

Exploratory Character Edit Method - A Black-box Textual Adversarial Attack on Machine Learning Models in Web Security

Abstract—Researchers are making efforts on detection and prevention in web security by deploying machine learning models, but such models are vulnerable to adversarial examples. It is unclear how much risk adversarial perturbation counts for the safety of the detection of web application. In most cases, attackers generate such adversarial examples by assuming knowledge of detailed model information. Moreover, textual adversarial attacks have not been fully researched. To address these issues, we introduce Exploratory Character Edit Method (ECEM), a textual adversarial attack requiring only detection labels, which is applicable to real-world black-box detection models, such as intrusion detection systems. It begins with a large adversarial perturbation and seeks to minimize the perturbation while staying adversarial. Experiments on various machine learning models with three different data sets show that this adversarial attack outperforms the Fast Gradient Sign Method (FGSM) and the Boundary Attack while being natural, legible to humans, and effective in evaluation and analyzing black-box classifiers. ECEM in particular and adversarial attacks in general not only pave a novel way in web security research, but also are promising to further research in the safety of machine learning models with textual inputs.

Index Terms—adversarial example, machine learning, black-box, web security, intrusion detection

I. INTRODUCTION

Instead of human reasoning, machine learning techniques become more and more popular in building smart systems of detection and prevention models in web security. However, they have been shown vulnerable to adversarial examples. Such examples are maliciously perturbed examples much similar to original examples in human perception, but can mislead models to make false decisions [1]. The vulnerability of machine learning models to adversarial examples implies a security risk in applications with real-world consequences, such as intrusion detection system, web application firewall and comment toxicity analyzer. The study of adversarial examples is thus necessary to ensure the security of current machine learning algorithms, to provide metrics for robustness, to investigate the potential risk, and to suggest methodology to improve the robustness of models [2].

Image adversarial attacks have been largely studied in computer visions. On the contrary, research in textual adversarial attacks in web security or in natural language processing is insufficient. The concept of the adversarial attack originally comes from models in computer vision [1]. A set of methods of generating adversarial examples to enable attacks have been proposed, such as FGSM [3] and Boundary Attack

[4]. However, applying image adversarial attacks to models for textual inputs requires an extra access to the embedding layer, which is impractical. Studies on textual adversarial attacks are an urgent necessity for countless applications. Major difficulties of textual adversarial attacks are reflected on three inherent proprieties: discreteness, perceivability and semantics. Textual inputs are discrete features while image inputs are continuous. Little modification of pixels is hard for a human to distinguish but changes of characters or words are obvious. This is also the reason of hardly changed semantics of images and possible change of semantics of textual inputs. Methods for textual adversarial attacks mainly entail searching in underlying semantic space, character flipping and word substitution.

Important criteria for an adversarial attack are the similarity metric, the attack goal and the threat model. There are three widely-used similarity metrics to quantify similarity between perturbed and original images, all of which are L_p distance, with $p = \infty, 2, 0$. They are considered today as an approximation of human perceptual distance. For text corpus, string similarity and Levenshtein Distance are two common metrics to assess the coherence between adversarial and original examples. Informally, the Levenshtein Distance between two words is the minimum number of single-character edits (insertions, deletions or substitutions) required to change one word into the other [5].

Generally speaking, the attack goal is either targeted or untargeted. An untargeted adversarial attack would be successful when the prediction of the input after perturbation is flipped to any one but the original class. On the other hand, a targeted adversarial attack has to ensure that the prediction class of the adversarial example is a specific target. Most of the adversarial attacks are capable of switching between two types of attacks via a change of loss function [6].

Threat models fall into two categories: white-box and black-box. Since most existing white-box attacks rely on detailed model information including the gradient of the loss. They are often referred as gradient-based attacks. Some examples are the Fast Gradient Sign Method (FGSM) [4], the Jacobian-based Saliency Map Attack (JSMA) [7] and the Carlini & Wagner attack (CW attack) [6]. A way to defend against gradient-based attacks is to mask the gradients, for example by adding non-differentiable elements.

Unlike the white-box setting, the black-box setting can only have access to the outputs of the victim model without

any knowledge of the model architecture or hyper-parameter information. Based on whether the attacker is able to obtain the full probability or the label of a given input, black-box attacks are further divided into score-based and decision-based. On a conceptual level, a score-based attack employs the probabilities to numerically estimate the gradient, including black-box variants of JSMA (Narodytska & Kasiviswanathan [8]) and of the Carlini & Wagner attack [9]. A decision-based attack is a direct attack that solely rely on the final decision of the model, such as Boundary Attack [4].

Another type of attack is the transfer-based attack, which demands no model information but requires information on the training data. The data are used to train a fully observable substitute model from which adversarial perturbations can be synthesized with white-box attacks [10]. Transfer-based attacks rely on the empirical observation that adversarial examples often transfer between models.

This paper emphasizes the consequences of adversarial attacks and proposes a decision-based textual adversarial attack for machine learning models in web security. The generation of adversarial examples adopts a greedy search which starts with a large perturbation and approaches the original example as close as possible. The attack algorithm is referred as Exploratory Character Edit Method. Each iteration involves one character edit and the total number of edits is bounded by a small value. Figure 1 shows the overview of the ECEM.

ECEM can be used to generate either targeted or untargeted adversarial examples, depending on the attack’s willingness. Our experiments demonstrate that ECEM requires significantly fewer model queries than Boundary Attack and possesses the advantage of generating much more similar adversarial examples than FGSM. Furthermore, ECEM shows its universality to various machine learning models and different text corpus. However, Boundary Attack has to procure the values after the embedding layer and FGSM requires in addition the loss function being derivable. The fact that these two attacks are designed for images will cause unexpected results. In summary, contributions of this paper are:

- Consequences of adversarial examples are analyzed in real-world applications.
- ECEM, a black-box textual adversarial attack is proposed which requires no machine learning model structure and weights information.
- Performance of ECEM, FGSM and Boundary Attack is evaluated on various data sets.
- Practical applicability of ECEM to general web applications is shown.

The paper proceeds as follows. Section II presents the background of techniques in web security and problems of machine learning, followed by definitions of some basic knowledge. Section III sketches the problem description and an exhaustive method. Our approach, ECEM is presented in Section IV. In Section V, experiments of three adversarial attacks are evaluated in different conditions, followed by the case study in IDS. In Section VI we discuss the related work, and conclusion and future work are included in Section VII.

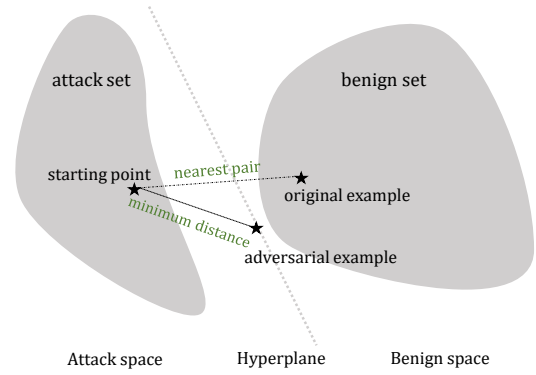


Fig. 1. Overview of Exploratory Character Edit Method. For an original example, a nearest starting point in the opposite set is selected and approaches to the original example as much as possible. The last example would be the nearest adversarial example as the output.

II. BACKGROUND

This section proceeds as follows. We firstly introduce the various attacks in web security. Then the concept and application of Intrusion Detection System (IDS) will be presented. In the third part, we will discuss the machine learning techniques and adversarial attacks. And finally, significant measurements in text corpus will be included.

A. Attacks in Web Security

The OWASP Top 10 is a standard awareness document for developers and web application security [11]. It represents a broad consensus about the most critical security risks to web applications. Top 10 web application security risks are as follows: injection, broken authentication, sensitive data exposure, XML external entities, broken access control, security misconfiguration, cross-site scripting, insecure deserialization, using components with known vulnerabilities and insufficient logging & monitoring. Here we introduce three of them to give a general idea. Readers who are interested in this topic could take advantage of the reference.

1) *Injection*: Injection flaws, such as SQL, NoSQL, OS, and LDAP injection, occur when untrusted data is sent to an interpreter as part of a command or query. The attacker’s hostile data can trick the interpreter into executing unintended commands or accessing data without proper authorization.

2) *Sensitive Data Exposure*: Many web applications and APIs do not properly protect sensitive data, such as financial, healthcare, and PII. Attackers may steal or modify such weakly protected data to conduct credit card fraud, identity theft, or other crimes. Sensitive data may be compromised without extra protection, such as encryption at rest or in transit, and requires special precautions when exchanged with the browser.

3) *Cross-Site Scripting (XSS)*: XSS flaws occur whenever an application includes untrusted data in a new web page without proper validation or escaping, or updates an existing web page with user-supplied data using a browser API that can create HTML or JavaScript. XSS allows attackers to execute

scripts in the victim's browser which can hijack user sessions, deface web sites, or redirect the user to malicious sites.

B. Intrusion Detection System

Web applications are becoming increasingly popular and complex in all sorts of environments, ranging from e-commerce applications to banking. As a consequence, web applications handle large amounts of sensitive data, which renders applications even more attractive to malicious users. The repercussion of many attacks might be devastating, like identity supplanting, sensitive data hijacking, access to unauthorized information, web page's content modification, command execution, etc. Therefore, it is fundamental to protect web applications and to adopt the suitable security methods.

A network intrusion attack can be any use of a network that compromises its stability or the security of information that is stored on computers connected to it. The goal of intrusion detection is to build a system which would automatically scan network activity and detect such intrusion attacks. Once an attack is detected, the system administrator is informed and can take corresponding action.

IDS could be categorized into Network Intrusion Detection System (NIDS) and Host-based Intrusion Detection Systems (HIDS). NIDS focuses on monitoring network traffic to identify attacks and malicious activities, while HIDS monitors critical operating system files for suspicious activities.

Two main methods exist for IDS detection process: signature-based detection and anomaly-based detection. Signature-based detection distinguishes patterns of attacks and malicious activities, and supervises anything else that does not match the patterns as well. Anomaly-based detection, on the other hand, is based on established models of benign behaviors, and regards any activity deviating from the models as malevolent.

Since recognized attacks and malicious activities can be captured as patterns, the signature-based IDS has high accuracy and lower rate of missing reports for known attacks. Nevertheless, for new or unknown attacks, also named Zero-day Attacks, the signature-based approach suffers often from low accuracy and high rate of missing reports. Thus, how to continuously refine rules that keep up with the attack patterns is a critical step towards improving the performance of IDS.

Anomaly-based techniques detect unknown attacks by learning the existing patterns for benign access. As any activity that does not match the existing benign patterns would be estimated as anomaly, regrettably, performance and accuracy suffer from high false positive rate.

C. Machine Learning Techniques and Adversarial Attack in Web Security

Machine learning is the ability to automate solving problems and tasks by learning immanent patterns from large amounts of data. In order to achieve this goal, two methods of learning are used: supervised learning (including classification, support vector machines, neural networks) and unsupervised learning (including clustering, dimensionality reduction, recommender

systems, distance, and normalization). Machine learning techniques are also deployed to solve web security problems. Attack signatures could be learned by well-configured machine learning models. Specifically for an IDS, detection models based on machine learning could be either anomaly-based or signature-based, or a mixture which combines advantages of two modes [12].

The existence of adversarial examples arises much concern on the safety of deployed machine learning models. Adversarial examples aim at causing target model to make mistakes on prediction. An adversarial attack could be effective to both online and offline learning models. No matter it is an intentional or unintentional adversarial attack, evaluating adversarial examples has become a trend of building a robust machine learning model and understanding their shortcomings.

In machine learning, systems which employ offline learning alter no longer their prediction function when the initial training phase has been completed. In web security, injecting tons of adversarial examples into such systems may significantly increase false positive rate and drop the prediction accuracy. The worst case is for example that some malicious accesses finally disguise themselves and successfully cheat the intrusion detection systems. In [13], a backdoor could be implemented for neural networks to make them perform well during training step while poorly during evaluation step.

On the other hand, online learning is a method in which the prediction model is trained by some data and should be updated to a better predictor with more data in the future. By flooding with adversarial examples, machine learning models could be fundamentally changed or dysfunctional. Figure 2 demonstrates how a machine learning model changes its prediction function after learning new data. In [10], the error rates of systems from Amazon and Google are respectively up to 96.19% and 88.94% after adversarial attack. Also in 2016, Microsoft's AI chatbot, Tay, became extremely racist in less than 24 hours after learning comments on Twitter. [14]

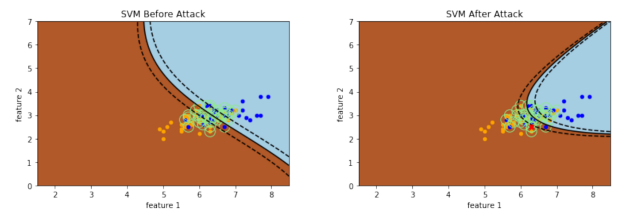


Fig. 2. Fundamental Change of Prediction Function after Learning New Data. Online learning is supposed to eliminate bias and overfitting but models could also be poisoned by malicious intention.

D. Measurements in Text Corpus

In linguistics, a text corpus is a large and structured set of texts, and in machine learning, it is the equivalent of "dataset". To depict a text corpus, techniques for measuring the coherence of textual elements are indispensable. We present two essential measurements which calculate the similarity and distance between two inputs. Both measurements illustrate

the concept of “closeness”, based on which the quality of adversarial examples would be evaluated.

1) *Similarity Metrics*: Out of generality, words, texts, logs, URLs and all the textual inputs could be considered as strings in computer science. We have applied Ratcliff-Obershelp pattern recognition for strings [15] as an essential similarity metric throughout the paper, explained as follows.

Definition 1. Ratcliff-Obershelp pattern recognition: Compute the similarity of two strings as the number of matching characters divided by the total number of characters in the two strings. Matching characters are those in the longest common subsequence plus, recursively, matching characters in the unmatched region on either side of the longest common subsequence. [15]

Note $\rho(\cdot, \cdot)$ the similarity between 2 strings, $|\cdot|$ the number of characters of a string, $\cdot \cap \cdot$ the matching characters of two strings, and S_a and S_b as two strings. The formula is shown as equation (1)

$$\rho(S_a, S_b) = \frac{2 \times |S_a \cap S_b|}{|S_a| + |S_b|} \quad (1)$$

2) *Levenshtein Distance*: Levenshtein Distance is a string metric for measuring the difference between two sequences, named after the Soviet mathematician Vladimir Levenshtein [5]. The definition is as follows.

Definition 2. Levenshtein Distance: The Levenshtein Distance between two words is the minimum number of single-character edits (insertions, deletions or substitutions) required to change one word into the other.

Given two strings S_a and S_b , denote $\psi_{S_a, S_b}(i, j)$ as the Levenshtein Distance between the first i characters in S_a and the first j characters of S_b , thus Levenshtein Distance equals $\psi_{S_a, S_b}(|S_a|, |S_b|)$, where $|\cdot|$ is the length of a string. Its formula is defined in equation (2)

$$\psi_{S_a, S_b}(i, j) = \begin{cases} \max(i, j) & , \text{if } \min(i, j) = 0 \\ \min \begin{cases} \psi_{S_a, S_b}(i-1, j) + 1 \\ \psi_{S_a, S_b}(i, j-1) + 1 \\ \psi_{S_a, S_b}(i-1, j-1) + \mathbb{1}_{S_a[i] \neq S_b[j]} \end{cases} & , \text{otherwise} \end{cases} \quad (2)$$

where $S_a[i]$ represents i -th character in S_a and $\mathbb{1}_{S_a[i] \neq S_b[j]}$ is an indicator function such that it is 1 if $S_a[i] \neq S_b[j]$, otherwise it is 0.

III. PROBLEM DESCRIPTION AND THE EXHAUSTIVE METHOD

In this section, we are going to frame the generation of an adversarial example as an optimization problem, and then propose an exhaustive method to solve it. We prove that the complexity of the method reach unfortunately to $O(n!)$. In the next section, we are going to introduce ECEM as an optimized solution that reduces the complexity to $O(n)$. Before presenting the problem description, it is appropriate to stipulate different elements in the framework.

A. Element Stipulation

1) *Machine Learning Model*: We assume a classification model including a prediction function F , accepting an input S as a string and producing an output y which is treated as a probability distribution on classes. The input S will be embedded into a numeric expression x since textual inputs could not be directly used for numerical calculation and the model produces in addition a classifier C that assigns an input S with the class of maximum probability: $C(S) = \arg \max_c (F(S))_c$. Since we are pursuing a decision-based attack, only the classifier C is accessible, but not the underlying model F .

2) *Training Data*: Training data is crucial for a machine learning model. Without knowledge of training data, a synthetic training data set could be generated using a random feature vector method [16] or more sophisticated techniques, such as hill-climbing [17]. In addition, [18] and [19] have shown the possibility of extracting the training data from a machine learning model. With the training data set, attackers could benefit from more information for their attack techniques.

3) *Starting Point*: The strategy for choosing a starting point is essential. As mentioned in Section I, the attacked goal could be either targeted or untargeted. As an exhaustive method, we are going to choose a random starting point solely satisfying this criteria. For a targeted attack, the starting point is supposed to be chosen as the target class and stays the same during the iteration, while for an untargeted setting, choose any starting point whose label is different from the original one and it is apparently unnecessary to stay the same as long as it is not the original class. Since the starting point is randomly chosen, it could be completely “unsimilar” to the original example and the Levenshtein Distance could be as long as the larger length of two textual inputs. If the original example is extremely long and complex, like a comment of a hundred words, and if the starting point is not wisely chosen, the distance could be several thousand.

4) *Character Edit*: After choosing a starting point, the Levenshtein Distance n to the original example could be calculated and we could know the minimum set of character edits for transforming the starting point to the original example. As shown in Figure 3, such operations for transforming String A to String B are as follows, where the index begins with 0:

- Replace ‘b’ with ‘u’ at 1st position
- Delete ‘d’ at 3rd position
- Insert ‘l’ at 6th position
- Insert ‘q’ at 10th position

```
String A: abcdSSmnop
String B: aucSSlmnopq
          ↓
String A: abcdSS mnop
String B: auc SSlmnopq
```

Fig. 3. Minimum Operations for Transforming String A to String B. Red parts are matching characters in two strings while black parts and spaces indicate required character edits.

B. Optimization Framework

Note the machine learning model as a classifier C . The generation of an adversarial example could be formulated as an optimization problem as follows:

$$\begin{aligned} \min_{S'} d(S, S') \\ \text{s.t. } \exists (S_i)_n \in \mathbb{S}, (\phi_i)_{n-1} \in \Phi, \\ \forall i \in [0, n-1], S_{i+1} = \phi_i(S_i), \\ \forall (i, j) \in [0, n]^2, C(S_i) = C(S_j), C(S_i) \neq C(S), \\ \text{and } S_n = S' \end{aligned} \quad (3)$$

where S is the original string, S' is the adversarial example, S_0 is the chosen starting point, \mathbb{S} is the space of all strings, and Φ is the set of character operation – i.e., insertion, deletion and substitution of a character at a given position. d is the Levenshtein Distance function that $d(S_a, S_b) = \psi_{S_a, S_b}(|S_a|, |S_b|)$, mentioned in Section II-D2. A targeted setting needs to add an extra condition on the class of the adversarial example.

C. Exhaustive Method

A nearest adversarial example is pursued by minimizing the Levenshtein Distance. The exhaustive method would be evaluating all the candidates of adversarial examples and choose the best one. Concrete instructions of a targeted attack would be presented as follows:

- 1 Pick an original example.
- 2 Randomly choose a starting point of the targeted class.
- 3 Calculate the difference between the starting point and the original example.
- 4 Try all the combinations of character edits and calculate the candidates for potential adversarial examples.
- 5 Choose the best one among candidates as the final adversarial example.

At each step, a character edit is chosen from a proposal set, which reduces the Levenshtein Distance of the adversarial towards the original example. As an exhaustive method, all the combinations of possible numbers of character edits ranging from 1 to $n-1$ have to be evaluated. The complexity reaches thus to $O(n!)$. Taking the example presented in the previous part, the Levenshtein Distance between the starting point and the original example could be several thousand. We know that $1000! \approx 10^{249}$, which takes too much time in practical. Hence the choice of the starting point is essential and we need a “smart” search strategy to accelerate the calculation process. In the Section IV, we introduce ECEM which reduces the complexity significantly to $O(n)$. Figure 4 demonstrates four possible candidates for potential adversarial examples. We can choose the best one among these candidates.

IV. EXPLORATORY CHARACTER EDIT METHOD

In Section III, we have presented the problem description and proposed an exhaustive method. The strategy for wisely choosing a starting point and the searching tactic are crucial to the complexity. In this section, we are going to introduce Exploratory Character Edit Method as a generative algorithm

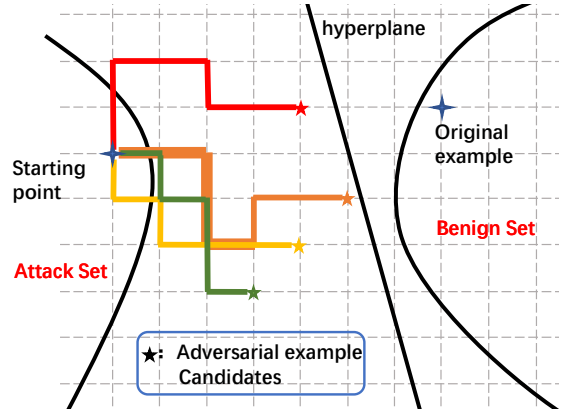


Fig. 4. Candidates for Adversarial Examples.

A starting point has several potential candidates for adversarial examples.

for textual adversarial examples. The following content develops as follows: the strategy for choosing the starting point, the search tactic, the detailed algorithm description and the proof of increasing improvement of adversarial quality.

A. Strategy for Choosing Starting Point

The idea for wisely choosing a starting point is trying to minimize the Levenshtein Distance. The distance determines the value of n . As always, the attack goal needs to be satisfied. For a targeted adversarial attack, choose an example with the target class as the starting point, while for an untargeted attack, choose the example with any class other than that of the original example.

In case of generating a single adversarial example for an original example, human wisdom could be helpful. Based on the inherent information of the original example, it is not difficult to choose a starting point in another class. After several times of attempts, a sufficiently “close” starting point could be chosen for the generation of the adversarial example.

A stronger strategy is still requisite to cope with large amounts of text corpus, and a dictionary of data for different classes would be helpful. The generation of the dictionary could rely on the training data set, obtained either by training data extraction or a synthetic generation as mentioned in Section III-A2. Using the training data set, we could separate data by their classes. The dictionary could be improved with more and more new samples. In accord with the attack goal, we choose the starting point in a class with the minimum Levenshtein Distance to the original example. The larger the dictionary is, the “closer” the starting point could be.

B. Searching Tactic

In the exhaustive method, all the potential candidates for the adversarial examples have been calculated and the best one among them would be chosen as the output. The idea to reduce the complexity of the algorithm is to choose randomly one candidate as the final adversarial example. In fact, the calculation of other candidates are unnecessary. Figure 5 shows the greedy searching of an adversarial example.

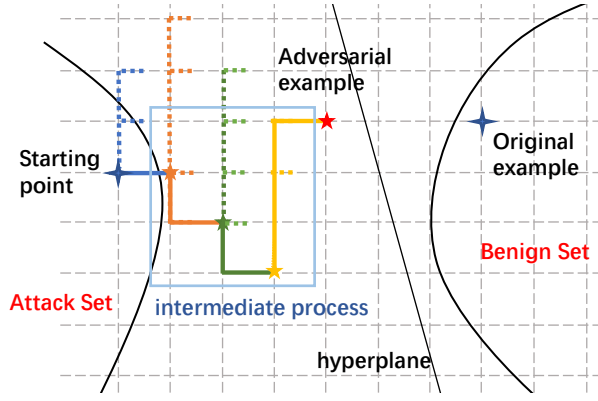


Fig. 5. Greedy Searching for an Adversarial Example. Solid lines indicate the path of greedy search and dotted lines show other possibilities in each step.

Instead of searching all the combinations of character edits, a greedy searching setting is adopted. As a rejective sampling, during each iteration, an operation would be chosen randomly from the possible operation set according to the Levenshtein Distance. If this operation could not make the example stay adversarial, drop it and choose another one. Once this character edit is done, the operation would be removed from the set. The searching ends when none of the operations in the set could keep the class label adversarial. Hence the approximately nearest adversarial example would be the last example. With this greedy searching, we trade off a little bit of “quality” for speed. Experiments further demonstrate that competitive candidates in the exhaustive method have in most cases less than five character differences, so the greedy searching is practical. Furthermore, since the greedy searching follows a random walk in the text corpus, the attack could launch another generation if the current adversarial example is not ideal.

Each iteration reduces the Levenshtein Distance by 1 and when the distance is large, the number of finding a feasible edit would be exactly 1 because the example is surrounded by other examples with the same class label and it never crosses the hyperplane, which would be further supported in Section V. The greedy setting successfully reduces the calculation complexity from $O(n!)$ to $O(n)$, with n as the distance between the starting point and the original example.

C. Generating Adversarial Examples

ECM begins with a large adversarial perturbation and seeks to minimize the perturbation while staying adversarial – i.e., the label should never be the original class. As mentioned before, the attack goal could be either targeted or untargeted. We should choose a corresponding starting point according to the previous strategy. The class label should always be maintained during the searching for a targeted attack. On the contrary for an untargeted setting, the class label could be any class other than that of the original example. The class label is not obliged to be maintained the same all the time.

The detailed targeted adversarial attack is presented in Algorithm 1, while for the untargeted setting, there would be only a tiny change of Line 12 to “if $C(S_{new}) \neq L_S$ then”.

Algorithm 1: Exploratory Character Edit Method

Input: Starting point S_0 ; Original string S ; Classifier C of the victim model

Output: Approximately nearest adversarial example S'

- 1: Let L_S be the class label of S passing C
 - 2: Let $L_{S'}$ be the class label of S' passing C
 - 3: Verify $L_S \neq L_{S'}$
 - 4: Let d be the Levenshtein Distance function
 - 5: $N \leftarrow d(S, S_0)$
 - 6: Generate Φ as the operation set of minimum required number of operations for transforming S_0 to S
 - 7: **for** $i = 0$ to $N - 1$ **do**
 - 8: Shuffle Φ
 - 9: $S_{i+1} \leftarrow None$
 - 10: **for** ϕ in Φ **do**
 - 11: $S_{new} \leftarrow \phi(S_i)$
 - 12: **if** $C(S_{new}) = L_{S'}$ **then**
 - 13: $S_{i+1} = S_{new}$
 - 14: Remove ϕ from Φ
 - 15: **Break**
 - 16: **end if**
 - 17: **end for**
 - 18: **if** $S_{i+1} = None$ **then**
 - 19: $S' = S_i$
 - 20: **Break**
 - 21: **end if**
 - 22: **end for**
 - 23: **return** S'
-

D. Improvement of Adversarial Quality

The adversarial quality could be measured by similarity and distance. Each iteration decreases the Levenshtein Distance by 1. We are going to prove that such iteration also increases the string similarity, no matter an insertion, a deletion or a substitution of a character. Recall the string similarity defined in Section II-D1 and suppose the similarity before the character edit equals:

$$\rho(S, S_i) = \frac{2 \times |S \cap S_i|}{|S| + |S_i|}$$

Insertion: This edit would insert a character from S into S_i , so the length of S_i and the number of matching characters would both increase 1.

$$\rho(S, S_{i+1}) = \frac{2 \times (|S \cap S_i| + 1)}{|S| + |S_i| + 1} \geq \rho(S, S_i)$$

Deletion: This edit would delete an unnecessary character of S_i , so the length of S_i would decrease 1 while the number of matching characters remains the same.

$$\rho(S, S_{i+1}) = \frac{2 \times |S \cap S_i|}{|S| + |S_i| - 1} \geq \rho(S, S_i)$$

Substitution: This edit would replace an error character of S_i by the exact character of S at the same position, so the length of S_i remains the same while the number of matching characters would increase 1.

$$\rho(S, S_{i+1}) = \frac{2 \times (|S \cap S_i| + 1)}{|S| + |S_i|} \geq \rho(S, S_i)$$

V. EXPERIMENT AND EVALUATION

In this section, we compare three adversarial attacks: ECEM, FGSM and Boundary attacks on machine learning based classification models simulated with three data sets. We firstly evaluate the quality of adversarial examples, then we show the efficiency and effectiveness of the three attacks. Furthermore, a close look on the impact of adversarial URL requests generated by ECEM on a machine learning based IDS is discussed.

A. Data Sets

Three data sets are explored for our experiment: HTTP data set CSIC 2010, Malicious & Non-Malicious URL and Toxic Comment Classification.

The HTTP data set CSIC 2010 includes web requests, widely used to test IDS systems for web applications. It was developed at the "Information Security Institute" of CSIC (Spanish Research National Council) [20]. The CSIC 2010 dataset contains 72000 normal requests and more than 25000 anomalous requests. Each record contains the following features: URL, User-Agent, Pragma, Cache-control, Accept, Accept-Encoding, Accept-Charset, Accept-Language, Host, Cookie, and Connection.

Malicious & Non-Malicious URL contains 411247 URLs with labels indicating whether they are malicious for accessing. 82% records are benign and the rest are malevolent. This is an open data set from Kaggle.com [21]. Each record contains only two features: URL and the label.

Toxic Comment Classification is a data set of comments from Wikipedia's talk page edits. These comments contains different types of toxicity like threats, obscenity, insults, and identity-based hate. Building accurate detection models based on this data set is a challenge from Kaggle.com [22]. The training data set has 159571 records and the test data set has 153164 records. Each record contains the comment and following labels: toxic, severe_toxic, obscene, threat, insult, identity_hate.

B. Adversarial Attacks on a CNN-based model

In this part, ECEM, FGSM and Boundary Attack are used to attack a CNN-based detection model [23]. HTTP data set CSIC 2010 is used for model training, testing and validation. The resulting adversarial examples are evaluated according to the following criteria:

- String similarity
- Levenshtein Distance
- Score
- Success rate

We have generated adversarial examples for those whose original labels are benign, which could be used to fool the model and poison an online-learning model. In this case, the false positive rate would be largely increased.

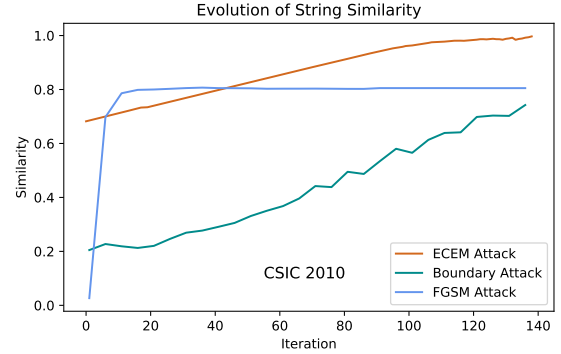


Fig. 6. Average String Similarity between Adversarial and Original Example. Both ECEM and Boundary Attack approach steadily to 1, while a rapid increase can be found with FGSM, but its similarity converges to 0.82. ECEM starts at a much higher value that reaches 1 with much fewer iterations.

The evolution of average string similarity between adversarial and original example is presented in Figure 6. We can observe from the figure that all three attacks start from a low value and increase with the iteration. Since ECEM chooses a nearest example in the "attack" set as the starting point, the string similarity begins with a higher value than FGSM and Boundary Attack. String similarity of ECEM converges to 1, which can be reached only when two strings are exactly the same, while FGSM is bounded by approximately 0.82. The advantage of FGSM is that within 20 iterations, the string similarity is able to achieve a surprisingly high value, but further effort is worthless to render it any larger. Boundary Attack also increases the string similarity step by step, but apparently the query number is much larger than ECEM to achieve the same quality of result.

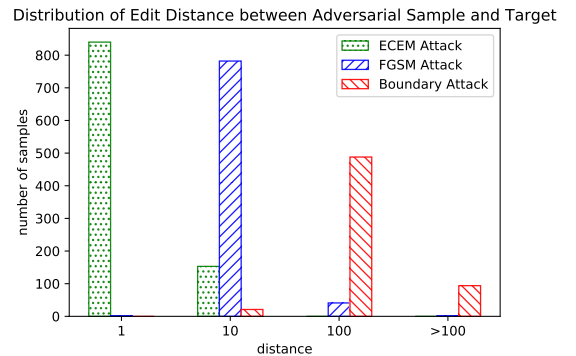


Fig. 7. Distribution of Levenshtein Distance between the Final Adversarial and Original Example.

Figure 7 further validates that ECEM outperforms FGSM and Boundary Attack. Here we presents the distribution of Levenshtein Distance between the best adversarial and original

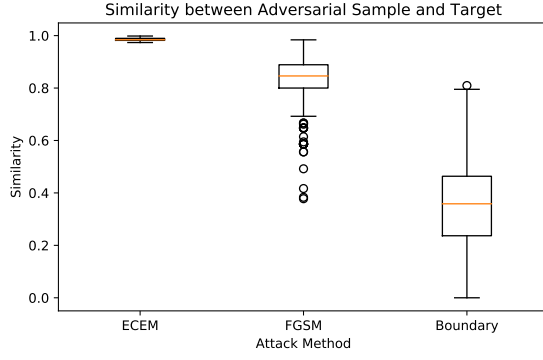


Fig. 8. Distribution of String Similarity between Final Adversarial and Original Example.

example – e.g., the number of iterations could be as much as the attack wants to generate the most similar adversarial example. Surprisingly, almost all the adversarial examples of ECEM have at most 10 characters different than their original examples, and a large part of them have only one character difference. FGSM can achieve a similar performance, but unfortunately, few adversarial examples could have only one character difference. Character difference is crucial for textual inputs since the less it is, the more possible the semantic meaning remains unchanged. Figure 8 shows that from the string similarity perspective, the variance of ECEM is much smaller than FGSM and Boundary Attack.

The prediction of a class is always related to the score of the last layer of the machine learning model. A classifier often calculates the probability of the different classes based on the score. For a binary classifier for example, the prediction will be an “attack” label only if the probability is higher than 50%.

ECEM is a decision-based adversarial attack which has access only to the result of the classification. For the research, we are interested in how the probability of the “attack” label evolves during the iteration. To be clear, this probability is not used for the generation of adversarial example.

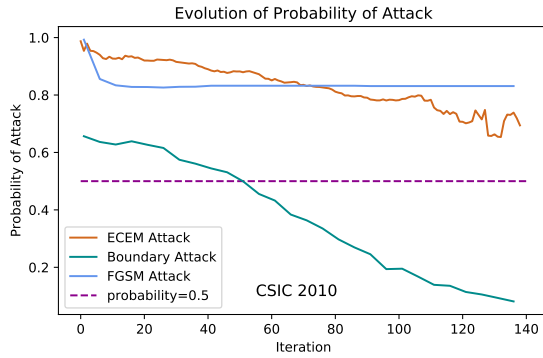


Fig. 9. Average Probability of “Attack” Label. Both ECEM and FGSM stop above the delimiter 50%, while a rapid decrease could be found with FGSM, but its probability converges to 0.8. Boundary Attack crosses the delimiter, indicating a failed adversarial attack.

The evolution of average probability of “attack” label is presented in Figure 9. We can observe from the graph that all three attacks start from a large value and decrease with the iteration, because they all start from an example labelled “attack”. Since the adversarial example is situated much closely to the hyper-plane of the prediction function, the probability of the “attack” label should be also close to 50%. Attack-label probability of ECEM indeed converges to 0.5, according to the figure, while FGSM is bounded by approximately 0.8. Again, the advantage of FGSM is that within 20 iterations, the probability is able to achieve a low value, but further effort is worthless to render it any lower. The deadly problem of Boundary Attack is that with further iteration, the example stays no longer adversarial.

Another important factor for an adversarial attack is the success rate. As we mentioned before, the boundary attack is not capable of assuring the success of generating an adversarial example for a text with the evolution of the iteration. Table I shows the success rates of ECEM and Boundary Attack launched on 1000 samples and 10000 samples respectively with 25, 50 and 75 iterations. The success rate of ECEM is always assured as 100%.

TABLE I
SUCCESS RATE

# iterations	1000 samples		10000 samples	
	B-A ^a	ECEM	B-A	ECEM
25	60.73%	100%	60.40%	100%
50	50.25%	100%	50.67%	100%
75	33.23%	100%	32.55%	100%

- ^aBoundary Attack

In this part, we have compared ECEM with FGSM (a white-box setting) and Boundary Attack (a black-box setting). The result of the experiment has shown following advantages:

- low query times requirement
- high similarity between original and adversarial examples
- 100% success rate
- no hyper-parameter tuning
- no knowledge of model architecture

C. Validation of ECEM on More Models

In order to further validate the performance of ECEM, we test it on three more machine learning models, namely, logistic regression, linear support vector machine and random forest. The models on CSIC achieve the following accuracy: LC 92.66%, Linear SVC 92.85% and RF 98.47%.

The similarity curve and the probability curve for four detection methods on HTTP data set CSIC 2010 are given in Figure 10 and 11. We observe that our method consistently performs well under different machine learning detection models. The string similarity finally converges to 1, and the probability of “attack” label decreases until a value above 50%.

Although these models have a good performance on predicting attacks, LC and linear SVC have some tiny fluctuation

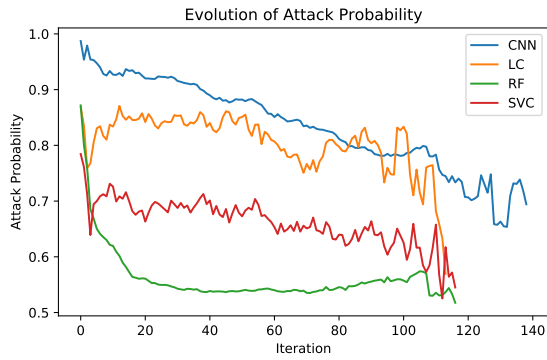


Fig. 10. Evolution of Probability of “Attack” Label with Machine Learning Models. Models include CNN-based, logistic regression, random forest and linear support vector, attacked by ECEM. All line start form a high value and goes down to 0.5.

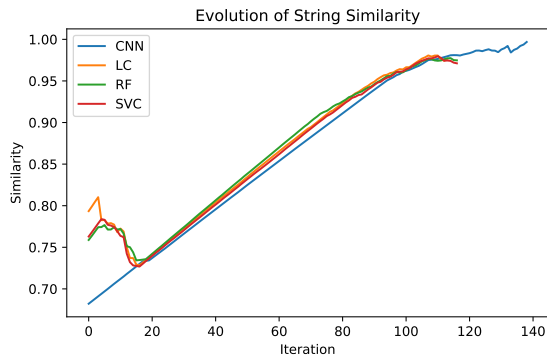


Fig. 11. String Similarity between Adversarial and Original Examples with Machine Learning Models. Models include CNN-based, logistic regression, random forest and linear support vector, attacked by ECEM. All lines start from a high value and approach to 1.

on predicting attack probability. Besides, RF decreases the probability rapidly within several iterations. Hence, we hypothesize that the performance may be influenced by models’ architectures.

In this part, we have validated ECEM on different victim models. The result of the experiment has shown following advantages:

- universally applicable
- high similarity between original and adversarial examples

D. Case Study with IDS

In web security, IDS is an effective detection mechanism for most applications. As machine learning techniques are more and more utilized to replace traditional IDS, in this part, we are going to apply our ECEM to machine learning models specially configured for IDS. Adversarial examples for attacks in web security are demonstrated as follows.

As mentioned in Section II-A, sensitive data access and injection are ranked in top 10 risks by OWASP Foundation. Figure 12 demonstrates a case study of critical file access. The

Sensitive Data Access

Starting point:
/tienda1/publico/productos.jsp.old → An old version file, dangerous!

Adversarial example:
/tienda1/publico/productos.jsp.d → No longer accessible, benign, but predicted dangerous.

Original example:
/tienda1/publico/productos.jsp → Current version, benign.

Fig. 12. An Adversarial Example in Sensitive Data Access.

original example is an access to the current version. It is classified as benign by IDS. We show that the adversarial example has only two-character difference from the original example. And its starting point of ECEM is a URL attempting to access an old-version file. Such behavior poses a considerable risk. At the same time the adversarial example should be considered benign, because a request with its URL can no longer reach the old-version file. On the contrary, it is predicted dangerous by IDS.

Injection

Starting point:
/tienda1/publico/vaciar.jsp?B2=Va
ciar+carrito%7C → Injection, malicious.

Adversarial example:
/tienda1/publico/vaciar.jsp?B2=Va
ciar+carrito7C → No longer injection, but predicted malicious.

Original example:
/tienda1/publico/vaciar.jsp?B2=Va
ciar+carrito → Normal access, benign

Fig. 13. An Adversarial Example of Injection.

Figure 13 presents us a case study of injection attack. The starting point is a typical injection attack since the underlining part “%7C” is the special character “|” in URL. The original example is a normal access. The adversarial example eliminates the character “%”, so the combination “%7C” represents no longer the special character “|”. Thus the adversarial example is not an injection attack and should be considered benign. Nevertheless, it is predicted malicious by the IDS.

The two case studies are conducted on data access and injection attacks. By generating large amounts of adversarial examples of different kinds, the existing false positive rate of IDS would be increased in large scale. By injecting such records, an online learning IDS could be poisoned as a fundamental change of prediction function.

VI. RELATED WORK

Recent research has shown interest in the deployment of machine learning techniques in web security problems. [24]–[26] Islam et al. [27] investigated the possibilities of modeling spammer behavioral patterns by Naive Bayesian classifier, Decision Tree Induction and Support Vector Machines. Calzavara et al. [28] proposed a detection method based on supervised

learning for web authentication. Clincy et al. [29] leveraged Genetic Algorithm to generate attack signatures, building Intrusion Detection System for web services.

The existence of adversarial examples for neural networks (Szegedy et al. [1]; Biggio et al. [30]) was initially a theoretical concern. Recent research has demonstrated the applicability of adversarial examples in real world. Adversarial examples on a printed page remain adversarial when captured using a cell phone camera in an approximately axis-aligned setting. [31] The research of adversarial examples continues. Engstrom et al. argue that the existence of adversarial examples is tied to the unrobust features. [32]

There are numbers of image adversarial attacks. Methods for transforming ordinary two-dimensional images into adversarial examples, including techniques such as the L-BFGS attack (Szegedy et al. [1]), FGSM (Goodfellow et al. [3]), and the CW attack (Carlini & Wagner [6]), are well-known. Later, Brendel et al. [4] introduced the first decision-based adversarial attack, Boundary Attack, which requires no information about model architecture or weights. More research for black-box adversarial attack is proposed such as HopSkipJumpAttack (Chen et al. [2]) and Spatical Attack (Engstrom et al. [33]) Adversarial examples generated by these techniques could be also transferred into real world. [31] Athalye et al. [34] also shows the possibility to render adversarial examples robust to noise in real world.

Compared with images, adversarial attack on text is more challenging, composed of discrete words, hard to adopt gradient-based methods to perturb it. Various methods are proposed to tackle the challenges, such as back-translation (Iyyer et al. [35]), searching in underlying semantic space (Zhao et al. [36]), character flipping (Ebrahimi et al. [37]) and word substitution (Ribeiro et al. [38]; Alzantot et al. [39]; Ren et al. [40]). Zang et al. [41] proposed an attack model which comprises a sememe-based word substitution strategy and the particle swarm optimization algorithm, to tackle the existing problems.

It has been shown that attack and defense often come hand-in-hand, and one's improvement depends on the other's progress. Defenses against adversarial examples could be roughly categorized into three types: detection-based defense, gradient and representation masking, and adversarial training. Detection-based approaches aim at differentiating an adversarial example from a set of benign examples using statistical tests or out-of-sample analysis. Recent works can be found in [42]–[47]. Papernot et al. [7] proposed a defense distillation with gradient masking. Representation masking aims to replace the internal representations in DNNs with robust representations to alleviate adversarial attacks. For example, Bradshaw et al. [48] proposed to integrate DNNs with Gaussian processes and RBF kernels for enhanced robustness. Adversarial training, also known as the data augmentation method, makes DNNs less sensitive to perturbations to the examples by jointly training with adversarial examples. Interested readers can refer to the recent works in [49]–[53] and the references therein.

VII. CONCLUSIONS AND FUTURE WORK

In this paper, we proposed a black-box textual adversarial example generation method, ECEM, to launch attacks against machine learning models in web security. This method requires no knowledge of model architecture or weights, nor the word or character embedding method. ECEM is proved to be more competitive than existing methods, such as FGSM and Boundary Attack. Our results show that application of image adversarial attack on text corpus does not provide ideal performance, and ECEM outperforms them in terms of similarity and iteration saving.

We foresee many avenues for further research in the domain of machine learning techniques in web security. Adversarial examples could be beneficial for studying properties of models. This textual adversarial attack could be transferred to numerous applications in network security: intrusion detection system, firewall and so on. In addition, it can also be applied to attacking other applications with textual inputs in natural language processing. Furthermore, defense mechanism could be another promising research direction for the robustness and safety of models.

Machine learning techniques facilitate common life in all kinds of areas, but we should be careful about the deployment of machine learning models since they suffer fatal vulnerabilities, including adversarial examples. The investigation on formal verification methods is imperative to ensure their general security when machine learning models are applied.

REFERENCES

- [1] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," dec 2013.
- [2] J. Chen, M. I. Jordan, and M. J. Wainwright, "HopSkipJumpAttack: A Query-Efficient Decision-Based Attack," pp. 1–23, 2019.
- [3] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, pp. 1–11, 2015.
- [4] W. Brendel, J. Rauber, and M. Bethge, "Decision-based adversarial attacks: Reliable attacks against black-box machine learning models," *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, pp. 1–12, 2018.
- [5] V. I. Levenshtein, "Binary Codes Capable of Correcting Deletions, Insertions and Reversals," *Soviet Physics Doklady*, vol. 10, p. 707, Feb 1966.
- [6] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 39–57, May 2017.
- [7] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," 2015.
- [8] N. Narodytska and S. P. Kasiviswanathan, "Simple black-box adversarial perturbations for deep networks," 2016.
- [9] 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," *AISec 2017 - Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, co-located with CCS 2017*, pp. 15–26, 2017.
- [10] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, "Practical black-box attacks against machine learning," in *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, ASIA CCS '17*, (New York, NY, USA), p. 506–519, Association for Computing Machinery, 2017.
- [11] D. Wichers and J. Williams, "Owasp top-10 2017," *OWASP Foundation*, 2017.

- [12] P. Y. Simard, "Using machine learning to break visual human interaction proofs (hips)," in *Advances in Neural Information Processing Systems 17, Neural Information Processing Systems (NIPS'2004)*, pp. 265–272, MIT Press, 2004.
- [13] T. Gu, B. Dolan-Gavitt, and S. Garg, "Badnets: Identifying vulnerabilities in the machine learning model supply chain," *arXiv preprint arXiv:1708.06733*, 2017.
- [14] G. C. Amy Tenner, "Microsoft's ai twitter bot goes dark after racist, sexist tweets," <https://www.reuters.com/article/us-microsoft-twitter-bot-idUSKCN0WQ2LA>, 2016.
- [15] J. Ratcliff and D. Metzner, "Ratcliff-overshelp pattern recognition," <https://xlinux.nist.gov/dads/HTML/ratcliffOvershelp.html>, 1998.
- [16] F. Tramèr, F. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart, "Stealing machine learning models via prediction apis," in *25th USENIX Security Symposium (USENIX Security 16)*, (Austin, TX), pp. 601–618, USENIX Association, Aug. 2016.
- [17] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, "Membership inference attacks against machine learning models," in *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 3–18, May 2017.
- [18] G. Ateniese, G. Felici, L. V. Mancini, A. Spognardi, A. Villani, and D. Vitali, "Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers," 2013.
- [19] N. Papernot, P. McDaniel, A. Sinha, and M. Wellman, "Towards the science of security and privacy in machine learning," 2016.
- [20] I. S. Institute, "Http dataset csic 2010," <https://www.isi.csic.es/dataset/>, Jan. 2012.
- [21] Kaggle.com, "Malicious & non-malicious url," <https://www.kaggle.com/antonyj453/urldataset>, Nov. 2017.
- [22] Kaggle.com, "Toxic comment classification challenge," <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/overview>, Mar. 2018.
- [23] X. Gong, H. Zhu, R. Deng, F. Wang, and J. Lu, "Crafting and detecting adversarial web requests," in *The 4th IEEE International Conference on Smart Cloud 2019, December 2019*, 2019.
- [24] M. Babiker, E. Karaarslan, and Y. Hoscan, "Web application attack detection and forensics: A survey," *6th International Symposium on Digital Forensic and Security, ISDFS 2018 - Proceeding*, vol. 2018-Janua, pp. 1–6, 2018.
- [25] B. W. Masduki, K. Ramli, F. A. Saputra, and D. Sugianto, "Study on implementation of machine learning methods combination for improving attacks detection accuracy on intrusion detection system (ids)," in *2015 International Conference on Quality in Research (QiR)*, pp. 56–64, 2015.
- [26] G. Tsudik and M. Almishari, "Machine-learning based security and privacy techniques for the modern web," 2012.
- [27] M. S. Islam, A. A. Mahmud, and M. R. Islam, "Machine learning approaches for modeling spammer behavior," *asia information retrieval symposium*, vol. 6458, pp. 251–260, 2010.
- [28] S. Calzavara, G. Tolomei, M. Bugliesi, and S. Orlando, "Quite a mess in my cookie jar!: leveraging machine learning to protect web authentication," in *Proceedings of the 23rd international conference on World wide web*, pp. 189–200, 2014.
- [29] V. Clincy and H. Shahriar, "Web service injection attack detection," *2017 12th International Conference for Internet Technology and Secured Transactions, ICITST 2017*, no. Line 7, pp. 173–178, 2018.
- [30] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrndić, P. Laskov, G. Giacinto, and F. Roli, "Evasion attacks against machine learning at test time," *Lecture Notes in Computer Science*, p. 387–402, 2013.
- [31] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," 2016.
- [32] L. Engstrom, J. Gilmer, G. Goh, D. Hendrycks, A. Ilyas, A. Madry, R. Nakano, P. Nakkiran, S. Santurkar, B. Tran, D. Tsipras, and E. Wallace, "A Discussion of 'Adversarial Examples Are Not Bugs, They Are Features'," *Distill*, vol. 4, no. 8, 2019.
- [33] L. Engstrom, B. Tran, D. Tsipras, L. Schmidt, and A. Madry, "Exploring the landscape of spatial robustness," 2017.
- [34] A. Athalye, L. Engstrom, A. Ilyas, and K. Kevin, "Synthesizing robust adversarial examples," *35th International Conference on Machine Learning, ICML 2018*, vol. 1, pp. 449–468, 2018.
- [35] M. Iyyer, J. Wieting, K. Gimpel, and L. Zettlemoyer, "Adversarial example generation with syntactically controlled paraphrase networks," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, (New Orleans, Louisiana), pp. 1875–1885, Association for Computational Linguistics, June 2018.
- [36] Z. Zhao, D. Dua, and S. Singh, "Generating natural adversarial examples," 2017.
- [37] J. Ebrahimi, A. Rao, D. Lowd, and D. Dou, "Hotflip: White-box adversarial examples for text classification," 2017.
- [38] 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)*, (Melbourne, Australia), pp. 856–865, Association for Computational Linguistics, July 2018.
- [39] M. Alzantot, Y. Sharma, A. Elgohary, B.-J. Ho, M. Srivastava, and K.-W. Chang, "Generating natural language adversarial examples," 2018.
- [40] S. Ren, Y. Deng, K. He, and W. Che, "Generating natural language adversarial examples through probability weighted word saliency," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, (Florence, Italy), pp. 1085–1097, Association for Computational Linguistics, July 2019.
- [41] Y. Zang, C. Yang, F. Qi, Z. Liu, M. Zhang, Q. Liu, and M. Sun, "Textual adversarial attack as combinatorial optimization," 2019.
- [42] R. Feinman, R. R. Curtin, S. Shintre, and A. B. Gardner, "Detecting adversarial samples from artifacts," 2017.
- [43] K. Grosse, P. Manoharan, N. Papernot, M. Backes, and P. McDaniel, "On the (statistical) detection of adversarial examples," 2017.
- [44] X. Huang, M. Kwiatkowska, S. Wang, and M. Wu, "Safety verification of deep neural networks," 2016.
- [45] J. H. Metzen, T. Genewein, V. Fischer, and B. Bischoff, "On detecting adversarial perturbations," 2017.
- [46] W. Xu, D. Evans, and Y. Qi, "Feature squeezing: Detecting adversarial examples in deep neural networks," *Proceedings 2018 Network and Distributed System Security Symposium*, 2018.
- [47] W. Xu, D. Evans, and Y. Qi, "Feature squeezing mitigates and detects carlini/wagner adversarial examples," 2017.
- [48] J. Bradshaw, A. G. de G. Matthews, and Z. Ghahramani, "Adversarial examples, uncertainty, and transfer testing robustness in gaussian process hybrid deep networks," 2017.
- [49] J. Jin, A. Dundar, and E. Culurciello, "Robust convolutional neural networks under adversarial noise," 2015.
- [50] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," 2017.
- [51] F. Tramèr, A. Kurakin, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," 2017.
- [52] V. Zantedeschi, M.-I. Nicolae, and A. Rawat, "Efficient defenses against adversarial attacks," 2017.
- [53] S. Zheng, Y. Song, T. Leung, and I. Goodfellow, "Improving the robustness of deep neural networks via stability training," 2016.