

Tackling skewed data in online fake review detection

Satheesh Joseph
satheeshrishi@berkeley.edu

Catherine Mou
catherine041616@berkeley.edu

Yi Zhang
yizhang7210@berkeley.edu

Abstract—The increasing need and value of the digital text and the trustworthiness behind it inspired our imagination to handle a niche problem within the space of opinion spam analysis. Recently many researches have been conducted to increase the detection accuracy. To the best of our knowledge, an important issue that has been less studied is to tackle heavily skewed data in online fake review detection. In this paper, we study a methodology to tackle the imbalanced data issue and the accuracy effects on different state-of-the-art models. For this research, we use publicly available Yelp review data for hotels and restaurants for our experiments.

I. INTRODUCTION

Online fake review detection is a relatively well-studied subject, especially in recent years, given its outsized impact on consumers engaging in e-commerce activities online. Positive reviews bring a meaningful increase in sales volume to the products [3], and vice versa for negative reviews.

As a result, there has been an increase in opinion spamming activities, and detecting fake reviews has become an essential requirement for online marketplaces to maintain the integrity and fairness of their platform.

However, there has been one persisting challenge [8], [9], [10], [13] in this area of research – the lack of a substantial body of actually proven fake reviews, directly leading to significantly imbalanced datasets.

Our work aims to tackle this problem of data imbalance by borrowing ideas from the Generative Adversarial Network (GAN). Our **hypothesis** is that it's possible to make up the data imbalance by generating fake reviews from a **language model** trained and/or fine-tuned on actual fake reviews. We will look to validate this approach by then training a review detection model on a balanced dataset that includes the generated fake

reviews and achieving comparable results to state-of-the-art research [9].

Section II provides the background on the existing research in fake review detection. Section III lays out our methods of constructing the fake review generator. Section IV discusses the experiments we run to validate the usefulness of the generated fake reviews in training a detection model. And we draw our final conclusions in section V.

II. BACKGROUND

There is no shortage of research tackling the problem of fake review detection. A recent survey [4] does a great job laying out the landscape of the various techniques and data sets used for fake review detection.

However, most of the publicly available data sets are very skewed with many much real review examples than fake ones. According to this survey, all large datasets (>20k reviews) from Yelp [7] contain less than 15% actual fake reviews. There are a few other widely used public datasets crawled from TripAdvisor, but they are of a much smaller scale, with the fake review training set generated via a manual process from Amazon Mechanical Turk.

The only balanced dataset of moderate volume is crawled from Yelp by Barbado et. al. [2], however it has not been widely adopted in the research community as a benchmark for detecting fake reviews.

This lack of trusted labeled data can apply a significant limit on the performance of models that try to detect fake reviews. Moreover, the unbalanced datasets tend to encourage the model to produce biased predictions. Given this state of the related work, and inspired by Stanton et. al. [8] who used GAN techniques to generate behavioral features (e.g. number of reviews, percentage

of positive reviews) for online Yelp reviewers, we believe that similar techniques can be used to generate the reviews themselves.

With sufficient representativeness, we believe the generated fake reviews can serve as additional training examples that can help with the detection model to distinguish between genuine reviews and fake ones. To paraphrase Tolstoy – genuine reviews are all alike; every fake review is fake in its own way.

III. METHODS

Generative Adversarial Network (GAN) was originally used for image data augmentation, but the development of Transformer-based models has enabled coherent and human-like *text* generation, and those models had demonstrated to be effective for NLP data augmentation to overcome the data scarcity problem.

The general structure of a Generative Adversarial Network consists of a *generative* network that generates candidates as well as a *discriminative* network that evaluates them. Here we apply this technique to generate fake review texts.

For the *generative* network, we here use the well-established, pre-trained GPT-2 language model.

GPT-2 was pre-trained on web scraped industrial-scale corpus, it was trained with a batch size of 512, well-defined sentence length, vocabulary size of 50,000, and can be used for various NLP tasks including text generation. It is a statistical tool to generate the next words in sequence based on preceding words, at every stage, it will take the previously generated data as additional input when generating the next output. GPT-2 has outperformed other language models when it comes to generating text based on small input contents like hotel/restaurant reviews. GPT-2 has the ability to adapt to the context of the text, so it can generate realistic and coherent output.

To generate domain-specific fake Yelp reviews, we finetuned the pre-trained GPT-2 model by adding layers on top and training with our raw data set. GPT-2 text generation pipeline predicts words that will follow a specified text prompt, so we gave it a simple prompt such as “the hotel is”, or “the restaurant is”.

However, this led to low-quality generated text that was very repetitive at the end. Furthermore, the generated text was frequently not full sentences. With multiple trials and errors, we eventually made a few improvements that led to higher quality generated fake reviews.

- We sampled some of the original actual fake reviews as the prompt supplied to GPT-2 for text generation, instead of the short static prompt. This is to give GPT-2 more context to generate longer and more realistic-looking reviews.
- We trained a simple Neural Network using ELMo embedding with LSTM and 2 dense layers on the original data set as the *discriminative* network. Even though this Network itself has not performed well in distinguishing genuine and fake reviews, it does do a good job distinguishing coherent, relevant reviews about hotels/restaurants and irrelevant sentences, making it a good candidate for the discriminative network in the GAN. And finally, we only picked reviews that sufficiently confused this network (with a predicted probability of more than 80%) to be included in the final generated fake review set.

IV. EXPERIMENT & DISCUSSION

A. Experimental Setup

For our experiments, we use the publicly available data set originally obtained by [6] and [5]. This data set, containing 5858 reviews for hotels and 67019 reviews for restaurants on Yelp, is also used by a number of prior research papers for benchmarking, notably [11] and [9]. A more detailed statistics of the raw data set is in Table I. As mentioned above, it is indeed highly imbalanced with a much smaller number of fake reviews.

TABLE I. SUMMARY OF EXPERIMENTAL DATA

Subject	Hotels	Restaurants
Total # Reviews	5858	67019
Total # Genuine Reviews	5078	58716
Total # Fake Reviews	780	8303
% Fake Reviews	13.3%	12.4%

To validate our hypothesis on the usefulness of the generated fake reviews, we set up a 2-stage experiment.

Firstly, we use the methods laid out in section III to generate fake reviews for both hotels and restaurants. Secondly, we add the generated fake reviews to the training data to obtain a balanced training set, then run a number of classification models on the mixed training set and compare it against our benchmark results.

We construct the data sets in the following ways to be comparable with the prior research papers [11], [9]. A summary is given in Table II.

For a balanced test set:

- we first limit the pool of reviews to the first review per reviewer after 2012-01-01
- we take all the fake reviews (because there are fewer)
- we sample the same number of reviews from the genuine reviews
- this gives us a balanced, non-duplicated test set

For a balanced training set:

- we first limit the pool of reviews to be the ones prior to 2012-01-01
- we take all the actual fake reviews
- we include all generated fake reviews
- we sample the same number of reviews from the genuine reviews
- this gives us a balanced, non-duplicated training set

TABLE II. SUMMARY OF TRAINING/TEST DATA

Subject	Hotels	Restaurants
Training set size	5070	
# Genuine Reviews	2535	
# Actual Fake Reviews	561	
# Generated Reviews	1974	
% Fake Reviews in training set	50%	
Test set size	432	
% Fake Reviews in test set	50%	

B. Models

We established a number of models all based on Neural Networks to compare our results with the benchmark.

Model 1: Our baseline model with GloVe embedding and 1 layer of LSTM.

Model 2: Our main model with GloVe embedding, 1 layer of Bidirectional LSTM, and 3 dense layers.

Model 3: A BERT-based model.

For each model, we'll be training on 4 different training sets for comparison:

Set 1: raw, imbalanced training set

Set 2: balanced training set by under-sampling genuine reviews to the number of fake reviews

Set 3: balanced training set by over-sampling fake reviews with replacement to the number of genuine reviews

Set 4: balanced training set by including generated fake reviews per Table II

C. Results

We report the best results of each of our models trained on each training set listed above after hyperparameter tuning via grid search. We will also report the results of a few state-of-the-art classification algorithms from the benchmark research [9].

Their brief descriptions are as follows:

LF: uses linguistic features from the review content only, by extracting bigrams from the reviews data.

CNN: uses the same bigram features but is trained using a Convolutional Neural Network.

LF+BF: is a concatenation of linguistic as well as *behavioral* features from the review, including its length, rating, and other reviews by the same reviewer.

bfGAN: is the state-of-the-art algorithm using a number of generated behavioral features as well as the review content.

We've selected classification *accuracy* as well as *f1-score* as the reporting metrics to be compatible with prior research. A summary of them all the experimental results is presented in Table III for the Hotels data set and Table IV for the restaurant data set.

TABLE III. EXPERIMENTAL RESULTS - HOTELS

accuracy/f1	Set 1	Set 2	Set 3	Set 4
Model 1	0.51/0.35	0.61/0.61	0.58/0.56	0.59/0.59
Model 2	0.56/0.62	0.60/0.61	0.57/0.62	0.59/0.60
Model 3	0.5/0.34	0.57/0.57	0.53/0.54	0.52/0.56
LF [9]	0.56/0.62			
CNN [9]	0.59/0.56			
LF + BF [9]	0.61/0.58			
bfGAN [9]	0.83/0.83			

TABLE IV. EXPERIMENTAL RESULTS - RESTAURANTS

accuracy/f1	Set 1	Set 2	Set 3	Set 4
Model 1	0.57/0.53	0.63/0.62	0.62/0.61	0.60/0.60
Model 2	0.57/0.50	0.63/0.62	0.62/0.61	0.59/0.58
Model 3	0.5/0.33	0.52/0.46	-	-
LF [9]	0.56/0.65			
CNN [9]	0.57/0.58			
LF + BF [9]	0.59/0.60			
bfGAN[9]	0.76/0.75			

D. Discussion

WIP: - Model size matters, with the data set, we can't afford models that are too large, they can easily overfit

- Non-repetitive, balanced data set matters, Training set 1 and 3 perform notably worse than the training set 2 and 4

- Non-linguistic features matter. Our model is doing better than the pure bigram-based model and is comparable to CNN and LF+BF-based models, but if we take into account other behavioral features, it should continue to improve.

V. CONCLUSION

REFERENCES

[1] <https://github.com/huggingface/transformers>.

- [2] Rodrigo Barbado, Oscar Araque, and Carlos A Iglesias. A framework for fake review detection in online consumer electronics retailers. *Information Processing & Management*, 56(4):1234–1244, 2019.
- [3] Nga N Ho-Dac, Stephen J Carson, and William L Moore. The effects of positive and negative online customer reviews: do brand strength and category maturity matter? *Journal of marketing*, 77(6):37–53, 2013.
- [4] Rami Mohawesh, Shuxiang Xu, Son N. Tran, Robert Ollington, Matthew Springer, Yaser Jararweh, and Sumbal Maqsood. Fake reviews detection: A survey. *IEEE Access*, 9:65771–65802, 2021.
- [5] Arjun Mukherjee, Vivek Venkataraman, Bing Liu, and Natalie Glance. What yelp fake review filter might be doing? In *Seventh international AAAI conference on weblogs and social media*, 2013.
- [6] Arjun Mukherjee, Vivek Venkataraman, Bing Liu, Natalie Glance, et al. Fake review detection: Classification and analysis of real and pseudo reviews. *UIC-CS-03-2013. Technical Report*, 2013.
- [7] Shebuti Rayana and Leman Akoglu. Collective opinion spam detection: Bridging review networks and metadata. In *Proceedings of the 21th acm sigkdd international conference on knowledge discovery and data mining*, pages 985–994, 2015.
- [8] Gray Stanton and Athirai A Irissappane. Gans for semi-supervised opinion spam detection. *arXiv preprint arXiv:1903.08289*, 2019.
- [9] Xiaoya Tang, Tiejun Qian, and Zhenni You. Generating behavior features for cold-start spam review detection with adversarial learning. *Information Sciences*, 526:274–288, 2020.
- [10] Jingdong Wang, Haitao Kan, Fanqi Meng, Qizi Mu, Genhua Shi, and Xixi Xiao. Fake review detection based on multiple feature fusion and rolling collaborative training. *IEEE Access*, 8:182625–182639, 2020.
- [11] Xuepeng Wang, Kang Liu, and Jun Zhao. Handling cold-start problem in review spam detection by jointly embedding texts and behaviors. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 366–376, 2017.
- [12] Jason Wu. How to build an ai text generator: Text generation with a gpt-2 model, 2020.
- [13] Chunyuan Yuan, Wei Zhou, Qianwen Ma, Shangwen Lv, Jizhong Han, and Songlin Hu. Learning review representations from user and product level information for spam detection. In *2019 IEEE International Conference on Data Mining (ICDM)*, pages 1444–1449. IEEE, 2019.