# Discrete Adversarial Attacks and Submodular Optimization with Applications to Text Classification

Qi Lei<sup>\*</sup>, Lingfei Wu<sup>†\*</sup>, Pin-Yu Chen<sup>†</sup>, Alexandros G. Dimakis<sup>\*</sup>, Inderjit S. Dhillon<sup>\*‡</sup>, and Michael Witbrock<sup>†</sup>

\* UT Austin † IBM Research ‡ Amazon {leiqi@ices, dimakis@austin, dhillon@cs}.utexas.edu {wuli@us., pin-yu.chen@, witbrock@us.}ibm.com

April 8, 2019

#### Abstract

Adversarial examples are carefully constructed modifications to an input that completely change the output of a classifier but are imperceptible to humans. Despite these successful attacks for continuous data (such as image and audio samples), generating adversarial examples for discrete structures such as text has proven significantly more challenging. In this paper we formulate the attacks with discrete input on a set function as an optimization task. We prove that this set function is submodular for some popular neural network text classifiers under simplifying assumption. This finding guarantees a 1-1/e approximation factor for attacks that use the greedy algorithm. Meanwhile, we show how to use the gradient of the attacked classifier to guide the greedy search. Empirical studies with our proposed optimization scheme show significantly improved attack ability and efficiency, on three different text classification tasks over various baselines. We also use a joint sentence and word paraphrasing technique to maintain the original semantics and syntax of the text. This is validated by a human subject evaluation in subjective metrics on the quality and semantic coherence of our generated adversarial text.

## 1 Introduction

Adversarial examples are carefully constructed modifications to an input that completely change the output of a classifier but are imperceptible to humans. Spam filtering and the carefully-crafted emails designed to fool these early classifiers are the first examples of adversarial machine learning going back to 2004 [1, 2]; see also the comprehensive survey by Biggio et al. [3]. Szegedy et al. [4] discovered that deep neural network image classifiers can be fooled with tiny pixel perturbations; exploration of this failure of robustness has received significant attention recently, see e.g. [5, 6, 7, 8, 9, 10, 11, 12]. Adversarial training [5, 13] seems to be the state of the art in defense against adversarial attacks, but creating robust classifiers remains challenging, especially for large image classifiers, see e.g. Athalye at al. [14].

Despite these successful attacks for continuous data (such as image and audio samples), generating adversarial examples for discrete structures such as text and code has proven significantly more challenging in two aspects:

One challenge is how to develop a fast yet (provably) effective attacking scheme. Gradient-based adversarial attacks for continuous data no longer directly apply to discrete structures. Although some variants are proposed when the model is differentiable to the embedding layer [15, 16, 17, 18], this line of methods achieve efficiency but suffer from poor success rate.

<sup>\*</sup>Both authors contributed equally to this work

Meanwhile, another natural idea is to find feasible replacement for individual features like words or characters. However, since the space of possible combinations of substitutions grows exponentially with the length of input data, finding the optimal combination of substitutions is intractable. Recent heuristic attacks on NLP classifiers operate by greedy character-level or word-level replacements [17, 19, 20]. However, greedy methods are usually slow, and it's theoretically not understood when they achieve good performance.

The other issue is how to maintain the original functionality of the input. Specifically for text, it remains challenging to preserve semantic and syntactic properties of the original input from the point of view of a human. Existing methods either require to change too many features, or change the original meaning. For instance, [19] alters up to 50% of words in each input document to achieve a 30% success rate. [18] attacks the document by replacing with completely different words. [21] inserts irrelevant sentences to the original text. Such changes can be easily detected by humans.

In this paper we argue that these limitations can be be resolved with the framework we propose. We highlight our main contributions as follows:

We propose a general framework for discrete attacks. We apply our framework to designing adversarial attacks for text classifiers but our techniques can be applied more broadly. For instance, the attacks include but are not limited to malware detection, spam filtering, or even discrete attacks defined on continuous data, e.g., segmentation of an image.

We formulate the attacks with discrete input on a set function as an optimization task. This problem, however, is provably NP-hard even for convex classifiers. We unify existing gradient-based as well as greedy methods using a general combinatorial optimization via further assumptions. We note that gradient methods solve a relaxed problem in polynomial time; while greedy algorithm for creating attacks has a provable 1-1/e approximation factor assuming the set function is submodular. We theoretically show that for two natural classes of neural network text classifiers, the set functions defined by the attacks are submodular. We specifically analyze two classes of classifiers: The first is word-level CNN without dropout or softmax layers. The second is a recurrent neural network (RNN) with one-dimensional hidden units and arbitrary time steps.

Nevertheless, greedy methods can be very time consuming when the space of attacks is large. We show how to use the gradient of the attacked classifier to guide the combinatorial search. Our proposed gradient-guided greedy method is inspired by the greedy coordinate descent Gauss-Southwell rule from continuous optimization theory. The key idea is that we use the magnitude of the gradient to decide which features to attack in a greedy fashion.

We extensively validate the proposed attacks empirically. With the proposed optimization scheme, we show significantly improved attack performance over most recent baselines. Meanwhile we propose a joint sentence and word paraphrasing technique to simultaneously ensure retention of the semantics and syntax of the text.

## 2 Related Work

Broadly speaking, adversarial examples refer to minimally modified natural examples that are spurious but perceptually similar and that lead to inconsistent decision making between humans and machine learning models. An example is automatically classifying an adversarial stop sign image (according to humans) as a speed limit sign. For continuous data such as images or audio, generating adversarial examples is often accomplished by crafting additive perturbations of natural examples, resulting in visually imperceptible or inaudible noise that misleads a target machine learning model. These small yet effective perturbations are difficult for humans to detect, but will cause an apparently well-trained machine learning model to misbehave; in particular, neural networks have been shown to be susceptible to such attacks [4], giving rise to substantial concern about safety-critical and security-centric machine learning applications.

For classifiers with discrete input structures, a simple approach for generating adversarial examples is to replace each feature with similar alternatives. Such features for text classification tasks are usually individual words or characters. Such attacks can be achieved using continuous word embeddings or with respect to some designed score function; this approach has been applied to attack NLP classifiers [15, 16, 22, 23, 24, 25, 18, 19, 26, 27, 20] and sequence-to-sequence models [17, 28, 29, 30]. The work in [31] considers semantically

Task: Sentiment Analysis. Classifier: LSTM. Original: 100% Positive. ADV label: 100% Negative.

I suppose I should write a review here since my little Noodle-oo is currently serving as their spokes dog in the photos. We both love Scooby Do's. They treat my little butt-faced dog like a prince and are receptive to correcting anything about the cut that I perceive as being weird. Like that funny poofy pompadour. Mohawk it out, yo. Done. In like five seconds my little man was looking fabulous and bad ass. Not something easily accomplished with a prancing pup that literally chases butterflies through tall grasses. (He ended up looking like a little lamb as the cut grew out too. So adorable.) The shampoo they use here is also amazing. Noodles usually smells like tacos (a combination of beef stank and corn chips) but after getting back from the Do's, he smelled like Christmas morning! Sugar and spice and everything nice instead of frogs and snails and puppy dog tails. He's got some gender identity issues to deal with. The pricing is also cheaper than some of the big name conglomerates out there The price is cheaper than some of the big names below. I'm talking to you Petsmart! I've taken my other pup to Smelly Dog before, but unless I need dog sitting play time after the cut, I'll go with Scooby's. They genuinely seem to like my little Noodle monster.

Task: Fake-News Detection. Classifier: LSTM. Original label: 100% Fake. ADV label: 77% Real

Man Guy punctuates high-speed chase with stop at In-N-Out Burger drive-thru Print [Ed.—Well, that's Okay, that 's a new one.] A One man is in custody after leading police on a bizarre chase into the east Valley on Wednesday night. Phoenix police began has begun following the suspect in Phoenix and the pursuit continued into the east Valley, but it took a bizarre turn when the suspect stopped at an In-N-Out Burger restaurant's drive-through near Priest and Ray Roads in Chandler. The suspect appeared to order food, but then drove away and got out of his pickup truck near Rock Wren Way and Ray Road. He then ran into a backyard ran to the backyard and tried to get into a house through the back door get in the home.

Task: Spam Filtering. Classifier: WCNN. Original label: 100% None-spam. ADV label: 100% Spam

- >> Hi All,
- >> I'm new to R from a C and Octave/Matlab background. I am trying to >> construct I 'm trying to build some classes in R to which I want to attach pieces of data.
- >> First, is attr(obj, 'member name') >> this? >> No, it isn't. You seem to be trying to deduce new-style classes from a >> representation used before R 2.4, (actually, still used) >> but in any case it would not be » sensible. Please-consult Contact John M. Chambers. Programming with Data. >> Springer, New York, 1998, and/or William N. Venables and Brian D. Ripley. >> S Programming. Springer, New York, 2000, or for a shorter online resource: >> http://www.stat.auckland.ac.nz/S-Workshop/Gentleman/Methods.pdf >> Unfortunately, all of those references are at least 4 years out of > date when it comes to S4 methods. Is there any comprehensive > reference of the current implementation of the S4 OO system apart from > the source code? Not that I know of, and it is it's a moving target. (E.g. I asked recently about some anomalies in the S4 bit introduced for 2.4.0 and what the intended semantics are.) I've said before that I believe we can only help solve some of the efficiency issues with S4 if we have a technical manual. It is unfair to pick out S4 here, but the 'R Internals' manual is an attempt to document important implementation ing details (mainly by studying the code), and that has only got most of the way through src/main/\*.c.

Figure 1: Examples of generated adversarial examples. The color red denotes sentence-level paraphrasing, and blue denotes word-level paraphrasing.

equivalent rules for debugging NLP models, but under the same input structure. This is a natural but limited practice to only consider attacks within one input structure, namely word or characters, but no joint attacks, nor the effect incurred from sentences. Unlike prior work, we conduct a joint sentence and word paraphrasing technique. It considers sentence-level factors and allows more degrees of freedom in generating text adversarial examples, by exploring the rich set of semantically similar paraphrased sentences.

Jia and Liang studied adversarial examples in reading comprehension systems by inserting additional sentences [21], which is beyond the concept of this paper since the approach changes the original meanings. Another related line of research, although not cast as adversarial examples, focuses on improving model robustness against out-of-vocabulary terms [32] or obscured embedding space representations [33].

# 3 Preliminary

In this paper, we propose a general framework for generating adversarial examples with discrete input data. A collection of such data and corresponding attacks are presented in Table 1.

To present our mathematical formulation, we start by introducing some notation.

**Input Structure.** Let the input  $\mathbf{x} = [x_1, x_2, \dots, x_n] \in \mathcal{X}^n$  be a list of n features (might be padded). For text environment, the feature space  $\mathcal{X}$  can be the character, word, phrase, or sentence space. For the problem of malware detection,  $\mathbf{x}$  is a concatenation of code pieces.

input data	task
document	text classification
$\operatorname{code}$	malware detection
url address	malicious website check

	 like	ιO	eai	lulicii	111	เมเธ	cale.
search	fancy		have	food	from	the	restaurant cafeteria
space	love	$\times$	get	× brunch	$\times$ at $\times$	that >	cafeteria
Space	adore						eatery
C					K		
transformation					<i>&gt;&gt;</i>		
indexing:	l∈	$\{0, 1,$	$\cdots, k$ –	$\{1\}^n$	l = [0, 2]	0, 1, 1	[2, 0, 0, 3]
-		-		-			

Figure 2: An illustration of the transformation indexing when applying to a text sentence. In this example, the transformation denoted as 1 modifies the original sentence to the new one shown in the red boxes.

**Table 1:** Applications to the framework.

**Remark 1.** For concrete usage, we use  $w \in W$  to denote word space, and  $s \in S$  to denote sentences to distinguish the differences.

**Embedding** V. The embedding layer is a key transition from discrete input data into continuous space, which could then be fed into the classifier. For text domain, we typically use the bag-of-words embedding or word-to-vector embedding.

For a bag-of-words embedding,  $V: \mathcal{X}^n \to \mathbb{R}^D$  represents a document as the statistics of word counts, i.e., the summation of each word's one-hot representation. Meanwhile, word-to-vector embeddings characterize different words as D-dimensional vectors, i.e.,  $V(x) \in \mathbb{R}^D, \forall x \in \mathcal{X}$ . When there's no ambiguity, we also use  $V: \mathcal{X}^n \to \mathbb{R}^{n \times D}$  to denote the concatenation of word vectors of the input document as a list of words.

**Transformation Indexing.** Suppose each feature  $x \in \mathcal{X}$  has (at most) k-1 possible replacements, denoted by  $x^{(i)}, i \in [k-1] (\equiv \{1, 2, \cdots, k-1\})$ . For future use, we also define  $x^{(0)} = x, \forall x \in \mathcal{X}$ . A valid transformation T is the combined replacement of each individual feature  $x_i, i \in [n]$ . Therefore we index T by a vector  $\mathbf{l} \in \{0, 1, \cdots, k-1\}^n$ , and  $l_i$  indicates the index of each replacement i. Namely,  $T_{\mathbf{l}}(\mathbf{x} = [x_1, x_2, \cdots, x_n]) = [x_1^{(l_1)}, x_2^{(l_2)}, \cdots, x_n^{(l_n)}]$ . An example with word replacement in the text classification environment can be found in Figure 2.

Classifier output  $C_y$ . We consider a targeted attack, i.e., we want to maximize the output probability C over a specific target label y.

In this paper, we use a regular lower-case symbol to denote a scalar or a single feature, and use a bold lower-case symbol for a vector or a list of features.

#### 3.1 Problem Setup

In most scenarios, we only allow transformations on at most m features, then the constraint is  $||\mathbf{l}||_0 \leq m$ . Therefore we present the adversarial attack problem formally:

**Problem 1.** For some input data  $\mathbf{x} \in \mathcal{X}^n$  and target label y, we try to find a feasible transformation  $T_{\mathbf{l}^*}$ , where  $\mathbf{l}^* \in \{0, 1, \dots k-1\}^n$  is the index so that:

$$\mathbf{l}^* = \underset{\|\mathbf{l}\|_0 \le m}{\arg\max} C_y \left( V \left( T_{\mathbf{l}}(\mathbf{x}) \right) \right). \tag{1}$$

Or similarly, we want to find the set of features to attack, i.e.,

$$S^* = \operatorname*{arg\,max}_{|S| < m} f(S),\tag{2}$$

where we defined the set function  $f: 2^{[n]} \to \mathbb{R}$ ,  $f(S) = \max_{supp(1) \subset S} C_y(V(T_1(\mathbf{x})))$ .

The set function f(S) represents the classifier output for the target label y if we apply a set of transformations S. We are therefore searching over all possible sets of up to m replacements to maximize the probability of the target label output of a classifier.

Remark 2. In this paper, we focus on replacements via word and sentence paraphrasing for empirical studies. However, our formulation is general enough to represent any set of discrete transformations. Possible transformations include replacement with the nearest neighbor of the gradient direction [18] and word vectors [19], or flipping characters within each word [17]. We will also conduct a thorough experimental comparisons among different choices.

## 4 Theoretical Analysis

First, notice that the original problem is computationally intractable in general:

**Proposition 1.** For a general classifier  $C_y$ , the problem 1 is NP-hard. Specifically, even for some convex  $C_y$ , the problem 1 can be polynomially reduced from subset sum and hence is NP-hard.

Details and all proofs referenced to in this paper can be found in the appendix.

## 4.1 Unifying Related Methodology via Further Assumptions

Fortunately, with further assumptions it becomes possible to solve problem 1, above, in polynomial time. Some existing heuristics are proposed to generate adversarial examples for the text classification problem. Though usually not specifically proposed in the relevant literature, we unify the underlying assumptions for these heuristics to succeed in polynomial time in this section.

One possible assumption is that the original function  $C_y$  is smooth, which could afterwards be approximated by its first-order Taylor expansion:

$$C_y(V(T_1(\mathbf{x}))) = C_y(\mathbf{v}) + \langle \nabla C_y(\mathbf{v}), V(T_1(\mathbf{x})) - \mathbf{v} \rangle + \mathcal{O}(\|V(T_1(\mathbf{x})) - \mathbf{v}\|_2^2)$$

where  $\mathbf{v} = V(\mathbf{x})$ . Therefore, Problem 1 can be relaxed as follows:

**Problem 2.** Given gradient  $\nabla C_y(\mathbf{v})$ , where  $\mathbf{v} = V(\mathbf{x})$ , maximize function  $C_y$  by its first-order Taylor expansion:

$$\mathbf{l}^* = \underset{\|\mathbf{l}\|_0 \le m}{\arg \max} V(T_{\mathbf{l}}(\mathbf{x}))^\top \nabla C_y(\mathbf{v}). \tag{3}$$

Problem 2 is similar to the Frank-Wolfe method [34] in continuous optimization and is easy to solve:

**Proposition 2.** Problem 2 can be solved in polynomial time for both bag-of-words and word to vector embeddings. Specifically,  $f(S) = \arg\max_{supp(1) \subset S} V(T_1(\mathbf{x}))^{\top} \Delta C_y(\mathbf{v})$  can be written as  $\sum_{i \in S} w_i$  for some w irrelevant to S, where  $\mathbf{v} = V(\mathbf{x})$ .

Related methods like [18] are attempts to solve problem 2. They propose to conduct transformations via replacement by synonyms chosen by (3). However, activations like ReLU break the smoothness of the function, and first order Taylor expansion only cares about very local information, while embeddings for word synonyms could be actually not that close to each other. Consequently, this unnatural assumption prevents related gradient-based attacks to achieve good performance.

Besides smoothness, another more natural assumption is that f(S) in the original problem 1 is submodular [35, 36]. Submodular is a property that is defined for set functions, which characterizes the diminishing returns of the function value change as the size of the input set increases.

**Definition 1.** [37] If  $\Omega$  is a finite set, a submodular function is a set function  $f: 2^{\Omega} \to \mathbb{R}$ , where  $2^{\Omega}$  denotes the power set of  $\Omega$ , which satisfies one of the following equivalent conditions.

- 1. For every  $X, Y \subseteq \Omega$  with  $X \subseteq Y$  and every  $x \in \Omega \setminus Y$  we have that  $f(X \cup \{x\}) f(X) \ge f(Y \cup \{x\}) f(Y)$ .
- 2. For every  $S, T \subseteq \Omega$  we have that  $f(S) + f(T) \ge f(S \cup T) + f(S \cap T)$ .
- 3. For every  $X \subseteq \Omega$  and  $x_1, x_2 \in \Omega \setminus X$  we have that  $f(X \cup \{x_1\}) + f(X \cup \{x_2\}) \geq f(X \cup \{x_1, x_2\}) + f(X)$ .

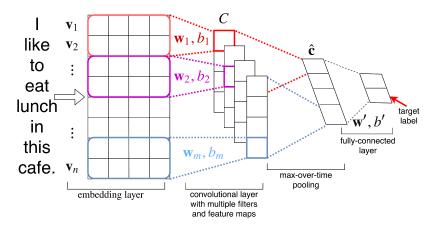


Figure 3: Model architecture of simplified W-CNN for an example sentence.

With the design of f(S) in Problem 1 to be monotone non-decreasing and if we further assume f to be submodular, our task becomes to maximize a monotone submodular function subject to a cardinality constraint [38]. Therefore, greedy method guarantees a good approximation of the optimal value of Problem 1:

Claim 1. In problem 1, f is monotone non-decreasing. Furthermore, if the function f is submodular, greedy methods achieve a (1-1/e)-approximation of the optimal solution in polynomial time.

Both our work and the optimization scheme from [19] propose some variants of greedy methods with the underlying submodular assumption.

The greedy method proposed in [19] selects candidate replacements directly by function value, one word at a time, which we will refer as the objective-guided greedy method. We will propose a more efficient yet comparable effective greedy method that is guided by the gradient magnitude in Section 5.2, and compare with the above two methods in Section 6.3. As an extension from the continuous optimization, our method uses the well-studied Gauss-Southwell rule [39] that is provably better than random selection. In each iteration, we determine and select the most important words by the gradient norm of words' embeddings, and then find the greediest transformation within the search space of the selected words. The advantage is that we are able to conduct multiple replacements in one iteration and thus take into consideration the joint effect of multiple words replacements. We will introduce our method, which we call Gradient-Guided Greedy Word Paraphrasing in Algorithm 3, and will show empirical performance comparison with the (objective-guided) greedy method [19] and the gradient method used in [18] in Section 6.3.

## 4.2 Submodular Neural Networks on the Set of Attacks

To argue that submodular is a natural assumption, we study and summarize the neural networks are submodular on the set of attacks.

In [40], it provides a class of submodular functions used in the deep learning community called deep submodular functions. Nevertheless the deep submodular functions are not necessarily applicable to our set function. We hereby formally prove the following two kinds of neural networks, that are ubiquitously used for text classification, indeed satisfy submodular property on the set of attacks under some conditions.

## 4.2.1 Simplified W-CNN [41]

Denote the stride as s, the number of grams (window size) h, and the word vector of the i-th word in a document as  $\mathbf{v}_i (\equiv V(x_i))$ . Then the output for the convolutional layer is a matrix  $C = [c_{ij}]_{i \in [n/s], j \in [m]}$  from n words and m filters:

$$c_{ij} = \phi(\mathbf{w}_j^{\top} \mathbf{v}_{s(i-1)+1:s(i-1)+h} + b_j), \ i = 1, 2, \dots n/s,$$

where  $\mathbf{w}_j \in \mathbb{R}^{Dh}$  is the *j*-th filter,  $b_j$  is the corresponding bias term and  $\phi$  is the non-linear, and non-decreasing activation such as ReLU, tanh and sigmoid function.  $\mathbf{v}_{i:j}$  denotes the concatenation of word vectors in the window of words *i* through *j*, namely  $[\mathbf{v}_i^{\top}, \mathbf{v}_{i+1}^{\top}, \cdots \mathbf{v}_j^{\top}]^{\top} \in \mathbb{R}^{D(j-i+1)}$ . Each filter  $\mathbf{w}_j$  is applied to individual windows of words to produce a feature map  $\mathbf{c}^j = [c_{1j}, c_{2j}, \cdots c_{n/s,j}]^{\top}$ .

Afterwards, a max-over-time pooling is applied to each feature map to form the penultimate layer  $\hat{\mathbf{c}} = [\hat{c}_1, \hat{c}_2, \dots \hat{c}_m]$ , where  $\hat{c}_i$  is the largest value in  $\mathbf{c}^j$ :

$$\hat{c}_j = \max_i c_{ij}.$$

Compared to the original [41] paper, we only omit the dropout and softmax layer, and instead consider the following WCNN classifier output for a target label:

$$C^{\text{WCNN}}(\mathbf{v}_{1:n}) = \mathbf{w}' \cdot \hat{\mathbf{c}} + b' \tag{4}$$

**Theorem 1.** We consider the simple version of W-CNN classifier described in (4), and suppose there's no overlapping between each window, i.e.,  $s \ge h$ , and  $\mathbf{w}'$  has all non-negative values. If further we only look at transformations that will increase the output, i.e.,  $\mathbf{w}_j^\top V(x_i^{(t)}) \ge \mathbf{w}_j^\top V(x_i), \forall i \in [n], j \in [m], t \in [k-1],$  then  $f^{WCNN}(S) = \max_{supp(1) \in S} C^{WCNN}(V(T_1(x)))$  is submodular.

The proof sketch is as follows. Every coordinate in  $\hat{\mathbf{c}}$  is a combination of max pooling over a modular function and is therefore submodular. And finally sums of submodular functions is still submodular.

Besides word-level CNN, another network that is popular in the NLP community is the recurrent neural network (RNN) or its variants. We will show that under some conditions, RNN satisfies submodular property.

#### 4.2.2 Recurrent Neural Network with One-dimensional Hidden Units

Consider a RNN with T time steps and each hidden layer is a single node. Then for all  $t \leq T$ , given the value of a previous hidden state  $h_{t-1} \in \mathbb{R}$  and an input word vector  $\mathbf{v}_{t-1} \in \mathbb{R}^D$  ( $\mathbf{v}_t \equiv V(x_t)$ ), RNN computes the next hidden state  $h_t$  and output vector  $\mathbf{o}_t \in \mathbb{R}$  as:

$$h_t = \phi(wh_{t-1} + \mathbf{m}^\top \mathbf{v}_{t-1} + b) \tag{5}$$

The classifier output is  $C^{RNN}(\mathbf{v}_{1:T}) = yh_T$ .

**Theorem 2.** For a recurrent neural network with T time steps and one-dimensional hidden nodes described in (5), if w and y are positive, and the activation is a non-decreasing concave function, then  $f^{RNN}(S) = \max_{supp(1) \in S} C^{RNN}(V(T_1(\mathbf{x})))$  is submodular.

This result is quite surprising, since the word vectors influence the network's output on different time steps and are by no means separable. In the proof, we first show that a same amount of change induced on an intermediate layer has a diminishing effect when the network is attacked on more features. Then together with the concavity and non-decreasing property of the network, we are able to finish the proof.

# 5 Adversarial Text Examples via Paraphrasing

In order to conduct adversarial attacks on models with discrete input data like text, one essential challenge is how to select suitable candidate replacements so that the generated text is both semantic meaning preserving and syntactically valid. Another key issue is how to develop an efficient yet effective optimization scheme to find good transformations. To solve the above two issues, we propose our methodology for generating adversarial examples for text.

## **Algorithm 1** Joint Sentence And Word Paraphrasing $(C_y, \mathbf{x}^{(0)}, P, \delta, \lambda_s, \lambda_w, \delta_s, \delta_w, \tau, k)$

- 1: **Input:** Classifier C associated with target label y, input document  $\mathbf{x}^{(0)}$ , language model P trained on the training set, syntactic threshold  $\delta$ , sentence and word paraphrasing ratio  $\lambda_s, \lambda_w$ , termination threshold  $\tau$ , WMD threshold  $\delta_s, \delta_w$ , limit number of paraphrases k.
- 2: Conduct sentence separation  $\mathbf{x}^{(0)} \to [s_1, s_2, \cdots s_l], s_i \in \mathcal{S}, 1 \leq i \leq l$ . (See Remark 1).
- 3: Create sentence neighboring set  $\mathbf{S} = \{S_1, S_2, \dots S_l\}$ , where each  $S_i \subset \mathcal{S}$  satisfies that  $|S_i| \leq k$  and  $WMD(s_i, s) \leq \delta_s, \forall s \in S_i$ .
- 4:  $\mathbf{x}^{(1)} \leftarrow Greedy \ Sentence \ Paraphrasing(C_y, \mathbf{x}^{(0)}, \mathbf{S}, \lambda_s, \tau) \ in \ Alg. \ 2.$
- 5: If  $C_u(V(\mathbf{x})) \ge \tau$  Return  $\mathbf{x}^{(1)}$
- 6: Conduct word separation  $\mathbf{x}^{(1)} \to [w_1, w_2, \cdots w_n], w_i \in \mathcal{W}, 1 \leq i \leq n$ .
- 7: Create word neighboring set  $\mathbf{W} = \{W_1, W_2, \dots W_n\}$ , where each  $W_i \subset \mathcal{W}$  satisfies that  $|W_i| \leq k$  and  $WMD(w_i, w) \leq \delta_w, |P(\mathbf{x}^{(1)}) P(\mathbf{x}'(w))| \leq \delta, \forall w \in W_i$ , where  $\mathbf{x}'(w)$  is text  $\mathbf{x}^{(1)}$  in which  $w_i$  is substituted by w.
- 8:  $\mathbf{x}^{(2)} \leftarrow Gradient \ Guided \ Greedy \ Word \ Paraphrasing(C_y, \mathbf{x}^{(1)}, \mathbf{W}, \lambda_w, \tau) \ in \ Alg. \ 3.$
- 9: Return  $\mathbf{x}^{(2)}$

## Algorithm 2 Greedy Sentence Paraphrasing $(C_y, \mathbf{x}, \mathbf{S}, \lambda_s, \tau)$

```
1: Input: Document x as list of sentences [s_1, s_2, \cdots, s_l], sentence neighboring sets \mathbf{S} = \{S_1, S_2 \cdots S_n\},
    model C_y and parameters \lambda_s, \tau.
   while C_u(V(\mathbf{x})) \leq \tau and number of sentence paraphrased \leq \lambda_s l do
       Create candidate set M = \emptyset
 3:
       for j = 1, 2, \dots, l do
 4:
          for s \in S_i do
 5:
             Substitute s_j by s to get x' and add it to the candidate set M \leftarrow M \cup \{\mathbf{x}'\}.
 6:
 7:
          end for
          \mathbf{x} \leftarrow \arg\max_{\mathbf{x}' \in M} C_y(V(x'))
 8:
       end for
9:
10: end while
```

## 5.1 Joint Sentence and Word Paraphrasing

To coincide with the definition of adversarial examples for text, we first determine appropriate word and sentence paraphrasing methods that maintain the semantic meaning of the original text. Our scheme is to generate an initial set for word and sentence replacements with a well-studied paraphrasing corpus and then filter out discrepant choices based on their semantic and syntactic similarities to the original text. A similar mechanism was also used by [19] to generate word replacement candidates.

#### Paraphrasing Corpus.

For word paraphrasing, we use the Paragram-SL999 [42] of 300 dimensional paragram embeddings to generate neighboring paraphrasing for words. For sentences, we use the pretrained model from Wieting and Gimpel's Para-nmt-50m project [43] to generate sentence paraphrases.

We further specify semantic and syntactic constraints to ensure good quality in adversarial texts:

#### Semantic similarity.

We use the Word Mover Distance (WMD) [44] to measure semantic dissimilarity. For sentence pairs, WMD captures the minimum total semantic distance that the embedded words of one sentence need to "travel" to the embedded words of another sentence. While for words, WMD directly measures the distance between their embeddings.

#### Syntactic similarity.

Alongside the semantic constraint, one should also ensure that the generated sentence is fluent and natural. We make use of a language model as in [19],  $P: \mathcal{X}^n \to [0,1]$  to calculate the probability of the adversarial

## **Algorithm 3** Gradient Guided Greedy Word Paraphrasing $(C_y, \mathbf{x}, \mathbf{W}, \lambda_w, \tau)$

```
1: Input: Document x as a list of words [w_1, w_2, \cdots, w_n], word neighboring sets \mathbf{W} = \{W_1, W_2 \cdots W_n\},
    model C_y and parameters \lambda_w, \tau.
 2: Let N (that we set as 5) be the number of words to replace at most in each iteration
 3: while C_{\nu}(\mathbf{x}) \leq \tau and number of words paraphrased \leq \lambda_{\nu} n do
       Compute score for each word \mathbf{p}: p_i = \|\nabla_i C_{ij}(\mathbf{v})\|_2, where \mathbf{v} = V(\mathbf{x}) and \nabla_i denotes the gradient with
       respect to the embedding of the i-th word in x.
       Get the indices I = \{i_1, i_2, \dots i_N\}: the N largest indices in p.
 5:
 6:
       Create candidate set M = \{x\}
       for j \in I do
 7:
          Let the new candidate set \bar{M} \leftarrow \emptyset
 8:
          for \bar{\mathbf{x}} \in M do
9:
10:
              for w \in W_i do
                 Substitute the j-th word in \bar{\mathbf{x}} by w to get \mathbf{x}' and add it to the candidate set \bar{M} \leftarrow \bar{M} \cup \{\mathbf{x}'\}.
11:
12:
              end for
          end for
13:
          M \leftarrow M \cup \bar{M}
14:
       end for
15:
       \mathbf{x} \leftarrow \operatorname{arg\,max}_{\mathbf{x}' \in M} C_y(\mathbf{x}')
16:
17: end while
```

sentence, and require:

$$|\ln(P(\mathbf{x})) - \ln(P(\mathbf{x}'))| \le \delta,$$

where  $\mathbf{x}'$  is the adversarial sentence paraphrased from  $\mathbf{x}$ .

In Algorithm 1, we present the whole procedure of finding the neighboring sets to conduct our proposal joint sentence and word paraphrasing attack. While with more details, we show how to use the objective value as well as gradient information to guide the search in Algorithm 2 (for sentences) and 3 (for words).

#### 5.2 Gradient-Guided Greedy Method

In Section 3.1 we have demonstrated the difficulty of finding the best transformation from combinatorially many choices. Here we specify our proposal, gradient-guided greedy word paraphrasing, as shown in Algorithm 3. We can see that we first use gradient values to determine the index set of N words  $(w_{i_1}, w_{i_2}, \dots, w_{i_N})$  that we want to replace (steps 4-5). Then in steps 7-15 we create a candidate set of all possible transformations in  $W_{i_1} \times \dots \times W_{i_N}$ . Finally, we choose the best paraphrase combinations within the candidate set. In this way, we are able to conduct multiple replacements in one iteration and thus take into consideration the joint effect of multiple words replacements.

This method is based on an intuition derived from coordinate descent with the Gauss-Southwell rule [39] in the continuous optimization theory; normally, updating the coordinates with the highest absolute gradient values is provably faster than optimizing over random coordinates [45, 46]. We only conduct this method in word paraphrasing, since the gradient information of sentence embedding is less trustworthy. Usually sentence paraphrasing changes the number of words. The calculated gradient before paraphrasing step might not even correspond to the right position of the new sentence. Therefore it makes more sense to use the objective value only and goes back to our Algorithm 2.

# 6 Experiments

In this section, we provide empirical evidence of the advantages of our attack scheme via joint sentence and word paraphrasing on both two WCNN and LSTM models and various classification tasks. Our code for

replicating our experiments is available online<sup>1</sup>.

#### 6.1 Tasks and Models.

We focus on attacking the following state-of-the-art models which also echo our theoretical analysis:

#### • Word-level Convolutional Network (WCNN).

We implement a convolutional neural network [41] with a temporal convolutional layer of kernel size 3 and a max-pooling layer, followed by a fully connected layer for the classification output.

#### • Long Short Term Memory classifier (LSTM).

The LSTM Classifier [47] is well-suited to classifying text sequences of various lengths. We construct a one-layer LSTM with 512 hidden nodes, following the architecture used in [19, 48].

We carried out experiments on three different text classification tasks: fake-news detection, spam filtering and sentiment analysis; these tasks are also considered in [19]. The corresponding datasets include:

#### • Fake/Real News.

The fake news repository [49] contains 6336 clean articles of both fake and real news in a 1:1 ratio (5336 training and 1000 testing), with both left- and right-wing sites as sources.

#### • Trec07p (emails).

The TREC 2007 Public Spam Corpus (Trec07p) contains 75,419 messages of ham (non-spam) and spam in a 1:2 ratio. We preprocess the data and retain only the main content in each email. We randomly hold out 10% as testing data.

#### • Yelp reviews.

The Yelp reviews dataset was obtained from the Yelp Dataset Challenge in 2015. The polarity dataset we used was constructed for a binary classification task that labels 1 star as negative and 5 star as positive. The dataset contains 560,000 training and 38,000 testing documents.

#### 6.2 General Settings

For the training procedure, we use similar settings for the WCNN and LSTM classifier. We extracted the top 100,000 most frequent words to form the vocabulary. The first layer of both WCNN and LSTM is the embedding that transforms individual word into a 300-dimensional vector using the pretrained *word2vec* embeddings [50]. We randomly hold out 10% training data as validation set to choose the number of epochs and use a constant mini-batch size of 16.

We manually selected the hyperparameters for each dataset. We set the termination threshold  $\tau=0.7$ , and set a neighbor size k for possible paraphrases to be 15. We set the semantic similarity  $\delta_w=\delta_s=0.75^2$  for all datasets and syntactic bound  $\delta_2=2$  for news and yelp datasets, and  $\delta=\infty$  for Trec07p; the email dataset contains many corrupted words rendering the language model ineffective. For all datasets, we only allow  $\lambda_w=20\%$  word paraphrasing. We set the sentence paraphrasing ratio  $\lambda_s=20\%$  for yelp and news dataset, and for spam  $\lambda_s=60\%$ .

## 6.3 Accuracy comparisons.

After setting up the experimental environment, we now present the empirical studies in several aspects. In Table 2 we present the original and adversarial test accuracy on the three datasets with the two chosen models, where we allow 20% word replacements. We also include the presented adversarial accuracy from [19] for reference. Since the word neighboring sets for the two methods are different and the values are not

 $<sup>^1</sup>$ https://github.com/cecilialeiqi/adversarial\_text

 $<sup>^{2}</sup>$ We use the WMD similarity in python's spacy package. The similarity is in [0,1] basis where 1 means identical and 0 means complete irrelevant.

Dataset		WCNI	N		LSTM			
Dataset	Origin	Origin   ADV (ours)		ADV [19]		ADV (ours)	ADV [19]	
News	93.1%	35.4%	71.0%	70.5%*	93.3%	16.5%	37.0%	22.8%*
Trec07p	99.1%	48.6%	64.5%	63.5%*	99.7%	31.1%	39.8%	37.6%*
Yelp	93.6%	23.1%	39.0%	41.2%*	96.4%	30.0%	24.0%	29.2%*

Table 2: Classifier accuracy on each dataset. Origin and ADV respectively stand for the clean and adversarial testing results. For all datasets, we set word paraphrasing ratio to be  $\lambda_w = 20\%$  for our method (ADV(ours)). We include results from [19] for comparison. The first column indicates reported values in their paper; while the consequent column marked by asterisk is our implementation using greedy method in [19] and the same word neighboring set as our method. Both results use large  $\lambda_w = 50\%$  and allow many more word replacements.

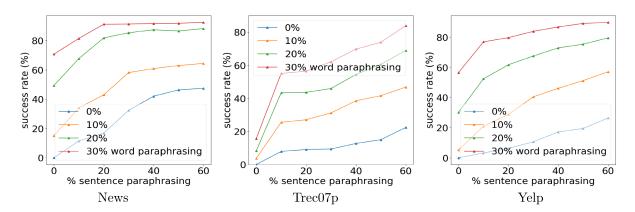


Figure 4: Success rate of attacking the LSTM classifier with different ratios of allowed paraphrasing.

directly comparable, one might argue that we have broaden the search space of words to make the problem easier. Therefore we also implemented the greedy mechanism in [19] using the same word replacement set as our method has chosen (marked by \*). Both the reported values from [19] and our implementation allow 50% word replacements. From Table 2 we can see that in both settings, we are able to successfully flip more prediction classes with fewer word paraphrases. We hereby conclude that joint sentence and word level paraphrasing is much more effective than mere word replacements. Meanwhile, since sentence-level attacks almost perfectly preserve the original meaning, our method can be less susceptible to humans. In the appendix we use some concrete examples to show the significantly improved quality of our generated adversarial texts compared to [19, 18].<sup>3</sup> In the examples, we can see that sometimes by simplifying or changing the language, or even by making the slightest changes like adding or erasing space, the sentence paraphrase can make a tremendous difference to the classifier output. Consequently, our method does far fewer word level alterations than other methods and greatly reduces the possibility of syntactic or grammar errors.

To further investigate the joint effect from combining sentence and word level attacks, we also study how each model is susceptible to different degrees of change permitted for both attack levels. Therefore we tested and presented the joint influence in Figure 4 for ratios of sentence paraphrasing  $\lambda_s$  ranging from 0% to 60%, as well as for allowed word paraphrasing percentages  $\lambda_w$ : 0%, 10%, 20% and 30%. In all datasets, sentence paraphrasing is especially effective when we allow only a few word paraphrases. For instance, in the sentiment analysis task, we could only successfully attack around 5% reviews by paraphrasing 10% of words. But after conducting 60% sentence paraphrasing beforehand, the success rate increases to almost 60%.

<sup>&</sup>lt;sup>3</sup>Since the former code is not available online, we implemented their algorithms. We use their chosen parameters to generate the adversarial examples to compare the quality of sentences in the appendix. While in Section 5.2 we use the same word neighboring sets for all algorithms to make a fair comparison of the optimization schemes.

Method	object	tive-guided	<u> </u>		nethod [18]	ours (Alg. 3)		
		$\lambda_w = 5\%$	$\lambda_w = 20\%$	$\lambda_w = 5\%$	$\lambda_w = 20\%$	$\lambda_w = 5\%$	$\lambda_w = 20\%$	
News	SR:	26.2%	28.4%	9.93%	12.8%	39.7%	45.4%	
news	time:	0.79	1.46	0.13	0.21	0.26	0.31	
Trec07p	SR:	5.1%	24.9%	0.86%	3.4%	12.9~%	45.3 %	
11eco7p	time:	0.19	0.33	0.03	0.05	0.07	0.09	
Yelp	SR:	12.7%	45.0%	4.2%	9.1%	20.7%	55.9%	
reib	time:	0.15	0.21	0.02	0.03	0.02	0.05	

Table 3: Attack success rate (denoted by SR) and time comparisons of each optimization mechanism. The performance is reported on the WCNN classifier. Here objective-guided greedy indicates the greedy method used in [19], and the gradient method is the one suggested in [18]. We can see that even when only applying Algorithm 3, our optimization method is more effective among others.

## 6.4 Optimization Method Comparisons for Word-level Attacks.

To investigate the effectiveness of our proposed gradient-guided greedy method, we implement and compare the time consumption and success rate with Algorithm 3 and the other two techniques: the gradient method [18] and the objective-guided greedy method [19]. To make a fair comparisons of the optimization schemes, we do not conduct sentence level paraphrasing in any of the methods, and we use the same hyperparameters and settings as suggested in Section 6.1. We observe that our scheme is especially more appealing to WCNN, partially because we used 5% dropout for inference. Recent work [51] indicates dropout not only works for training but also for inference as a Bayesian approximation. The small alteration of one word replacement per iteration [19] is not significant enough to be considered as true gains or the noise from the dropout. While our method replaces 5 words per iteration to capture more difference, thus it is easier to distinguish the change from the dropout randomness. From Table 3 we can see that our method requires only 1/5 to 1/3 time cost relative to the objective-guided greedy method and also achieves better success rate. On the other hand, gradient method fails to produce good performance when we allow a small set of word replacements.

Dataset		Task I		Task II			
	News	Trec07p	Yelp	News	Trec07p	Yelp	
					$3.23 \pm 0.31$		
Adversarial	50.0%	80.0%	100.0%	$3.13 \pm 0.50$	$3.10 \pm 0.40$	$2.10 \pm 1.05$	

Table 4: Human-subject validation. Task I measures classification accuracy while Task II the subjective likelihood that each example was crafted by a human (scale from 1 to 5). We used five participants, each shown n = 60 text examples, half original and half generated using our algorithm. The quality of the generated adversarial text (Task II) is near equal to the original and in fact, slightly higher for the Yelp dataset, but this finding is not necessarily statistically significant.

#### 6.5 Human Evaluation Validation

Despite the significantly higher attack proportion of our text examples, our aim is to deliver a message that is faithful to and coherent with the original text. To evaluate the quality of these generated text examples, we presented a number of original and adversarial text pairs (randomly shuffled before the test) to five human evaluators. The evaluators were asked to complete two tasks: I) Assign the correct label to each text sample; II) Rate each text sample with respect the the likelihood that was crafted by a human (scale from 1 to 5). We adopted a majority vote for task I, and averaged the results from five evaluators for task II. As shown in Table 4, we found that human evaluators tend to achieve similar performance for each kind of text in both

tasks, indicating that text examples generated via joint sentence and word paraphrasing are indeed coherent and faithful to the original texts in the relevant respects.

Dataset		LSTM		WCNN			
Dataset	News	Trec07p	Yelp	News	Trec07p	Yelp	
Test (before) Test (after)	93.3% 94.5%	99.7% 99.5%	96.4% 97.3%	93.1% 93.8%	99.1% 99.2%	93.6% 94.9%	
ADV (before) ADV (after)	16.5% 32.7%	31.1% 50.1%	30.0% 46.7%	35.4% 40.0%	48.6% 54.2%	23.1% 44.4%	

**Table 5:** Performance of adversarial training.

## 6.6 Adversarial Training.

Finally, we investigated whether our adversarial examples can help improve model robustness. For each dataset, we randomly selected 20% of the training data and generated adversarial examples from them using Algorithm 1. We then merged these adversarial examples with corrected labels into the training set and retrained the model. We present the testing and adversarial accuracy before and after this adversarial training process in Table 5. Under almost all circumstances, adversarial training improved the generalization of the model and made it less susceptible to attack.

## 7 Conclusion

In this paper, we propose a general framework for discrete adversarial attacks. Mathematically, we formulate the adversarial attack as an optimization task on a set of attacks. We then theoretically prove that greedy method guarantees a 1-1/e approximation factor for two classes of neural network for text classification task. Empirically, we propose a gradient-guided greedy method that inherits the efficiency of gradient method and ability to attack of greedy method. Specifically, we investigate joint sentence and word paraphrasing to generate attacking space that maintain the original semantics and syntax for text adversarial examples.

**Acknowledgements.** I.D. acknowledges the support of NSF via IIS-1546452 and CCF-1564000. A.D. acknowledges the support of NSF Grants 1618689, DMS 1723052, CCF 1763702, ARO YIP W911NF-14-1-0258 and research gifts by Google, Western Digital and NVIDIA.

## References

- [1] N. Dalvi, P. Domingos, S. Sanghai, D. Verma et al., "Adversarial classification," in *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining.* ACM, 2004, pp. 99–108.
- [2] D. Lowd and C. Meek, "Adversarial learning," in *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining.* ACM, 2005, pp. 641–647.
- [3] B. Biggio and F. Roli, "Wild patterns: Ten years after the rise of adversarial machine learning," arXiv preprint arXiv:1712.03141, 2017.
- [4] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," arXiv preprint arXiv:1312.6199, 2013.

- [5] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," arXiv preprint arXiv:1412.6572, 2014.
- [6] S. M. Moosavi Dezfooli, A. Fawzi, and P. Frossard, "Deepfool: a simple and accurate method to fool deep neural networks," in *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition* (CVPR), no. EPFL-CONF-218057, 2016.
- [7] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," in *Security and Privacy (SP)*, 2016 IEEE Symposium on. IEEE, 2016, pp. 582–597.
- [8] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *IEEE Symposium* on Security and Privacy, 2017, pp. 39–57.
- [9] I. Evtimov, K. Eykholt, E. Fernandes, T. Kohno, B. Li, A. Prakash, A. Rahmati, and D. Song, "Robust physical-world attacks on deep learning models," arXiv preprint arXiv:1707.08945, vol. 1, 2017.
- [10] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh, "ZOO: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models," in *ACM Workshop on Artificial Intelligence and Security*, 2017, pp. 15–26.
- [11] P.-Y. Chen, Y. Sharma, H. Zhang, J. Yi, and C.-J. Hsieh, "EAD: elastic-net attacks to deep neural networks via adversarial examples," AAAI, 2018.
- [12] D. Su, H. Zhang, H. Chen, J. Yi, P.-Y. Chen, and Y. Gao, "Is robustness the cost of accuracy?—a comprehensive study on the robustness of 18 deep image classification models," *ECCV*, 2018.
- [13] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," arXiv preprint arXiv:1706.06083, 2017.
- [14] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," arXiv preprint arXiv:1802.00420, 2018.
- [15] N. Papernot, P. McDaniel, A. Swami, and R. Harang, "Crafting adversarial input sequences for recurrent neural networks," in *Military Communications Conference*, MILCOM 2016-2016 IEEE. IEEE, 2016, pp. 49–54.
- [16] J. Li, W. Monroe, and D. Jurafsky, "Understanding neural networks through representation erasure," arXiv preprint arXiv:1612.08220, 2016.
- [17] J. Ebrahimi, A. Rao, D. Lowd, and D. Dou, "Hotflip: White-box adversarial examples for NLP," arXiv preprint arXiv:1712.06751, 2017.
- [18] Z. Gong, W. Wang, B. Li, D. Song, and W.-S. Ku, "Adversarial texts with gradient methods," arXiv preprint arXiv:1801.07175, 2018.
- [19] V. Kuleshov, S. Thakoor, T. Lau, and S. Ermon, "Adversarial examples for natural language classification problems," 2018. [Online]. Available: https://openreview.net/forum?id=r1QZ3zbAZ
- [20] P. Yang, J. Chen, C.-J. Hsieh, J.-L. Wang, and M. I. Jordan, "Greedy attack and Gumbel attack: Generating adversarial examples for discrete data," arXiv preprint arXiv:1805.12316, 2018.
- [21] R. Jia and P. Liang, "Adversarial examples for evaluating reading comprehension systems," arXiv preprint arXiv:1707.07328, 2017.
- [22] T. Miyato, A. M. Dai, and I. Goodfellow, "Adversarial training methods for semi-supervised text classification," arXiv preprint arXiv:1605.07725, 2016.

- [23] S. Samanta and S. Mehta, "Towards crafting text adversarial samples," arXiv preprint arXiv:1707.02812, 2017.
- [24] B. Liang, H. Li, M. Su, P. Bian, X. Li, and W. Shi, "Deep text classification can be fooled," arXiv preprint arXiv:1704.08006, 2017.
- [25] Y. Yao, B. Viswanath, J. Cryan, H. Zheng, and B. Y. Zhao, "Automated crowdturfing attacks and defenses in online review systems," in ACM Conference on Computer and Communications Security, 2017, pp. 1143–1158.
- [26] J. Gao, J. Lanchantin, M. L. Soffa, and Y. Qi, "Black-box generation of adversarial text sequences to evade deep learning classifiers," arXiv preprint arXiv:1801.04354, 2018.
- [27] M. Alzantot, Y. Sharma, A. Elgohary, B.-J. Ho, M. Srivastava, and K.-W. Chang, "Generating natural language adversarial examples," arXiv preprint arXiv:1804.07998, 2018.
- [28] C. Wong, "Dancin seq2seq: Fooling text classifiers with adversarial text example generation," arXiv preprint arXiv:1712.05419, 2017.
- [29] Z. Zhao, D. Dua, and S. Singh, "Generating natural adversarial examples," *ICLR*; arXiv preprint arXiv:1710.11342, 2018.
- [30] M. Cheng, J. Yi, H. Zhang, P.-Y. Chen, and C.-J. Hsieh, "Seq2sick: Evaluating the robustness of sequence-to-sequence models with adversarial examples," arXiv preprint arXiv:1803.01128, 2018.
- [31] M. T. Ribeiro, S. Singh, and C. Guestrin, "Semantically equivalent adversarial rules for debugging NLP models," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, vol. 1, 2018, pp. 856–865.
- [32] Y. Belinkov and Y. Bisk, "Synthetic and natural noise both break neural machine translation," arXiv preprint arXiv:1711.02173, 2017.
- [33] N. Mrkšić, D. O. Séaghdha, B. Thomson, M. Gašić, L. Rojas-Barahona, P.-H. Su, D. Vandyke, T.-H. Wen, and S. Young, "Counter-fitting word vectors to linguistic constraints," arXiv preprint arXiv:1603.00892, 2016.
- [34] M. Frank and P. Wolfe, "An algorithm for quadratic programming," Naval Research Logistics (NRL), vol. 3, no. 1-2, pp. 95–110, 1956.
- [35] H. Narayanan, Submodular functions and electrical networks. Elsevier, 1997, vol. 54.
- [36] S. Fujishige, Submodular functions and optimization. Elsevier, 2005, vol. 58.
- [37] A. Schrijver, Combinatorial optimization: polyhedra and efficiency. Springer Science & Business Media, 2003, vol. 24.
- [38] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher, "An analysis of approximations for maximizing submodular set functions—i," *Mathematical Programming*, vol. 14, no. 1, pp. 265–294, 1978.
- [39] J. Nutini, M. Schmidt, I. Laradji, M. Friedlander, and H. Koepke, "Coordinate descent converges faster with the Gauss-Southwell rule than random selection," in *International Conference on Machine Learning*, 2015, pp. 1632–1641.
- [40] J. Bilmes and W. Bai, "Deep submodular functions," arXiv preprint arXiv:1701.08939, 2017.
- [41] Y. Kim, "Convolutional neural networks for sentence classification," arXiv preprint arXiv:1408.5882, 2014.

- [42] J. Wieting, M. Bansal, K. Gimpel, K. Livescu, and D. Roth, "From paraphrase database to compositional paraphrase model and back," arXiv preprint arXiv:1506.03487, 2015.
- [43] J. Wieting and K. Gimpel, "Pushing the limits of paraphrastic sentence embeddings with millions of machine translations," in arXiv preprint arXiv:1711.05732, 2017.
- [44] M. Kusner, Y. Sun, N. Kolkin, and K. Weinberger, "From word embeddings to document distances," in *International Conference on Machine Learning*, 2015, pp. 957–966.
- [45] Q. Lei, K. Zhong, and I. S. Dhillon, "Coordinate-wise power method," in *Advances in Neural Information Processing Systems*, 2016, pp. 2064–2072.
- [46] Q. Lei, I. E. Yen, C.-y. Wu, I. S. Dhillon, and P. Ravikumar, "Doubly greedy primal-dual coordinate descent for sparse empirical risk minimization," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70.* JMLR. org, 2017, pp. 2034–2042.
- [47] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [48] X. Zhang, J. J. Zhao, and Y. LeCun, "Character-level convolutional networks for text classification," *CoRR*, vol. abs/1509.01626, 2015. [Online]. Available: http://arxiv.org/abs/1509.01626
- [49] G. McIntire, "Fake news dataset," https://github.com/GeorgeMcIntire/fake real news dataset, 2017.
- [50] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," arXiv preprint arXiv:1301.3781, 2013.
- [51] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *international conference on machine learning*, 2016, pp. 1050–1059.

## A Proofs

## A.1 Proof of Proposition 1

Proof of Proposition 1. We will show that even for a very simple function f, it could be reduced from subset sum problem when  $k \geq 2$ .

For instance, let

$$f(S) = \underset{\text{supp}(1) \subset S}{\arg \max} \left\| \sum_{i=1}^{n} V(x_i^{(l_i)}) - \mathbf{v} \right\|_2^2,$$

where the target is to find the best  $\ell_2$  approximation of some target vector  $\mathbf{v}$  from the embedding vectors.

For simplicity, denote the embedding vector of each paraphrased words to be  $V(x_i^{(j)}) = \mathbf{v}_i^{(j)}, 1 \le i \le n, 0 \le j \le k-1$ . Suppose there is an algorithm that solves the above problem in time polynomial to n. Then we will now show that the subset sum problem has a solution in polynomial time. Let the n numbers to be  $s_1, s_2, \dots s_n$ , and the target to be W. Then we let  $\mathbf{v}_i^{(0)} = [s_i, 0, 0, \dots, 0]$ , and  $\mathbf{v}_i^{(j)} = \mathbf{0}, j = 1, \dots k-1$ , with target  $\mathbf{v} = [W, 0, 0, \dots, 0]$ . Then just check if the best approximation of  $\mathbf{v}$  is exactly  $\mathbf{v}$  will suffice the subset sum problem. Therefore it contradicts with the fact that subset sum is in NP-complete class.

## A.2 Proof of Proposition 2

Proof of Proposition 2. Define set function  $h(S) = \arg\max_{\sup(\mathbf{l}) \subset S} V(T_{\mathbf{l}}(\mathbf{x}))^{\top} \nabla C_y(\mathbf{v})$ , where  $\mathbf{v} = V(\mathbf{x})$ . Denote  $\mathbf{g} = \nabla C_y(\mathbf{v})$ . When V is bag-of-words embedding, we denote the embedding of each paraphrased word in  $\mathbf{x} = [x_1, x_2, \cdots, x_n]$  as  $V(x_i^{(j)}) = \mathbf{e}_{d_{ij}}$ . Here for any i,  $\mathbf{e}_i$  is defined as the one-hot vector with 1 in index i and 0 elsewhere. Then  $V(T_{\mathbf{l}}(\mathbf{x})) = \sum_{i=1}^{n} \mathbf{e}_{d_{il_i}}$ .

$$\begin{split} h(S) &= \underset{\text{supp}(\mathbf{l}) \subset S}{\arg\max} \ V(T_{\mathbf{l}}(\mathbf{x}))^{\top} \mathbf{g} \\ &= \underset{\text{supp}(\mathbf{l}) \subset S}{\arg\max} \sum_{i=1}^{n} \mathbf{e}_{d_{il_{i}}}^{\top} \mathbf{g} \\ &= \underset{\text{supp}(\mathbf{l}) \subset S}{\arg\max} \sum_{i=1}^{n} g_{d_{il_{i}}} \\ &= \mathbb{1}_{S}^{\top} \mathbf{w}, \end{split}$$

where  $w_i = \max_{0 \le t \le k-1} g_{d_{it}}$ .

When V is d-dimentional word2vec embedding, the embedding  $V(\mathbf{x}) = [V(x_1)^\top | V(x_2)^\top | \cdots | V(x_n)^\top]^\top \in \mathbb{R}^{nd}$ . Denote  $\hat{\mathbf{g}}_i = \mathbf{g}_{(id-d+1):id}$  to be the gradient with respect to the word  $w_i$ .

$$h(S) = \underset{\text{supp}(\mathbf{l}) \subset S}{\arg \max} V(T_{\mathbf{l}}(\mathbf{x}))^{\top} \mathbf{g}$$
$$= \underset{\text{supp}(\mathbf{l}) \subset S}{\arg \max} \sum_{i=1}^{n} V(x_{i}^{(l_{i})})^{\top} \hat{\mathbf{g}}_{i}$$
$$= \mathbb{1}_{S}^{\top} \mathbf{w},$$

where  $w_i = \max_{0 \le t \le k-1} V(x_i^{(t)})^{\top} \hat{\mathbf{g}}_i$ . Therefore for both bag-of-words embedding and word2vec embedding, h is a modular (linear) set function, and Problem 2 is solvable in polynomial time.

#### A.3Proof of Claim 1

Proof of Claim 1. Clearly for any  $S \subset V \subset [n]$ ,

$$f(S) = \max_{\text{supp}(\mathbf{l}) \subset S} C_y(V(T_{\mathbf{l}}(\mathbf{x}))) \le \max_{\text{supp}(\mathbf{l}) \subset V} C_y(V(T_{\mathbf{l}}(\mathbf{x})))$$
 (since  $S \subset T$ )
$$= f(V)$$

Therefore the set function f is non-decreasing. Since the problem of maximizing a monotone submodular function subject to a cardinality constraint admits a 1-1/e approximation algorithm[38], Problem 1 can be solved in time polynomial to n with greedy method. 

#### Proof of Theorem 1

Proof of Theorem 1. We start from a simple case, h = 1, i.e., a unit kernel size, and we look at a single feature corresponding to one filter, i.e.  $\hat{c}_j = \max_{i=1}^n c_{ij}$ .

To further incorporate the transformation to the input, we rewrite  $\hat{c}_j$  as a function of the transformation index 1.

$$\hat{c}_j(\mathbf{l}) \equiv \max_{i=1}^n \phi(\mathbf{w}_j^\top V(x_i^{(l_i)}) + b_j) = \max_{i=1}^n v_{ij}^{(l_i)},$$

where  $\mathbf{w}_j$  is the j-th filter and we denote  $v_{ij}^{(k)} = \phi(\mathbf{w}_j^\top V(x_i^{(k)}) + b_j)$  for simplicity. Let S, T denote two sets that satisfy  $S \subset T \subset [n]$ . For any two vectors  $\mathbf{l}^S$  and  $\mathbf{l}^T$  satisfy that  $l_i^S = l_i^T, \forall i \in S$ , and  $\operatorname{supp}(\mathbf{l}^S) = S$ ,  $\operatorname{supp}(\mathbf{l}^T) = T$ . With the assumption that  $\mathbf{w}_j^\top V(x_i) \leq \mathbf{w}_j^\top V(x_i^{(t)})$ , and since the activation function is non-decreasing, we have  $v_{ij}^{(0)} \leq v_{ij}^{(t)}, \forall i \in [n], j \in [m], t \in [k-1]$ , and hereby  $\hat{c}_j(\mathbf{l}^S) \leq \hat{c}_j(\mathbf{l}^T)$ . Therefore for any new element's position s and its replacement index t, we have

$$\begin{split} & \hat{c}_{j}(\mathbf{l}^{S} + t\mathbf{e}_{s}) - \hat{c}_{j}(\mathbf{l}^{S}) = \max\{v_{sj}^{(t)} - \hat{c}_{j}(\mathbf{l}^{S}), 0\} \\ & \geq \max\{v_{sj}^{(t)} - \hat{c}_{j}(\mathbf{l}^{T}), 0\} \\ & = & \hat{c}_{j}(\mathbf{l}^{T} + t\mathbf{e}_{s}) - \hat{c}_{j}(\mathbf{l}^{T}). \end{split} \tag{since } \hat{c}_{j}(\mathbf{l}^{S}) \leq \hat{c}_{j}(\mathbf{l}^{T})$$

Since the final output probability is a positive weighted summation of each  $\hat{c}_i$ , it also satisfies

$$C^{\text{WCNN}}(\mathbf{l}^S + t\mathbf{e}_s) - C^{\text{WCNN}}(\mathbf{l}^S) \geq C^{\text{WCNN}}(\mathbf{l}^T + t\mathbf{e}_s) - C^{\text{WCNN}}(\mathbf{l}^T)$$

Taking the max over all  $\mathbf{l}^S$ ,  $\mathbf{l}^T$  we have:

$$f(S + \{s\}) = \max_{t=1}^{k-1} \max_{\text{supp}(\mathbf{l}^S) = S} C^{\text{WCNN}}(\mathbf{l}^S + t\mathbf{e}_s)$$

Therefore

$$f(S + \{s\}) - f(S)$$

$$= \max_{t=1}^{k-1} \left\{ \max_{\text{supp}(\mathbf{l}^S)=S} \left\{ C^{\text{WCNN}}(\mathbf{l}^S + te_s) - C^{\text{WCNN}}(\mathbf{l}^S) \right\} \right\}$$

$$\geq \max_{t=1}^{k-1} \left\{ \max_{\text{supp}(\mathbf{l}^T)=T} \left\{ C^{\text{WCNN}}(\mathbf{l}^T + te_s) - C^{\text{WCNN}}(\mathbf{l}^T) \right\} \right\}$$

$$= f(T + \{s\}) - f(T).$$
(from (6))

The case when  $2 \le h \le s$  is essentially the same with h = 1 since each window has no overlapping. We could simply replace  $\mathbf{v}_1$  by  $\mathbf{v}_{1:h}$  and conduct the same analysis. 

#### A.5 Proof of Theorem 2

Proof of Theorem 2. Recall that the hidden state node  $h_i$  is defined recursively as:

$$h_0 = C,$$
 (C is constant)  
 $h_i = \phi(wh_{i-1} + \mathbf{m}^\top V(x_{i-1}) + b).$ 

And the classifier output is  $C^{RNN}(V(\mathbf{x})) = yh_T$ .

For simplicity, we denote  $v_i^{(j)} \equiv \mathbf{m}^\top V(x_i^{(j)}) + b$ . Since we will only look for the transformation that maximizes the classifier output, without loss of generality, we assume  $v_i^{(j)} \geq v_i^{(0)}, \forall i \in [T], j \in [k-1]$ .

For a fixed input  $\mathbf{x} = [x_1, x_2, \dots, x_T]$  and transformation index  $\mathbf{l}$ , we want to study how changing an intermediate hidden state affects the consecutive layers' output. Therefore we represent the value of a j-th hidden state as a function of the i-th hidden node and the transformation label  $\mathbf{l}$ , that captures the network from i-th through j-th time steps, i.e.,

$$f_{i:j}(h_i, \mathbf{l}) = \phi\left(w \cdots \phi(wh_i + v_i^{(l_i)}) + \cdots + v_{j-1}^{(l_{j-1})}\right).$$

Finally we want to study the whole network's output  $yf_{0:T}(C,1)$ . We first prove the following lemma:

#### Lemma 1.

$$f_{i:j}(h_i + \delta, \mathbf{l}) - f_{i:j}(h_i, \mathbf{l}) \ge f_{i:j}(h_i + \delta, \mathbf{l} + t\mathbf{e}_s) - f_{i:j}(h_i, \mathbf{l} + t\mathbf{e}_s), \tag{6}$$

for any  $0 \le i < j \le T, t \in [k-1], s \in [T], s \notin supp(1), \delta > 0$ .

Proof of Lemma 1.

$$f_{i:j}(h_i + \delta, \mathbf{l} + t\mathbf{e}_s)$$
=  $f_{s+1:j}(\phi(wf_{i:s}(h_i + \delta, \mathbf{l} + t\mathbf{e}_s) + v_s^{(t)}), \mathbf{l} + t\mathbf{e}_s)$   
=  $f_{s+1:j}(\phi(wf_{i:s}(h_i + \delta, \mathbf{l}) + v_s^{(t)}), \mathbf{l})$ 

Now we simplify the equation by define  $a(\delta,t) = \phi(wf_{i:s}(h_i+\delta,\mathbf{l})+v_s^{(t)}), \delta \in \mathbb{R}, t \in [k-1]$ . Therefore we could rewrite the four terms in Eqn. (6) as:  $f_{s+1:j}(a(\delta,0),\mathbf{l}), f_{s+1:j}(a(0,0),\mathbf{l}), f_{s+1:j}(a(\delta,t),\mathbf{l}), \text{ and } f_{s+1:j}(a(0,t),\mathbf{l}).$  Since  $\phi$  is concave and  $v_s^{(t)} \geq v_s^{(0)}$ , notice

$$a(\delta, t) - a(0, t) \le a(\delta, 0) - a(0, 0).$$
 (7)

Now since  $f_{s+1:j}(\cdot, \mathbf{l})$  is a composite of concave function and is also concave, we have:

Lemma 1 could be extended to a more general form. Suppose two indices  $\mathbf{l}^S$  and  $\mathbf{l}^U$  satisfy supp $(\mathbf{l}^S) = S$ , supp $(\mathbf{l}^U) = U, S \subset U$ , and  $l_i^S = l_i^U, \forall i \in S$ . Since we could write  $\mathbf{l}^U$  as  $\mathbf{l}^S + \sum_{i \in U \setminus S} l_i^U \mathbf{e}_i$ , by repeatedly using Lemma 1 we have:

$$f_{i:j}(h_i + \delta, \mathbf{l}^S) - f_{i:j}(h_i, \mathbf{l}^S) \ge f_{i:j}(h_i + \delta, \mathbf{l}^U) - f_{i:j}(h_i, \mathbf{l}^U).$$

This conclusion basically claims an increase into an intermediate layer of the network will have smaller effect to the output when the network is attacked on more time steps. Then back to Theorem 2. Now consider we add a coordinate s to the set S and U,  $s \notin \text{supp}(S) \cup \text{supp}(U)$ .

$$f_{0:T}(C, \mathbf{l}^{S} + t\mathbf{e}_{s}) - f_{0:T}(C, \mathbf{l}^{S})$$

$$= f_{s:T}(\phi(wf_{0:s-1}(C, \mathbf{l}^{S}) + v_{s}^{(t)}), \mathbf{l}^{S}) - f_{s:T}(\phi(wf_{0:s-1}(C, \mathbf{l}^{S}) + v_{s}^{(0)}), \mathbf{l}^{S})$$

$$\geq f_{s:T}(\phi(wf_{0:s-1}(C, \mathbf{l}^{S}) + v_{s}^{(t)}), \mathbf{l}^{U}) - f_{s:T}(\phi(wf_{0:s-1}(C, \mathbf{l}^{S}) + v_{s}^{(0)}), \mathbf{l}^{U})$$

$$(\text{from (9) and since } \phi \text{ is non-decreasing, } v_{s}^{(t)} \geq v_{s}^{(0)})$$

$$\geq f_{s:T}(\phi(wf_{0:s-1}(C, \mathbf{l}^{U}) + v_{s}^{(t)}), \mathbf{l}^{U}) - f_{s:T}(\phi(wf_{0:s-1}(C, \mathbf{l}^{U}) + v_{s}^{(0)}), \mathbf{l}^{U})$$

$$(\text{since } f_{s:T}(\cdot, \mathbf{l}^{U}) \text{ is concave and similar analysis as (8))}$$

$$= f_{0:T}(C, \mathbf{l}^{U} + t\mathbf{e}_{s}) - f_{0:T}(C, \mathbf{l}^{U})$$

$$(9)$$

Finally, since

$$\max_{\operatorname{supp}(\mathbf{l}) \subset S \cup \{s\}} C^{\operatorname{RNN}}(V(T_{\mathbf{l}}(\mathbf{x}))) = \max_{\operatorname{supp}(\mathbf{l}^S) \subset S} \max_{t \in [k-1]} y f_{0:T}(C, \mathbf{l}^S + t\mathbf{e}_s),$$

we have:

$$\begin{aligned} & \max_{\sup_{l} (l) \subset S \cup \{s\}} C^{\text{RNN}}(V(T_{\mathbf{l}}(\mathbf{x}))) - \max_{\sup_{l} (l) \subset S} C^{\text{RNN}}(V(T_{\mathbf{l}}(\mathbf{x}))) \\ &= \max_{\sup_{l} (l) \subset S} (\max_{t \in [k-1]} y f_{0:T}(C, \mathbf{l}^S + t\mathbf{e}_s) - y f_{0:T}(C, \mathbf{l}^S)) \\ &\geq \max_{\sup_{l} (l) \subset U} (\max_{t \in [k-1]} y f_{0:T}(C, \mathbf{l}^U + t\mathbf{e}_s) - y f_{0:T}(C, \mathbf{l}^U)) \text{ (since (9) holds for any } t) \\ &= \max_{\sup_{l} (l) \subset U \cup \{s\}} C^{\text{RNN}}(V(T_{\mathbf{l}}(\mathbf{x}))) \\ &- \max_{\sup_{l} (l) \subset U} C^{\text{RNN}}(V(T_{\mathbf{l}}(\mathbf{x}))) \end{aligned}$$

## B Data statistics

Task Dataset #Train#Test Trec07p Spam filtering 67.9k7.5kYelp Sentiment analysis 560k38kNews Fake news detection 5.3k1.0k

Table 6: Statistics of each datasets

## C Comparisons with other methods with concrete examples

In this section, we provide some concrete examples to compare our method with the other related methods. The following six examples respectively show the combinations of three datasets (fake news, Trec07p, and yelp) as well as the two models we use (LSTM and WCNN).

We use red font to denote changes from sentence level paraphrasing and blue for word paraphrasing.

## C.1 Empirical example 1: Task - Fake news detection. Classifier-CNN.

Method: Ours. Origin: 100% Real. ADV: 71% Fake

6 Six detainees detained in raids in Belgium Brussels, Belgium (CNN) Police detained six people in raids raid Thursday night as when investigators raced were sent to uncover the network behind this week's terror attacks in the Belgian capital. The Belgian federal prosecutor's office didn't provide details about who had been detained in the Brussels raids, why they had been apprehended or whether they will face charges. It will be decided tomorrow if these people will remain in custody, the office said in a statement released late Thursday. Two people were taken into custody in Brussels' Jette neighborhood, one person was detained in a different part of the capital, and three people were in a vehicle in front of the federal prosecutor's office when authorities apprehended them, public broadcaster RTBF reported. So far, authorities have Authorities said they believe five men played a part people took a shot in Tuesday's bombings in Belgium that killed 31 people and injured wounded 330. Three of the attackers are dead. Two of them could still be on the loose. Investigators are combing over evidence from surveillance footage and the explosives stash they seized from an apparent hideaway in a suburb. Sweeps where investigators detain people first and ask questions later are likely to become an increasingly common tactic, CNN national security analyst Juliette Kayyem said. There will be lots more of them, she said. They are going to be what's called overbroad. They are going to just try to find people or evidence that may stop the next terrorism attack, and they will figure out who they have under custody. Khalid El Bakraoui, one of the terrorists who bombed a train near the Maelbeek metro station, is dead. Authorities believe a second unidentified person was also involved in that attack, a senior Belgian security source told CNN. But investigators don't know where that the suspect is – or whether he's he was dead or alive. Surveillance footage shows the man holding a large bag at the station, according to Belgian public broadcaster RTBF. It's not clear if he was among the at least 20 killed in that blast, RTBF said. Authorities have released a grainy image of another suspect who they believe is on the run. That man, they say, shown in photographs wearing a black hat, was one of three attackers at Brussels Airport. Authorities say he planted a bomb at the airport and left. The other two men in the photographs are believed to be the suicide bombers. Fair to ask whether 'we missed the chance' Did Belgian authorities miss a chance to stop at least one of the suspects involved in the attacks? Bakraoui had been sentenced to nine years in prison in Belgium back in 2010 for opening fire on police officers with a Kalashnikov during a robbery, according to broadcaster RTBF and CNN affiliate RTL. Needless to say, he didn't serve all that time. Given the facts, it is justified that ... people ask how it is possible that someone was released early and we missed the chance when he was in Turkey to detain him, said Jambon, whose offer to resign was rebuffed by Prime Minister Charles Michel. Investigators suspect Abdeslam planned to be part of an attack by the same ISIS cell that lashed out Tuesday, a senior Belgian counterterrorism official told CNN's Paul Cruickshank. Authorities looked Wednesday at the Brussels homes of the Bakraoui brothers. These These two searches findings were not eonclusive decisive, the federal prosecutor's office prosecutors said. Homes were searched Thursday in several areas in and around the city, officials said. One operation in the neighborhood of Schaerbeek stretched for hours into Friday morning. Investigators sealed off streets for several blocks. It was not immediately clear why such a large area had been cordoned. Masked teams in hazmat gear could be seen exiting a building and heading toward a police van. As investigations continue, a larger question looms: What could happen next? Not long ago, Western authorities believed ISIS was focused on taking territory in Syria and Iraq, not lashing out elsewhere. But U.S. officials now think the extremist group has been sending trained militants to Europe for some time. These men don't necessarily follow orders directly from ISIS headquarters. But they build on what they've learned, as well as a shared philosophy and approach, to develop their own terror cells and hatch their own plots. How many more ISIS militants are in Europe, poised to attack? That's not clear. For now, though, the top priority is tracking down the two men linked directly to Tuesday's terror.

Method: Greedy[19]. Origin: 100% Real. ADV: 79% Fake

6 7 detained detention in raids in Belgium Brussels, Belgium (CNN) Police cops detained deported six people in raids Thursday night as investigators investigation raced to uncover the network behind this week's terror terrorists attacks in the Belgian capital. The Belgian federal prosecutor's office didn't provide details about who had been detained in the Brussels raids, why they had been apprehended or whether they we will should face eyes charges. It will be decided tomorrow if these people will remain in custody, the office said

told in a statement stating released late Thursday. Two people were taken into custody in Brussels' Jette neighborhood, one person was detained detention in a different part of the capital, and three people were in a vehicle in front of the federal prosecutor's office when authorities apprehended them, public broadcaster RTBF reported. So far, authorities have said they believe five men played a part in Tuesday's bombings in Belgium what killed wounded 31 26 people individuals and injured 330. Three of the attackers are dead. Two of them could still be on of loose. Investigators Investigating are they combing over evidence from surveillance footage filmed and the explosives stash they never seized from an apparent obvious hideaway in a suburb. Sweeps where investigators detain people first and ask questions later are likely to become an increasingly commonly tactic, CNN national Security analysts Juliette Kayyem said say. There will be lots more of them, she knew said say. They are going to be what's called overbroad. They are going to just try to find people or and/or evidence findings that may stop the of next before Terrorism attack, and they will figure out who they have under custody. Khalid El Bakraoui, one of the terrorists who bombed a train near the Maelbeek metro station, is dead. Authorities believe a second unidentified person was also involved in that attack, a senior Belgian Security source sources told saying CNN. But investigation don't know where that suspect victim is be - or any whether not he's dead dying or and/or alive. Surveillance footage shows the man holding a large bag at the station, according to Belgian public broadcaster RTBF. It's not clear if he was among the at least 20 26 killed kill in that blast, RTBF said say. Authorities have released a another grainy image of the another suspect who they believe is on the run. That man, they say, shown in photographs photo wearing wear a one black red hat, was one of three twelve attackers attacker at Brussels Airport. Authorities say he planted a bomb at the airport and left. The other two men in the photographs are believed supposedly to be the suicide suicidal bombers. Fair to ask whether 'we missed the chance' Did Belgian authorities miss a chance to stop at least one of the suspects involved in the attacks? Bakraoui had been sentenced to nine years in prison in Belgium back in 2010 for opening fire on police policemen officers deputies with a Kalashnikov during a robbery, according to broadcaster RTBF and CNN affiliates RTL. Needless to say, he didn't serve all that time. Given the facts, it is justified that actually ... people everyone ask tell how it is possible that what someone was released early and we you missed of chance when he was in Turkey to detain him, said told Jambon, whose offer to resign was rebuffed by Prime Minister Charles Michel. Investigation suspected Abdeslam planning to be part of an attack enemy by the same different ISIS cells that lashed out Tuesday, a senior junior Belgian counter-terrorism unofficial told CNN's Paul Cruickshank. Authorities looked Wednesday at the Brussels homes of the Bakraoui brothers. Those Them two one searches were not conclusive, the federal prosecutor's Office said say. Homes Houses were searched Thursday in several areas in and around the city, officials authorities said say. One operation in the neighborhood of Schaerbeek stretched for hours into Friday morning. Investigators sealed off streets for several blocks. It was not immediately clear why such a large area had been cordoned. Masked teams in hazmat gear could be seen exiting a building and heading toward a police van. As investigations continue, a larger question looms: What could happen next? Not long ago, Western authorities believed ISIS was focused on taking territory in Syria and Iraq, not lashing out elsewhere. But U.S. officials now think the extremist group has been sending trained militants to Europe for some time. These men don't necessarily follow orders directly from ISIS headquarters. But they build on what they've learned, as well as a shared philosophy and approach, to develop their own terror cells and hatch their own plots. How many more ISIS militants are in Europe, poised to attack? That's not clear. For now, though, the top priority is tracking down the two men linked directly to Tuesday's terror.

Method: Gradient method[18]. Origin: 100% Real. ADV: 99.5% Real

6 detained in raids in Belgium Brussels, Belgium (CNN) Police detained six people in raids Thursday night as well investigators raced rode to uncover the of network networks behind this it week's terror attacks in the of Belgian capital. The Belgian federal prosecutor's office didn't provide details about who had been detained arrested in the of Brussels raids, why they have had been being apprehended or whether they will be face charges. It will should decided tomorrow if these people will be remain remains in custody, the office offices said in a statement released late Thursday. Two Three people were taken into custody in Brussels' Jette neighborhood, one another person was detained in a different part of the the capital, and three people

were had in a vehicle in front of the of federal prosecutor's office when authorities officials apprehended them, public broadcaster RTBF reported. So far, authorities have said they believe five men played playing a another part in Tuesday's bombings in Belgium that killed 31 29 people and injured 330. Three of the attacker be dead. Two of them could still be on the loose. Investigators are combing over evidence from surveillance footage filmed and the explosives stash they seized from an apparent hideaway in a another suburb. Sweeps where investigator detain people first second and ask questions later are likely to become an increasingly commonly tactic, CNN national security analyst Juliette Kayvem said. There will be lot more of them, she said guess. They are going to able be should what's ealled known overbroad. They are going to able just trying to find people or evidence that may stop the next Terrorism attack, and they will figures out up whom they have under custody. Khalid El Bakraoui, one of the terrorists who bombed a train near the Maelbeek metro station, is dead. Authorities Officials believe a second unidentified person was also involved in that attack, a senior Belgian Security sources told talked CNN. Blut Though investigators don't know think where there that because suspect is – or whether if he's dead dying or either alive. Surveillance footage shows the man holding a large bag at the station, according to Belgian public broadcaster RTBF. It's not clear if he was among the of at least 20 25 killed in that because blast, RTBF said guess. Authorities have 've released another grainy image of another one suspect who they have believe is on the run. That man, they say, shown in photographs wearing a black grey hat, was one another of three attackers at Brussels Airport. Authorities say he planted a bomb at the airport and left. The other two women in of photographs are believed to be the suicidal bombers. Fair to able ask tell whether 'we missed the chance' Did Belgian authorities miss a chance to stopping at least one of the suspects involved in of attacks? Bakraoui had came been being sentenced convicted to nine years in prison in Belgium Netherlands back in 2010 for opening closing fires on policemen officers policemen with another Kalashnikov during a robbery, according to broadcaster RTBF and CNN affiliate RTL. Needless to say, he didn't serve all that time. Given the facts, it is justified that ... people ask tell how what it is possible that someone was released early and we missed of chance when He was in Turkey to detaining him, said Jambon, whose offering to resign was rebuffed by Prime Minister Charles Michel. Investigators suspect Abdeslam planned to able be should part of the an attack by the same ISIS cell that lashed out Tuesday, a junior Belgian counter-terrorism unofficial told CNN's Paul Cruickshank. Authorities looked seemed Wednesday at the Brussels homes of the Bakraoui brothers. Those two searches were not conclusive, the federal prosecutor's office said. Homes were searched Thursday in several numerous areas in and around the city, officials said guess. One operation in the of neighborhood of Schaerbeek stretched for hours into Friday morning. Investigators sealed off streets for several blocks. It was not immediately clear why such a large areas had been being cordoned. Masked teams in hazmat gears could be seen exiting another buildings and heading towards another police van. As investigations continue, a larger sized questions looms: What could happen next? Not long ago, Western authorities believed ISIS was focused on taking territory in Syria and Iraq, not lashing out elsewhere. But U.S. officials authorities now think know the extremist group has being sending trained militants insurgents to Europe for some time. These men don't necessarily follow orders directly from ISIS headquarters. But they build on what they've learned, as well as a shared philosophy and approach, to develop their own terror cells and hatch their own your plots. How many more less ISIS militants are these in Europe, poised to attack? That's not clear. For now, though, the top priority is tracking down the two men linked directly to Tuesday's terror.

#### C.2 Empirical example 2: Task - Fake news detection. Classifier - LSTM.

Method: Ours. Origin: 100% Fake. ADV: 77% Real

Man Guy punctuates high-speed chase with stop at In-N-Out Burger drive-thru Print [Ed.—Well, that's a new one. Okay, that 's a new one.] A One man is in custody after leading police on a bizarre chase into the east Valley on Wednesday night. Phoenix police began has begun following the suspect in Phoenix and the pursuit continued into the east Valley, but it took a bizarre turn when the suspect stopped at an In-N-Out Burger restaurant's drive-thru drive-through near Priest and Ray Roads in Chandler. The suspect appeared to order food, but then drove away and got out of his pickup truck near Rock Wren Way and Ray Road. He then ran into a backyard ran to the backyard and tried to get into a house through the back door get in the home.

Mehod: greedy. Origin: 100% Fake. ADV: 86% Fake.

Man Guy punctuates high-speed chase with stopping at In-N-Out Burger drive-thru Print [Ed. - Well, that's a new one.] Another man is which in custody after earlier leading police officers on a bizarre chase into out the east north Valley on Wednesday night. Phoenix police arrested began begun following the of suspects in Phoenix and the pursuit pursuing continued into the eastern Valley, but though it took a bizarre turning when the of suspects stopping at an In-N-Out Burger restaurant's drive-thru nearby Priest and Ray Roads in Chandler. The suspect appeared to order food, but then drove away and got out of his pickup truck near Rock Wren Way and Ray Road. He then ran into a backyard and tried to get into a house through the back door.

Method: gradient method[18]. Origin: 100% Fake. ADV:  $1 - 2.5e^{-3}$  Fake.

Man punctuates high-speed chase with stopping at In-N-Out Burgers drive-thru Print [Ed. - Well, that's a new one.] A man is in custody after leading police arrest on a bizarre chase into the east Valley on Wednesday night. Phoenix police began following the suspect in Phoenix and the pursuit continued into the of east west Valley, but it took a bizarre turn then when that the of suspect stopped at an In-N-Out Burger restaurant's drive-thru nearby Priest and Ray Roads in Chandler. The suspect appeared to ordering food, but then drove away and got out of his pickup pick-up truck near Rocks Wren Chickadee Ways and Ray Road. He then ran into a backyard and tried to get into a house through the back again door.

## C.3 Empirical example 3: Task - Spam filtering. Classifier - WCNN.

Method: Ours. Origin: 100% Spam. ADV: 77% Ham

Become Fit For Life! HGH is a very complex molecule produced by the anterior lobe of the pituitary gland, which is located at the base of the brain. While it stimulates growth in children, it is important for maintaining a healthy body healthy bodies composition and well-being in adults. It is the primary hormone estrogen that controls many several of the body's organs and it stimulates tissue repair, brains functions, cell replacement, and enzyme function. Determining the levels of IGF-1 (Insulin Growth Factor) is how we measure HGH in the body. Receive a younger future potential with HGH

Method: Greedy[19]. Origin: 100% Spam. ADV: 71% Ham

Become Fit For Life! HGH is a very fairly complex molecule produced by the anterior lobe of the the of pituitary gland, which is has located at the base of the brain. While it that stimulates growth growing in children, it what is important significant for maintaining another healthy body bodies composition and well-being in adults. It is the primary secondary hormone progesterone that could controls many several of the body's organs and it that stimulates tissue repair, brains functions, cell replacement, and enzyme function. Determining Determine the levels of IGF-1 (Insulin Growth Factor) is which how understand we measure HGH in the of body. Receive a younger future with HGH

Method: Gradient method[18]. Origin: 100% Spam. ADV:  $1-2.7e^{-5}$  spam

Become Fit For Life! HGH is a very complex molecule produced by the anterior lobe of the pituitary gland, which that is located situated at the base of the brain. While it stimulates growth in children, it but is has important for maintaining a healthy bodies compositions and well-being in adults. It is the primary secondary hormones that controls many of the body's organs and it but stimulates tissues repair, brains functions, cell replacement, and enzyme function. Determining the levels of IGF-1 (Insulin Growthing Factor) is how we measure HGH in the of body. Receive a younger future with HGH

#### C.4 Empirical Example 4: Task - Spam filtering. Classifier - LSTM.

Method: Ours. Origin: 100% Ham. ADV: 87% Spam

I've always run jigdo-lite against my own mirror. It provides offers two couple things:

- 1) Proves I can you are able to build the ISOs from what I have mirrored locally.
- 2) Doesn't waste additional bandwidth. As long as the checksums match what is provided from the official ISO image masters site, I don't see what the difference would be. Anyone else do this? :) ^\_^ Will Simon Paillard wrote:
- > On Mon, Apr 09, 2007 at 08:43:07AM -0400, Jean-Francois Chevrette wrote:

```
> Hi.
> >
>> does anyone have another straightforward guide on how to use jigdo to build and mirror ISOs? I've been
reading both jigdo documentation and debian's webpage web-site on the subjet and it just won't work.
> Maybe with this one:
http://www.debian.org/CD/mirroring/#jigdomirror
> and the related links?
> Best regards,
   Method: Greedy[19]. Origin: 100% Ham. ADV: 90% Spam
I've always run jigdo-lite against my myself own mirror. It provides offers two five things:
1) Proves I u can reliably build the ISOs from what something I im have 've mirrored locally.
2) Doesn't waste additional extra bandwidth. As long as the checksums matches what is provided from the
unofficial ISO image master site, I thats don't see 'll what the difference would be. Anyone Somebody else do
this you? :)!!
Will Must Simon Paillard wrote:
> On Mon, Apr 09, 2007 at 08:43:07AM -0400, Jean-Francois Chevrette wrote:
> Hi,
> >
>> does anyone have another straightforward guide on how to able use jigdo to build and mirror ISOs? I've
been reading writing both other jigdo documentation and debian's webpage on the subjet and it just won't
work.
> Maybe with this one:
http://www.debian.org/CD/mirroring/#jigdomirror
> and the related links?
> Best regards,
   Method: Gradient method[18]. Origin: 100% Ham. ADV: 1-2.2e^{-15} Ham
I've always run jigdo-lite against my myself own mirror. It What provides two things:
1) Proves I you can reliably building the ISOs from what I have mirrored locally.
2) Doesn't waste additional bandwidth. As long as the checksums match what is provided from the unofficial
ISO image master site, I u don't see what the difference would could be. Anyone else do this? :):-)
Will Simon Paillard wrote:
> On Mon, Apr 09, 2007 at 08:43:07AM -0400, Jean-Francois Chevrette wrote:
> Hi,
> >
>> does anyone have a straightforward guide on how what to able use jigdo to build and mirror ISOs? I've
been reading both jigdo documentation and debian's webpage web-site on the subjet and it just won't work.
>> Maybe with this one:
http://www.debian.org/CD/mirroring/#jigdomirror
> and the related links?
> Best regards,
```

## C.5 Empirical Example 5: Task - Sentiment analysis. Classifier - CNN.

Method: Ours. Origin: 100% Positive. ADV: 93% Negative

This Starbucks location is located in the Bally's Grand Bazaar Shops. It's open 24/7 and it is huge. There is plenty of seating. Most of the seating is stadium type seating with benches. They also have an out door patio. The staff is very friendly and attentive to the guests. I do notice that they are under staffed sometimes

when they are busy. They 'll get your drinks out pretty fast though. Also, this location place is not owned by the easine property so they don't do n't charge outrageous prices like the location as a place on the an Linq promenade does. Definitely one of my favorite Starbucks stores. Stop by if your on the Strip.

Method: Greedy[19]. Origin: 100% Positive. ADV: 74% Negative

This Starbucks location is be located in the Bally's Grand Bazaar Shops. It's open 24/7 and it nothing is be huge. There Nothing is plenty of the seating. Most Extremely of the of seating is has stadium types seating seats with benches. They Have also will have never an out door patio. The staff is very friendly and attentive to the guests. I do notice that they are under staffed sometimes when they are busy. They get your drinks out pretty fast though. Also, this location is not owned by the casino so they don't charge outrageous prices like the location on the Linq promenade does. Definitely one of my favorite Starbucks stores. Stop by if your on the Strip.

Method: Gradient method[18]. Origin: 100% Positive. ADV:  $1-6.9e^{-12}$  Positive This It Starbucks Mcdonalds location is located in the Bally's Grand Bazaar Shops. It's open 24/7 and it is huge. There is plenty of seating. Most Many of the the seating is stadium type seating with benches. They also have 've an out up door patio. The staff is very friendly and attentive to the guests. I do notice that they are under staffed sometimes when they are busy. They getting your drinks out pretty fast though. Also, this location is not owned by the of casino so too they don't charge outrageous prices like think the location on the of Linq promenade seafront does. Definitely one of my favorite Starbucks stores. Stop by if unless your on the Strip.

## C.6 Empirical Example 6: Task - Sentiment analysis. Classifier - LSTM.

Method: Ours. Origin: 100% Positive. ADV: 93% Negative

I suppose I should write a review here since my little Noodle-oo is currently serving as their spokes dog in the photos. We both love Scooby Do's. They treat my little butt-faced dog like a prince and are receptive to correcting anything about the cut that I perceive as being weird. Like that funny poofy pompadour. Mohawk it out, yo. Done. In like five seconds my little man was looking fabulous and bad ass. Not something easily accomplished with a prancing pup that literally chases butterflies through tall grasses. (He ended up looking like a little lamb as the cut grew out too. So adorable.) The shampoo they use here is also amazing. Noodles usually smells like tacos (a combination of beef stank and corn chips) but after getting back from the Do's, he smelled like Christmas morning! Sugar and spice and everything nice instead of frogs and snails and puppy dog tails. He's got some gender identity issues to deal with. The pricing is also cheaper than some of the big name conglomerates out there The price is cheaper than some of the big names below. I'm talking to you Petsmart! I've taken my other pup to Smelly Dog before, but unless I need dog sitting play time after the cut, I'll go with Scooby's. They genuinely seem to like my little Noodle monster.

Method: Greedy[19]. Origin: 100% Positive. ADV: 88% Negative

I suppose I should write a review here since my little Noodle-oo is currently serving as their spokes dog in the photos. We both love Scooby Do's. They treat my little butt-faced dog like a prince and are receptive to correcting anything about the cut that I perceive as being weird. Like that funny humorous poofy pompadour. Mohawk it out, yo. Done. In like five seconds my little woman was looking fabulous and bad ass. Not something easily accomplished with a prancing pup that literally chases butterflies through tall grasses. (He ended up looking like a little lamb as the cut grew out too. So adorable.) The shampoo they use here is also amazing. Noodles usually smells like tacos (a combination of between beef stank and corn chips) but after getting back from the Do's, he smelled like Christmas morning! Sugar and spice and everything nice instead of frogs and snails and puppy dog tails. He's got some gender identity issues to deal with. The pricing is also cheaper than some of the big name conglomerates out there. I'm talking to you Petsmart! I've taken my other pup to Smelly Dog before, but unless I need dog sitting play time after the cut, I'll go with Scooby's. They genuinely seem to like my little Noodle monster.

**Method**: Gradient method [18]. **Origin**: 100% Positive. **ADV**: 93% Negative

I suppose I should write Write a review here since my little Noodle-oo is currently serving as their spokes dog in the photos. We both love Scooby Do's. They treat cure my little butt-faced dog like another prince knight and but are receptive to correcting anything about the cut that I perceive as that being weird. Like that funny poofy pompadour. Mohawk it out, yo. Done. In like five eleven seconds secs my little man was

looking fabulous and bad ass. Not something easily readily accomplished with a prancing strutting pup that literally chases butterflies through tall grasses. (He ended up looking like a little lamb beef as the The cut grew out too. So adorable.) The shampoo they use here is also amazing. Noodles usually smells like taeos quesadillas (a combination of beef stank and corn chips) but after getting back from the Do's, he smelled like Christmas morning! Sugar and spice cumin and everything nice instead of frogs and snails and puppy dog tails. He's got some those gender sexuality identity issues difficulties to deal contract with. The pricing is also cheaper than some of the big huge name conglomerates out there. I'm talking to you Petsmart! I've taken brought my other pup to Smelly Dog before, but unless I need dog sitting play time after the cut, I'll go with Scooby's. They genuinely nonetheless seem to like my little Noodle monster.