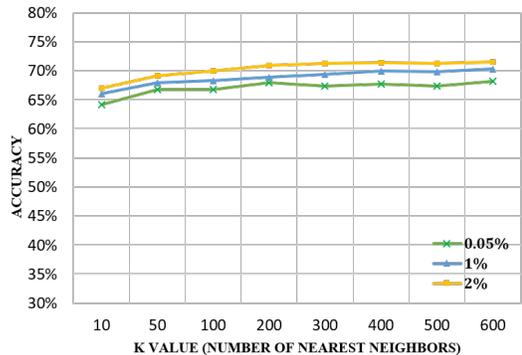


Fig. 3. Performance using extremely sparse (<5%) labeled articles and varying number of nearest neighbors.

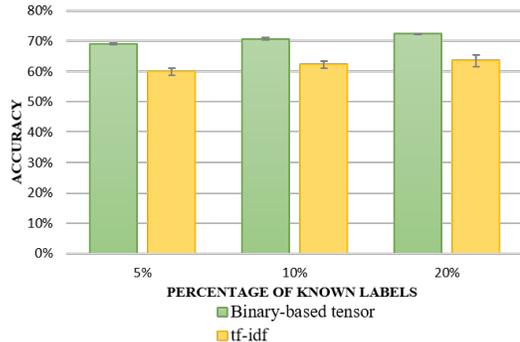


Fig. 4. Performance of our method using tensor-based article embedding compared to using a graph built from $tf-idf$ matrix.

we constructed a k -NN graph built from the term frequency inverse-document-frequency ($tf-idf$) representations. Figure 4 shows that our method with tensor embeddings consistently attains better accuracy than the alternative over varying known label percentages. This empirically suggests that binary-based tensor representations can better capture spatial/contextual nuances of news articles over vectorial representations. In addition, we evaluated our model using *Dataset1* and *Dataset2*. We compare the accuracy achieved by our method to the accuracy achieved by the following approaches:

SVM on content-based features as proposed in [7]. To this extent, we replicated the feature extraction from news content and used SVM in order to show the performance using different percentages of training data.

Logistic regression on content-based features proposed by [5]. We used their publicly available implementation. In particular, we run their method with linguistic (n -gram) feature extraction using different percentages of training data.

Figure 5 shows the results for *Dataset1*. Our approach demonstrates improved accuracy even with fewer labels – specifically, we achieved 75.43% accuracy using only the 30% of news labels while SVM(30%/70% train/test), SVM(5-fold cross-validation), and logistic regression (30%/70% train/test) attained 67.43%, 71% and 50.09% of accuracy, respectively. The accuracy achieved by SVM(5-fold cross-validation) was

reported by Horne et al. in [7]. For *Dataset2*, we run logistic regression and SVM, using 10%/90% train/test split. These approaches achieved an accuracy of 59.84% and 64.79%, respectively, compared to the 67.38% accuracy achieved by our approach, using the same percentage of labeled articles.

Our method is able to achieve this performance only having a small number of labeled news articles due to the quality of the tensor embeddings which define a favorable graph where the node labels are propagated.

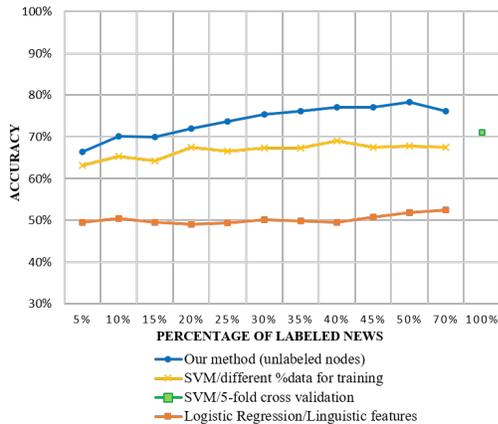


Fig. 5. Performance using *Dataset1* provided by Horne et al. [7]

V. RELATED WORK

In [5], the authors proposed a logistic regression classifier using linguistic (n -gram), credibility (punctuation, pronoun use, capitalization) and semantic features generated from the news content. In another work, authors used SVM on content-based features that are categorized into stylistic, complexity and psychological features in order to classify real, fake and satirical news [7]. In [18], the authors propose detecting rumors by building naïve-Bayes classifiers on content, network and microblog-specific features. [19] and [20] leverage temporal structure by using recurrent neural network (RNN) based models to represent text and user characteristics. In [21], the authors propose a Dynamic Series-Time Structure (DSTS) model for detecting rumors by capturing the social context of an event from content, user and propagation-based features.

VI. CONCLUSIONS

We propose a semi-supervised content-based method for detecting fake news articles, leveraging tensor-based article embeddings and guilt-by-association. Extensive experiments on over 63K real articles demonstrate that our method distinguishes fake from real news only using a fraction of labeled articles, performing on par or better than state-of-the-art.

ACKNOWLEDGMENT

Research was supported by a gift from Snap Inc, an Adobe Data Science Faculty Award, and by the Department of the Navy, Naval Engineering Education Consortium under award no. N00174-17-1-0005. Any opinions, findings, and conclusions or recommendations expressed here are those of the author(s) and do not necessarily reflect the views of the funding parties.

REFERENCES

- [1] A. Gupta, H. Lamba, P. Kumaraguru, and A. Joshi, "Faking sandy: Characterizing and identifying fake images on twitter during hurricane sandy," in *Proceedings of the 22Nd International Conference on World Wide Web, WWW '13 Companion*, (New York, NY, USA), pp. 729–736, ACM, 2013.
- [2] A. Gupta, H. Lamba, and P. Kumaraguru, "\$1.00 per rt #bostonmarathon #prayforboston: Analyzing fake content on twitter," in *2013 APWG eCrime Researchers Summit*, pp. 1–12, Sept 2013.
- [3] C. Silverman, "This analysis shows how fake election news stories outperformed real news on facebook.," BuzzFeed News, "November" 2016.
- [4] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake news detection on social media: A data mining perspective," *CoRR*, vol. abs/1708.01967, 2017.
- [5] M. Hardalov, I. Koychev, and P. Nakov, "In search of credible news," *Artificial Intelligence: Methodology, Systems, and Applications. AIMSA 2016 Lecture Notes in Computer Science*, p. 172180, 2016.
- [6] V. L. Rubin, N. J. Conroy, Y. Chen, and S. Cornwell, "Fake news or truth? using satirical cues to detect potentially misleading news," 2016.
- [7] B. D. Horne and S. Adali, "This just in: Fake news packs a lot in title, uses simpler, repetitive content in text body, more similar to satire than real news," *CoRR*, vol. abs/1703.09398, 2017.
- [8] M. Gupta, P. Zhao, and J. Han, *Evaluating Event Credibility on Twitter*, pp. 153–164.
- [9] Z. Jin, J. Cao, Y. Zhang, and J. Luo, "News verification by exploiting conflicting social viewpoints in microblogs," 2016.
- [10] Z. Jin, J. Cao, Y. G. Jiang, and Y. Zhang, "News credibility evaluation on microblog with a hierarchical propagation model," in *2014 IEEE International Conference on Data Mining*, pp. 230–239, Dec 2014.
- [11] S. Hosseinimotlagh and E. E. Papalexakis, "Unsupervised content-based identification of fake news articles with tensor decomposition ensembles," 2017.
- [12] R. A. Harshman, "Foundations of the PARAFAC procedure: Models and conditions for an "explanatory" multi-modal factor analysis," *UCLA Working Papers in Phonetics*, vol. 16, no. 1, p. 84, 1970.
- [13] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 3rd ed., 2011.
- [14] D. Koutra, T.-Y. Ke, U. Kang, D. Chau, H.-K. Pao, and C. Faloutsos, "Unifying Guilt-by-Association Approaches: Theorems and Fast Algorithms," in *Machine Learning and Knowledge Discovery in Databases (ECML/PKDD)*, vol. 6912 of *Lecture Notes in Computer Science*, pp. 245–260, 2011.
- [15] T. G. K. B. W. Bader et al., "Matlab tensor toolbox version 2.6. Available online," February 2015.
- [16] B.S. Detector. <http://bsdetecter.tech/>, 2017.
- [17] N. Shah, A. Beutel, B. Gallagher, and C. Faloutsos, "Spotting suspicious link behavior with fbox: An adversarial perspective," in *Data Mining (ICDM), 2014 IEEE International Conference on*, pp. 959–964, IEEE, 2014.
- [18] V. Qazvinian, E. Rosengren, D. R. Radev, and Q. Mei, "Rumor has it: Identifying misinformation in microblogs," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*, (Stroudsburg, PA, USA), pp. 1589–1599, Association for Computational Linguistics, 2011.
- [19] J. Ma, W. Gao, P. Mitra, S. Kwon, B. J. Jansen, K.-F. Wong, and M. Cha, "Detecting rumors from microblogs with recurrent neural networks," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI'16*, pp. 3818–3824, AAAI Press, 2016.
- [20] N. Ruchansky, S. Seo, and Y. Liu, "CSI: A hybrid deep model for fake news," *CoRR*, vol. abs/1703.06959, 2017.
- [21] J. Ma, W. Gao, Z. Wei, Y. Lu, and K.-F. Wong, "Detect rumors using time series of social context information on microblogging websites," in *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management, CIKM '15*, (New York, NY, USA), pp. 1751–1754, ACM, 2015.