

# Adversarial Texts with Gradient Methods

Zhitao Gong\*, Wenlu Wang\*, Bo Li†, Dawn Song†, Wei-Shinn Ku\*

\* {gong, wenluwang, weishinn}@auburn.edu

Auburn University, Auburn, AL, USA

† {crystalboli, dawnsong}@berkeley.edu

University of California, Berkeley, Berkeley, CA, USA

**Abstract**—Adversarial samples for images have been extensively studied in the literature. Among many of the attacking methods, gradient-based methods are both effective and easy to compute. In this work, we propose a framework to adapt the gradient attacking methods on images to text domain. The main difficulties for generating adversarial texts with gradient methods are: (i) the input space is discrete, which makes it difficult to accumulate small noise directly in the inputs, and (ii) the measurement of the quality of the adversarial texts is difficult. We tackle the first problem by searching for adversarial texts in the embedding space and then reconstruct the adversarial texts via nearest neighbor search. For the latter problem, we employ the Word Mover’s Distance (WMD) to quantify the quality of adversarial texts. Through extensive experiments on three datasets, IMDB movie reviews, Reuters-2 and Reuters-5 newswires, we show that our framework can leverage gradient attacking methods to generate very high-quality adversarial texts that are only a few words different from the original texts. There are many cases where we can change one word to alter the label of the whole piece of text. We successfully incorporate FGM and DeepFool into our framework. In addition, we empirically show that WMD is closely related to the quality of adversarial texts.

## I. INTRODUCTION

It has been shown that carefully crafted noise may trick the deep neural networks into wrong predictions with very high confidence [25]. Many followup work proposed cheap yet effective methods to find such adversarial noise, e.g., fast gradient method (FGM) [4], Jacobian-based saliency map approach (JSMA) [19], DeepFool [17], CW [3], etc. We have seen that the investigation of adversarial samples for image models has provided us new perspectives to understand the mechanism of deep neural networks, e.g., linear hypothesis [4], rethinking smoothness assumptions [25]. In addition, many algorithms have been proposed to enhance the robustness of the deep models, e.g., adversarial training [8]. However, most of the previous work focused on images. Only a few attempts have been made in text domain. We think it is worthwhile to investigate adversarial samples for text models as well. In this work, we propose a simple and effective framework to adapt the adversarial attacking methods for images to text domain. Specifically, we focus on gradient-based method since they are very fast in practice. There are two major problems we need to resolve before we can plugin the gradient based methods for generating adversarial images.

- 1) The input space is discrete. As a result, it is not possible to accumulate small noise computed with gradient methods directly in the input space. They work well in image

domain since image models usually take input in a continuous domain  $[0, 1]$ .

- 2) It is difficult to quantify the quality of adversarial texts. For adversarial images, we usually utilize  $p$ -norm distance, perceptual distance [10], etc. While for adversarial texts, there are no good metrics to measure the quality.

In this work, we propose a general framework in which we generate adversarial texts via slightly modified gradient-based attacking methods. We first search for adversarial texts in the text embedding space [15] via gradient-based methods and then reconstruct the adversarial texts via nearest neighbor search. In addition, we also empirically evaluate using Word Mover’s Distance (WMD) [9] as a quality measurement for the adversarial texts. The advantage of our framework is that no manual features are needed.

This paper is organized as follows. We briefly review recent work on generating adversarial images and texts in Section II. Our adversarial text framework is detailed in Section III. We evaluate our method on various text benchmarks and report the results in Section IV. We conclude this paper and provide directions for future work in Section V.

## II. RELATED WORK

The existence of adversarial samples was first discussed in [25]. There has been an abundance of work on attacking methods to generate adversarial images. These adversarial images raise security concerns about the wide application of deep neural networks [8]. As a result, many work have investigated defense against these adversarial samples. However, so far as we see in literature, the attacking is much easier and cheaper than defense.

For notation,  $x$  denotes the input,  $y$  the prediction,  $f$  the target model such that  $y = f(x)$ ,  $L$  the loss function,  $x^*$  the adversarial sample.  $\|\cdot\|_p$  denotes the  $p$ -norm. We slightly abuse the notation here,  $L_x$  denotes the loss with  $x$  as the input.

The attacking methods mainly fall into three categories, *gradient attack*, *optimization attack* and *model attack*. Generally speaking, gradient attack is faster than the others. However, the other two require much less knowledge about the model, thus more practical. In addition, optimization attacks are usually more effective and generate more subtle noise.

### A. Adversarial Image Method

1) *Gradient Attack*: This class of attacking methods rely on target model gradients, thus requiring full knowledge of

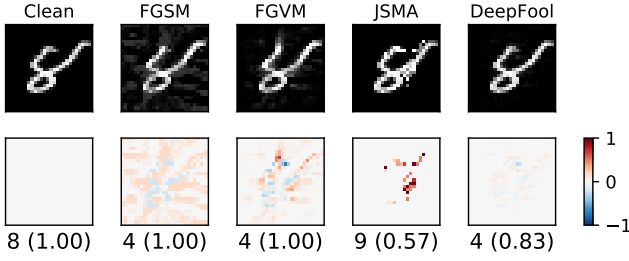


Fig. 1: Random MNIST adversarial images generated via different attacking algorithms. The upper image in *Clean* column is the original clean image. The upper images in the following columns are adversarial images generated by the corresponding attacking algorithm based on the first clean image, respectively. The lower image in each column is the difference between the adversarial image and the clean image, illustrated in heatmap. Below each column is the label predicted by the target model, along with probability in parenthesis.

the target model. Fast gradient method (FGM) [4] and its variants, iterative FGM and targeted FGM [8], add to the whole image the noise that is proportional to either  $\nabla L_x$  (FGVM) or  $\text{sign}(\nabla L_x)$  (FGSM). Jacobian-based saliency map approach (JSMA) [19], on the contrary, perturbs one pixel at a time. It chooses the pixel with the highest saliency score, which is calculated as  $-\nabla y_t \cdot \nabla y_o$  subject to  $\nabla y_t > 0$ , where  $y_t$  is the probability for the target class, and  $y_o$  is the sum of probabilities of all other classes. Intuitively, JSMA tries to increase the probability of the target class while decreasing others. DeepFool [17] iteratively finds the optimal direction in which we need to *travel* the minimum distance to cross the decision boundary of the target model. Although in non-linear case, the optimality is not guaranteed, in practice, however, DeepFool usually finds very subtle noise compared to other gradient methods.

Figure 1 shows adversarial image examples of four gradient methods on MNIST. As we can see, FGSM tends to generate more salient noise spread across the whole image. On the other hand, FGVM is slightly better since it uses gradients instead of the sign of gradients as noise. In practice, most of the absolute values of gradients are far less than 1. JSMA, on the contrary, increases the intensity of the most salient pixel until its value goes beyond the input domain. As a result, we expect to see a few very intense spots in the image. DeepFool, as shown in the last image, generates the most subtle noise.

2) *Optimization Attack*: This class of attacks is usually black-box attacks, using the target model only as an oracle. Generally speaking, this class minimizes  $\|x^* - x\|_p$  subject to  $f(x^*) \neq f(x)$  and  $x^*$  is in the input domain. Following this formulation, this is a box-constrained optimization problem, L-BFGS is used in [25] to solve it. In order to utilize more optimization methods, CW [3] proposes a refinement of the above formulation by minimizing  $\|s(x^*) - x\|_p - L_{s(x^*)}$ , where  $s$  is a squashing function that keeps  $x^*$  within the input

domain, e.g., *sigmoid* for images in the domain  $[0, 1]$ . Many advance attacking algorithms have been proposed based on the optimization formulation. [18] shows that, instead of applying different noise to each image, it is possible to apply the same noise, i.e., a universal perturbation, to different images, such that the resulting images still trick the target model in most cases. The one-pixel attack is also shown to be possible [24].

3) *Model Attack*: Similar to the optimization attack, this class also formulates the adversarial attack as an optimization problem. The difference is that, instead of performing the optimization directly, this class trains a separate model to map the input to noise or adversarial samples. Adversarial transformation network (ATN) [1] trains a separate model  $g$  that minimizes  $\beta \|x^* - x\|_p + \|f(x^*) - f(x)\|_{p'}$ , where  $g(x) = x^*$ . [27] proposes to first create a mapping between the input space and a random noise space, and then search in the noise space for potential adversarial samples which are mapped back to the input space. To create the mapping between input and noise space, the authors propose an autoencoder structure which consists of 1) an encoder  $G$ , a generator network that maps the random noise  $z$  to the input  $x$ ,  $G(z) = x$ , and 2) a decoder  $I$  (referred to as *inverter*), another generator network that maps the input to the random noise,  $I(x) = z$ . Generative Adversarial Network (GAN) [5] is used for both generator networks. The whole network is trained end-to-end by minimizing the loss  $\mathbb{E}_x \|G(I(z)) - x\|_p + \lambda \mathbb{E}_z \|I(G(x)) - z\|_p$ .

## B. Adversarial Text Method

Almost all the work in the previous section focus on image models. As we have discussed, the main problem to generate adversarial texts are the discrete input space and the lack of quality measurement. The aforementioned model attack [27] is a viable workaround for the first problem since the noise space is smooth. However, the disadvantage with their method is that they do not have an explicit control of the quality of the generated adversarial samples. As we have seen in [27], the generated adversarial images on complex dataset usually have large visual changes.

Most work [11, 23, 6] on attacking text models follow a similar strategy, 1) first identify the features (characters, words, sentences, etc.) that have the most influence on the prediction, and then 2) follow different strategies to perturb these features according to *manually* constructed perturbation candidates. This strategy is similar to JSMA, in which the intensity of the pixel with the highest saliency score is increased or decreased. The Jacobian value  $\nabla f$  or the loss gradient  $\nabla L$  are usually employed to construct a measurement for the feature importance, e.g.,  $\nabla L$  is used in [11] to select important characters and phrases to perturb. The perturbation candidates usually include typos, synonyms, antonyms, frequent words in each category, and other task-dependent features. For example, typos, synonyms, and important adverbs and adjectives are used as candidates for insertion and replacement in [23]. The strategies to apply the perturbation generally include *insertion*, *deletion*, and *replacement*.

TABLE I: Model accuracy under different parameter settings.  $\epsilon$  is the noise scaling factor. We report two accuracy measurements per parameter setting in the format  $acc_1/acc_2$ , where  $acc_1$  is the model accuracy on adversarial embeddings before nearest neighbor search,  $acc_2$  the accuracy on adversarial embeddings that are reconstructed by nearest neighbor search. In other words,  $acc_2$  is the model accuracy on generated adversarial texts.

Method	Dataset	$acc_1/acc_2$				
		$\epsilon$	0.40	0.35	0.30	0.25
FGSM	IMDB		0.1213 / 0.1334	0.1213 / 0.1990	0.1213 / 0.4074	0.1213 / 0.6770
	Reuters-2		0.0146 / 0.6495	0.0146 / 0.7928	0.0146 / 0.9110	0.0146 / 0.9680
	Reuters-5		0.1128 / 0.5880	0.1128 / 0.7162	0.1128 / 0.7949	0.1128 / 0.8462
FGVM		$\epsilon$	15	30	50	100
	IMDB		0.6888 / 0.8538	0.6549 / 0.8354	0.6277 / 0.8207	0.5925 / 0.7964
	Reuters-2		0.7747 / 0.7990	0.7337 / 0.7538	0.6975 / 0.7156	0.6349 / 0.6523
	Reuters-5		0.5915 / 0.7983	0.5368 / 0.6872	0.4786 / 0.6085	0.4000 / 0.5111
DeepFool		$\epsilon$	20	30	40	50
	IMDB		0.5569 / 0.8298	0.5508 / 0.7225	0.5472 / 0.6678	0.5453 / 0.6416
	Reuters-2		0.4416 / 0.6766	0.4416 / 0.5236	0.4416 / 0.4910	0.4416 / 0.4715
	Reuters-5		0.1163 / 0.4034	0.1162 / 0.2222	0.1162 / 0.1641	0.1162 / 0.1402

A slightly different strategy is used in [6]. The authors add to the samples *manually* constructed legit distracting sentences, which introduce fake information that does not contradict with the samples. This strategy, despite being effective, is not scalable.

In [11], the authors attempt applying FGM directly on character-level CNN [26]. Although the labels of the text pieces are altered, the texts are changed to basically random stream of characters.

### III. OUR METHOD

As we have discussed, all the previous work on generating adversarial texts rely on *manually* selected and *task-dependent* features, which is not practical. In this section, we propose a general framework that generates high-quality adversarial texts without human intervention.

#### A. Discrete Input Space

The first problem we need to resolve is how we can employ small noise to perturb the input. The general idea is simple. Instead of working on the raw input texts, we first embed these texts to vector space and search for adversarials in the embedding space via gradient methods, and then reconstruct the adversarial sentences via nearest neighbor search. Searching for adversarials in the embedding space is similar to searching for adversarial images. To make sure that the generated adversarial embeddings are meaningful, i.e., corresponding to actual tokens so that we can generate sentences from them, we use nearest neighbor search to round the perturbed vectors to nearest meaningful word vectors. The sentence reconstructing process can be seen as a strong *denoising* process. With appropriate noise scale, we would expect most of the words remain unchanged. This framework builds upon the following observations.

- 1) The input features (pixels, words) that are relatively more important for the final predictions will receive more noise, while others relatively less noise. This is actually the core

property of the adversarial image attacking methods. For example, in Figure 1, usually a subset of the features are perturbed.

- 2) The embedded word vectors preserve the subtle semantic relationships among words [15, 14]. For example,  $\text{vec}(\text{"clothing"})$  is closer to  $\text{vec}(\text{"shirt"})$  as  $\text{vec}(\text{"dish"})$  to  $\text{vec}(\text{"bowl"})$ , while  $\text{vec}(\text{"clothing"})$  is far away, in the sense of  $p$ -norm, from  $\text{vec}(\text{"dish"})$  since they are not semantically related [13]. This property assures that it is more likely to replace the victim words with a semantically related one rather than a random one.

Most of the attacking algorithms that apply to image models are applicable in our framework. In this work, however, we focus on gradient methods since they are usually faster.

#### B. Word Mover's Distance (WMD)

The second problem we need to resolve is the choice of quality metric for generated adversarial texts, so that we have a scalable way to measure the effectiveness of our framework. We employ the Word Mover's Distance (WMD) [9] as the metric. WMD measures the dissimilarity between two text documents as the minimum amount of distance that the embedded words of one document need to *travel* to reach the embedded words of another document. WMD can be considered as a special case of Earth Mover's Distance (EMD) [22]. Intuitively, it quantifies the semantic similarity between two text bodies. In this work, WMD is closely related to the ratio of number of words changed to the sentence length. However, we plan to extend our framework with paraphrasing and insertion/deletion, where the sentence length may change. In that case, WMD is more flexible and accurate.

### IV. EXPERIMENT

We evaluate our framework on three text classification problems. Section IV-A details on the data preprocessing. The adversarial attacking algorithms which we use are (FGM) [4]

Clean Text	Label	WMD ( $n/L$ )	Adversarial Text
Quick summary of the book : [...] The book was n't bad , but was soooooo cliché < br / > < br / > Now about the movie [...]	0→1	0.0317 (0.0050)	Quick summary of the book : [...] The book was n't bad , but was soooooo TahitiNut < br / > < br / > Now about the movie [...]
zulchzulu < SM > TO OFFER SPECIAL DIVIDEND South-mark Corp said it will issue its shareholders a special dividend right [...]	1→0	0.0817 (0.0125)	zulchzulu < SM > TO OFFER OFFERS SHARES South-mark Corp said it will issue its shareh olders a special dividend right [...]
U . K . MONEY MARKET GIVEN FURTHER 68 MLN STG HELP The Bank of England said it provided the market with a further [...]	3→2	0.0556 (0.0077)	U . K . MONEY MARKET GIVEN FURTHER 68 ARL STG HELP The Bank of England said it provided the market with a further [...]

Fig. 2: Adversarial texts generated via DeepFool. Refer to Section IV-D for notations and discussions.

and DeepFool [17]. We tried JSMA, however, due to the mechanism of JSMA, it is not directly applicable in our framework. We report in Section IV-D the original model accuracy, accuracy on adversarial embeddings, and accuracy on reconstructed adversarial texts in our experiment. Only a few examples of generated adversarial texts are shown in this paper due to the space constraint. The complete sets of adversarial texts under different parameter settings and the code to reproduce the experiment are available on our website<sup>1</sup>.

Computation-wise, the bottleneck in our framework is the nearest neighbor search. Word vector spaces, such as GloVe [20], usually have millions or billions of tokens embedded in very high dimensions. The vanilla nearest neighbor search is almost impractical. Instead, we employ the an approximate nearest neighbor (ANN) technique in our experiment. The ANN implementation which we use in our experiment is Approximate Nearest Neighbors Oh Yeah (annoy)<sup>2</sup>, which is well integrated into gensim [21] package.

#### A. Dataset

We use three text datasets in our experiments. The datasets are summarized in Table II. The last column shows our target model accuracy on clean test data.

TABLE II: Dataset Summary

Dataset	Labels	Training	Testing	Max Length	Accuracy
IMDB	2	25000	25000	400	0.8787
Reuters-2	2	3300	1438	160	0.9854
Reuters-5	5	1735	585	350	0.8701

1) *IMDB Movie Reviews*: This is a dataset for binary sentiment classification [12]. It contains a set of 25,000 highly polar (positive or negative) movie reviews for training, and 25,000 for testing. No special preprocessing is used for this dataset except that we truncate/pad all the sentences to a fixed maximum length, 400. This max length is chosen empirically.

2) *Reuters*: This is a dataset of 11,228 newswires from Reuters, labeled over 90 topics. We load this dataset through the NLTK [2] package. The raw Reuters dataset is highly

unbalanced. Some categories contain over a thousand samples, while others may contain only a few. The problem with such highly unbalanced data is that the texts that belong to underpopulated categories are almost always get classified incorrectly. Even though our model may still achieve high accuracy with 90 labels, it would be meaningless to include these underpopulated categories in the experiment since we are mainly interested in perturbation of those samples that are already being classified correctly. Keras<sup>3</sup> uses 46 categories out of 90. However, the 46 categories are still highly unbalanced. In our experiment, we preprocess Reuters and extract two datasets from it.

- 1) *Reuters-2* It contains two most populous categories, i.e., *acq* and *earn*. The *acq* category contains 1650 training samples and 719 test samples. Over 71% sentences in the *acq* category have less than 160 tokens. The *earn* category contains 2877 training samples and 1087 test samples. Over 83% sentences in *earn* category have less than 160 tokens. In order to balance the two categories, for *earn*, we use 1650 samples out of 2877 for training, and 719 for testing. The maximum sentence length of this binary classification dataset is set to 160.
- 2) *Reuters-5* It contains five categories, i.e., *crude*, *grain*, *interest*, *money-fx* and *trade*. Similar to Reuters-2, we balance the five categories by using 347 examples (the size of *interest* categories) for each category during training, and 117 each for testing. The maximum sentence length is set to 350.

#### B. Embedding

Our framework relies heavily on the *size* and *quality* of the embedding space. More semantic alternatives would be helpful to improve the quality of generated adversarial texts. As a result, we use the GloVe [20] pre-trained embedding in our experiment. Specifically, we use the largest GloVe embedding, *glove.840B.300d*, which embeds 840 billion tokens (approximately 2.2 million cased vocabularies) into a vector space of 300 dimensions. The value range of the word vectors are roughly  $(-5.161, 5.0408)$ .

<sup>1</sup><http://gongzhitaao.org/adversarial-text>

<sup>2</sup><https://github.com/spotify/annoy>

<sup>3</sup><https://keras.io/>



Clean Text	Label	WMD ( $n/L$ )	Adversarial Text
One of those TV films you saw in the seventies that scared the hell out of you when you were a kid but still gives you an eerie feeling . No great actors or expensive production but everytime that phone rings ... .. (IMDB)	0→1	0.3979 (0.0930)	One of those TV films you saw in the seventies that scared the hell out of you when you were a kid but not gives you considered unnerving feeling . No great actors and/or expensive production but everytime that phone rings ... ..
CIS & It ; CISIF . O > AGREES TO SECOND EXTENSION CIS Technologies Inc said that it and the Swiss Reinsurance Co of Zurich , Switzerland agreed to a second extension of two dates for the final part of their share purchase agreement . It said the June one election date has been extended to June 15 and the June 30 closing date has been changed to July 31 . (REUTERS-2)	0→1	0.9414 (0.1972)	CIS & It - Un-Idle . O > RESULT TO SECOND EXTENSION CIS Technologies Inc saying that it and the Swiss Reinsurance Co of Zurich run Switzerland stated to single second extension of two dates cost the final part the their share purchase agreement . When said the July one election date has been extended to June 15 both the June 30 closing date has been changed to July 29 .
FED EXPECTED TO ADD RESERVES The Federal Reserve will enter the government securities market to supply reserves via either a large round of customer repurchase agreements or by overnight or possibly four - day system repurchases , economists said . [...] tax payments swell Treasury balances at the Fed . Fed funds hovered at 6 - 3 / 4 pct after averaging 6 . 80 pct yesterday . (REUTERS-5)	2→3	0.6628 (0.1146)	FED EXPECTED TO ADD RESERVES The Federal Reserve 'll enter the government securities market able supply reserves via either a large round of customer repurchase agreements or directed overnight or possibly four - day system repurchases books economists said . [...] tax payments surf Treasury balances at the futures well Fed funds hovered at 6 - 3 Category 4 Q2 after averaged 6 ..... 80 pct yesterday .

Fig. 3: Adversarial texts generated via FGSM. Refer to Section IV-D for notations and discussions.

### C. Model

In this work, we focus on feedforward architectures. Specifically, we use CNN model for the classification tasks. The model structure is summarized in Figure 4.

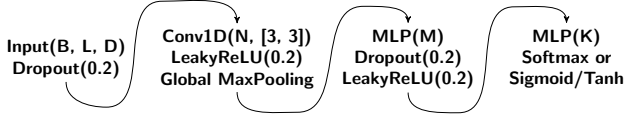


Fig. 4: CNN model for text classification.

Where  $B$  denotes batch size,  $L$  the maximum sentence length,  $D$  the word vector space dimension. In our experiment, we have  $B = 128$ , and  $D = 300$  since we are using the pre-trained embedding glove.840B.300d.

Note that for models trained for binary classification tasks, DeepFool assumes the output in the range  $[-1, 1]$ , instead of  $[0, 1]$ . Thus we have two slightly different models for each of the binary classification task (IMDB and Reuters-2), one with sigmoid output, and the other with tanh. The model with tanh output is trained with Adam [7] by minimizing the mean squared error (MSE), while all the other models are trained with Adam by minimizing the cross-entropy loss. Despite the small difference in architecture, sigmoid- and tanh-models on the same task have almost identical accuracy. As a result, in Table II, we report only one result for IMDB and Reuters-2.

All our models have  $N = 256$  and  $M = 512$ , except for the one with tanh output on the IMDB classification task, in which we have  $N = 128$  and  $M = 256$ . The reason that we change to a smaller model is that the larger one always gets

stuck during the training. We are not yet clear what causes this problem and why a smaller model helps.

### D. Results

The model accuracy on adversarial embeddings before and after the nearest neighbor search under different parameter settings are summarized in Table I.

In the adversarial text examples, to aid reading, we omit the parts that are not changed, denoted by [...] in the texts. The "(IMDB)" at the end of each clean text piece denotes the dataset that this piece of text belongs to. In addition to Word Mover's Distance (WMD), we also report the change rate,  $\frac{n}{L}$ , where  $n$  is the number of changed words,  $L$  the sentence length. The corresponding changed words are highlighted in the figures.

1) *Fast Gradient Method*: We first evaluate two versions of FGM, i.e., FGSM and FGVM. Their example results are shown in Figure 3 and Figure 5, respectively. For FGVM, it was proposed in [16] to use  $\frac{\nabla L}{\|\nabla L\|_2}$  to FGVM usually needs much larger noise scaling factor since most gradients are close to zero.

2) *DeepFool*: Adversarial examples are shown in Figure 2. We experiment with different overshoot values (also denoted as  $\epsilon$  in the table). Usually, for images, we tend to use very small overshoot values, e.g., 1.02, which creates just enough noise to cross the decision boundary. However, in our framework, the reconstructing process is a very strong denoising process, where much of the subtle noise will be smoothed. To compensate for this, we experiment with very large overshoot values. In practice, this works very well. As we can see, labels are altered by replacing just one word in many cases.

Clean Text	Label	WMD (n/L)	Adversarial Text
I was n't really interested in seeing Step Up , but my friend just kept bugging and bugging [...] where are you when we need you ? ! < br / > < br / > 4/10 (IMDB)	1→0	0.0089 (0.0027)	I was n't really interested in seeing Step Up , but my friend just kept bugging and bugging [...] where are you when we need you ? ! < br / > < br / > 7/10
SAN MIGUEL DEAL HIT BY MORE LAWSUITS A bid by San Miguel Corp ( SMC ) & It ; SANM . MN > to buy back 38 . 1 mln sequestered shares from United Coconut Planters Bank ( UCPB ) has been hit by two new lawsuits , sources in the Philippine food and brewery company said . A Manila court yesterday issued an injunction barring UCPB from selling the shares , which represent 31 pct of SMC ' s outstanding capital stock of 121 mln shares , until hearings on April 21 on a petition filed by Eduardo Cojuangco , a former chairman of both SMC and UCPB . Cojuangco said the Coconut Industry Investment Fund ( CIIF ) and 1 . 4 mln farmers were the rightful owners of the shares . Cojuangco said the shares were held in trust by UCPB and represented a blue chip investment . His petition said UCPB ' s plans to sell (REUTERS-2)	0→1	0.5851 (0.0750)	SAN MIGUEL DEAL YEAR RESULT MORE LOSSES A bid by San Miguel Corp ( SMC ) & It ; SANM . MN > to get back 38 . 1 mln sequestered earnings from United Coconut Compost Bank ( UCPB ) has been hit by two new lawsuits , sources in the Philippine food and brewery company said . A Manila court yesterday issued an injunction barring UME from selling the shares , which represent 31 pct of SMC ' s outstanding capital earnings of 121 mln shares , until hearings on April 21 on a petition filed by Eduardo Cojuangco , a former chairman of both SMC and Bolivariano . Cojuangco said the Coconut Industry Investment Fund [ Odyss ) and 1 . 4 mln farmers were the rightful owners of the shares . Cojuangco said the shares were held in trust by Hyvinkaa and represented a blue chip investment . His petition said UCPB ' s plans to sell
CHINA DAILY SAYS VERMIN EAT 7 - 12 PCT GRAIN STOCKS A survey of 19 provinces and seven cities showed vermin consume between seven and 12 pct of China ' s grain stocks , the China Daily said [...] (REUTERS-5)	1→3	0.1249 (0.0153)	CHINA DAILY SAYS VERMIN EAT 7 - 12 PCT CARRYING STOCKS A survey of 19 provinces and seven cities showed vermin consume between seven and 12 pct of China ' s gas stocks , the China Daily said [...]

Fig. 5: Adversarial texts generated via FGVM. Refer to Section IV-D for notations and discussions.

### E. Discussion

In contrary to the experiment in [11], our framework generates much better adversarial texts with gradient methods. One main reason is that the embedding space preserves semantic relations among tokens.

Based on the generated text samples, DeepFool generates the adversarial texts with the highest quality. Our experiment confirms that the DeepFool's strategy to search for the optimal direction is still effective in text models. On the other hand, the strong denoising process will help to smooth unimportant noise. FGVM is slightly better than FGSM, which is quite similar to what we saw in Figure 1. By using  $\text{sign } \nabla L$ , FGSM applies the same amount of noise to every feature it finds to be important, which ignores the fact that some features are more important than others. Since FGVM does not follow the optimal direction as DeepFool does, it usually needs larger perturbation. In other words, compared to DeepFool, FGVM may change more words in practice.

### V. CONCLUSION

In this work, we proposed a framework to adapt image attacking methods to generate high-quality adversarial texts in an end-to-end fashion, without relying on any manually selected features. In this framework, instead of constructing adversarials directly in the raw text space, we first search for adversarial embeddings in the embedding space, and then reconstruct the adversarial texts via nearest neighbor search. We demonstrate the effectiveness of our method on three texts benchmark problems. In all experiments, our framework can successfully generate adversarial samples with only a few

words changed. In addition, we also empirically demonstrate Word Mover's Distance (WMD) as a valid quality measurement for adversarial texts. In the future, we plan to extend our work in the following directions.

- 1) WMD is demonstrated to be a viable quality metric for the generated adversarial texts. We can employ the optimization and model attacking methods by minimizing the WMD.
- 2) We use a general embedding space in our experiments. A smaller embedding that is trained on the specific task may help to speed up the computation needed to reconstruct the texts.

### REFERENCES

- [1] Shumeet Baluja and Ian Fischer. "Adversarial Transformation Networks: Learning To Generate Adversarial Examples". In: *CoRR* abs/1703.09387 (2017). URL: <http://arxiv.org/abs/1703.09387>.
- [2] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc.", 2009.
- [3] Nicholas Carlini and David Wagner. "Towards Evaluating the Robustness of Neural Networks". In: *CoRR* abs/1608.04644 (2016). URL: <http://arxiv.org/abs/1608.04644>.
- [4] I. J. Goodfellow, J. Shlens, and C. Szegedy. "Explaining and Harnessing Adversarial Examples". In: *ArXiv e-prints* (Dec. 2014). arXiv: 1412.6572 [stat.ML].

- [5] I. J. Goodfellow et al. “Generative Adversarial Networks”. In: *ArXiv e-prints* (June 2014). arXiv: [1406.2661 \[stat.ML\]](#).
- [6] Robin Jia and Percy Liang. “Adversarial Examples for Evaluating Reading Comprehension Systems”. In: *arXiv preprint arXiv:1707.07328* (2017).
- [7] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *CoRR* abs/1412.6980 (2014). URL: <http://arxiv.org/abs/1412.6980>.
- [8] A. Kurakin, I. Goodfellow, and S. Bengio. “Adversarial Examples in the Physical world”. In: *ArXiv e-prints* (July 2016). arXiv: [1607.02533 \[cs.CV\]](#).
- [9] Matt Kusner et al. “From word embeddings to document distances”. In: *International Conference on Machine Learning*. 2015, pp. 957–966.
- [10] Baitao Li, E. Chang, and Ching-Tung Wu. “DPF - a perceptual distance function for image retrieval”. In: *Proceedings. International Conference on Image Processing*. Vol. 2. 2002, II-597-II-600 vol.2. DOI: [10.1109/ICIP.2002.1040021](#). URL: <https://doi.org/10.1109/ICIP.2002.1040021>.
- [11] Bin Liang et al. “Deep Text Classification Can Be Fooled”. In: *arXiv preprint arXiv:1704.08006* (2017).
- [12] Andrew L. Maas et al. “Learning Word Vectors for Sentiment Analysis”. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Portland, Oregon, USA: Association for Computational Linguistics, June 2011, pp. 142–150. URL: <http://www.aclweb.org/anthology/P11-1015>.
- [13] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. “Linguistic regularities in continuous space word representations.” In: *hlt-Naacl*. Vol. 13. 2013, pp. 746–751.
- [14] Tomas Mikolov et al. “Distributed Representations of Words and Phrases and Their Compositionality”. In: *CoRR* abs/1310.4546 (2013). URL: <http://arxiv.org/abs/1310.4546>.
- [15] Tomas Mikolov et al. “Efficient Estimation of Word Representations in Vector Space”. In: *CoRR* abs/1301.3781 (2013). URL: <http://arxiv.org/abs/1301.3781>.
- [16] Takeru Miyato et al. “Distributional Smoothing With Virtual Adversarial Training”. In: *stat* 1050 (2015), p. 25.
- [17] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. “Deepfool: a Simple and Accurate Method To Fool Deep Neural Networks”. In: *CoRR* abs/1511.04599 (2015). arXiv: [1511.04599](#). URL: <http://arxiv.org/abs/1511.04599>.
- [18] Seyed-Mohsen Moosavi-Dezfooli et al. “Universal Adversarial Perturbations”. In: *arXiv preprint arXiv:1610.08401* (2016).
- [19] Nicolas Papernot et al. “The Limitations of Deep Learning in Adversarial Settings”. In: *CoRR* abs/1511.07528 (2015). URL: <http://arxiv.org/abs/1511.07528>.
- [20] Jeffrey Pennington, Richard Socher, and Christopher Manning. “Glove: Global vectors for word representation”. In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014, pp. 1532–1543.
- [21] Radim Řehůřek and Petr Sojka. “Software Framework for Topic Modelling with Large Corpora”. English. In: *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*. <http://is.muni.cz/publication/884893/en>. Valletta, Malta: ELRA, May 2010, pp. 45–50.
- [22] Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. “The Earth Mover’s Distance As a Metric for Image Retrieval”. In: *International journal of computer vision* 40.2 (2000), pp. 99–121.
- [23] Suranjana Samanta and Sameep Mehta. “Towards Crafting Text Adversarial Samples”. In: *arXiv preprint arXiv:1707.02812* (2017).
- [24] J. Su, D. Vasconcellos Vargas, and S. Kouichi. “One Pixel Attack for Fooling Deep Neural networks”. In: *ArXiv e-prints* (Oct. 2017). arXiv: [1710.08864 \[cs.LG\]](#).
- [25] Christian Szegedy et al. “Intriguing Properties of Neural Networks”. In: *CoRR* abs/1312.6199 (2013). URL: <http://arxiv.org/abs/1312.6199>.
- [26] Xiang Zhang, Junbo Zhao, and Yann LeCun. “Character-level Convolutional Networks for Text Classification”. In: *Advances in Neural Information Processing Systems* 28. Ed. by C. Cortes et al. Curran Associates, Inc., 2015, pp. 649–657. URL: <http://papers.nips.cc/paper/5782-character-level-convolutional-networks-for-text-classification.pdf>.
- [27] Z. Zhao, D. Dua, and S. Singh. “Generating Natural Adversarial Examples”. In: *ArXiv e-prints* (Oct. 2017). arXiv: [1710.11342 \[cs.LG\]](#).