

---

# Adversarial Attack and Defense

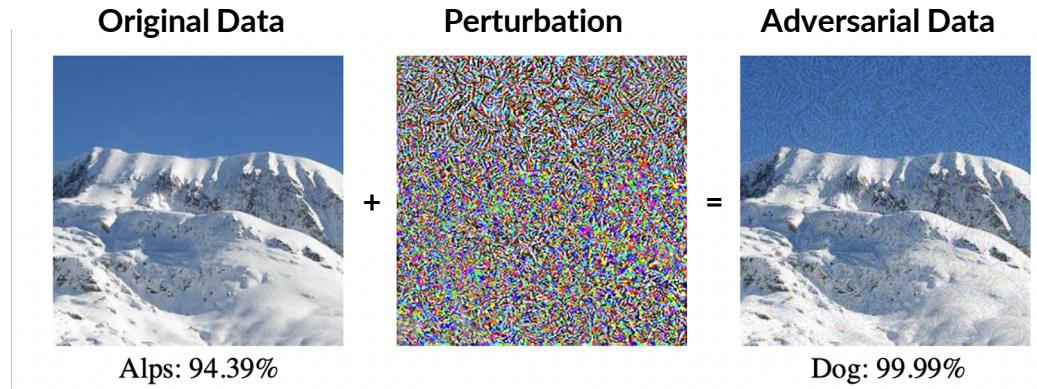
Jaechul Roh

YouTube Video Link: <https://youtu.be/maMC93Lf-mY>

---

# Introduction

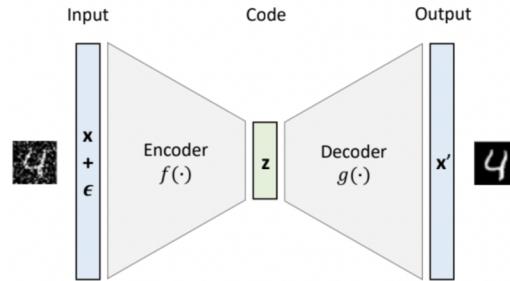
- Adversarial Examples input data with an imperceptible change
- Adversarial Examples = Original data ( $x$ ) + Perturbation with noise ( $\epsilon$ )
- Adversarial Attack induce misclassification in purpose to make machine learning models more **ROBUST**



---

## Course related material

Stacked Denoising Autoencoder = The **NOISY INPUT** will be inputted to denoising autoencoder, which will learn how to recover the original input ( $x$ ). Such method will help to create a **MORE ROBUST CODE**, so that the model will **NOT BE SENSITIVE TOWARDS NOISY INPUTS**.



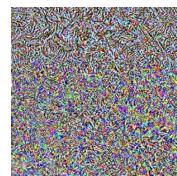
# Real-life adversarial attack example

(1) INPUT



Original Data

+

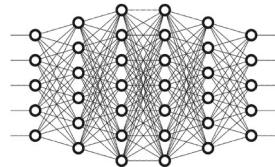


Noise

(2) ADVERSARIAL  
EXAMPLE



(3) DEEP LEARNING  
MODEL



"Green"

(4) OUTPUT

MOVE  
FORWARD

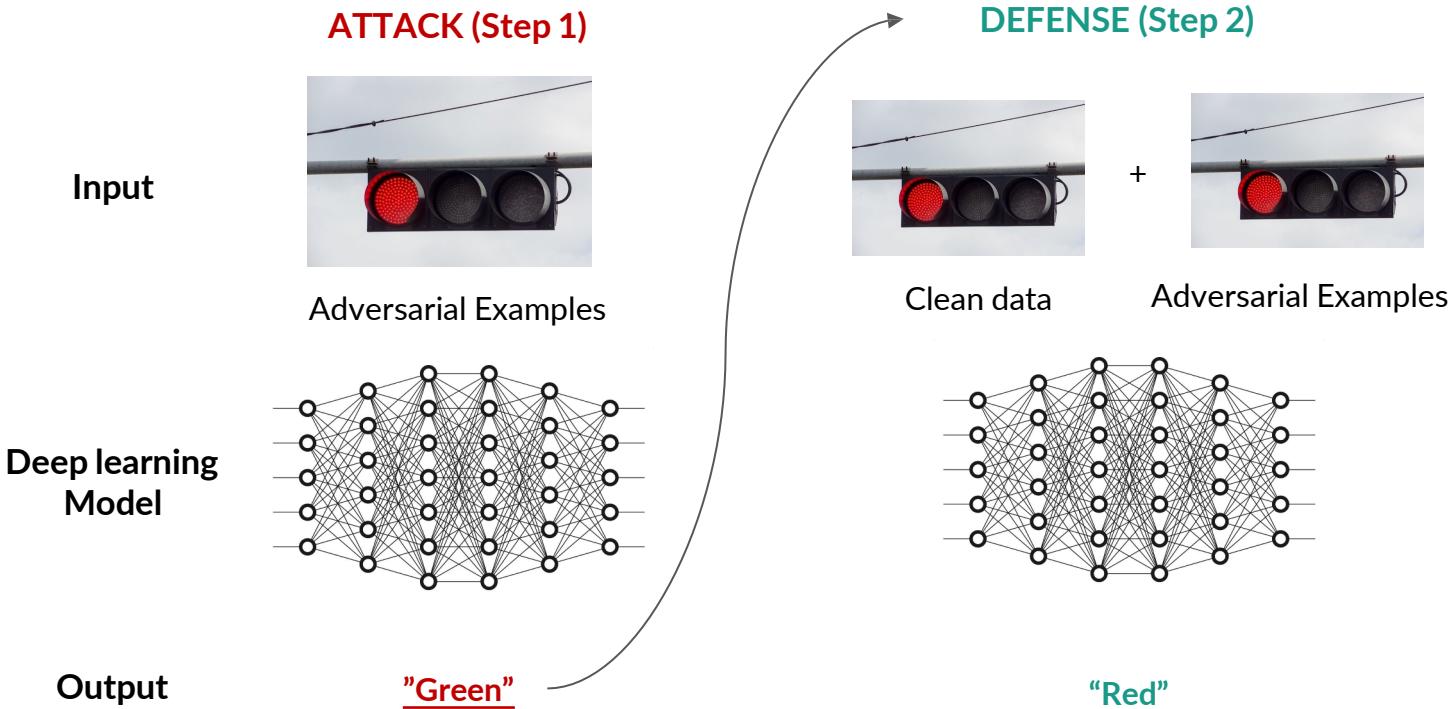
(5) ACTION



(6) CONSEQUENCE

# Adversarial Defense

---



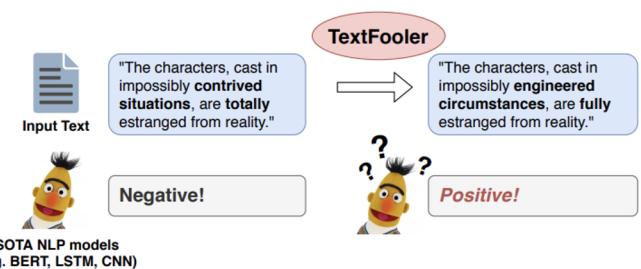
# Representative models/algorithms

- *Explaining and Harnessing Adversarial Examples (2015)*
  - Ian.J.Goodfellow, Jonathon Shlens & Christian Szegedy
- *Is BERT Really Robust? A Strong Baseline for Natural Language Attack on Text Classification and Entailment*
  - Di Jin, Zhijing Jin, Joey Tianyi Zhou, Peter Szolovits

(Image Credit: (Goodfellow et al. 2014b))

$$\begin{array}{c} \text{Original Image } x \\ \text{"panda"} \\ 57.7\% \text{ confidence} \end{array} + .007 \times \begin{array}{c} \text{Sign gradient } \text{sign}(\nabla_x J(\theta, x, y)) \\ \text{"nematode"} \\ 8.2\% \text{ confidence} \end{array} = \begin{array}{c} \text{Adversarial Image } x + \epsilon \text{sign}(\nabla_x J(\theta, x, y)) \\ \text{"gibbon"} \\ 99.3 \% \text{ confidence} \end{array}$$

Classification Task: Is this a *positive* or *negative* review?



---

# Explaining and Harnessing Adversarial Examples

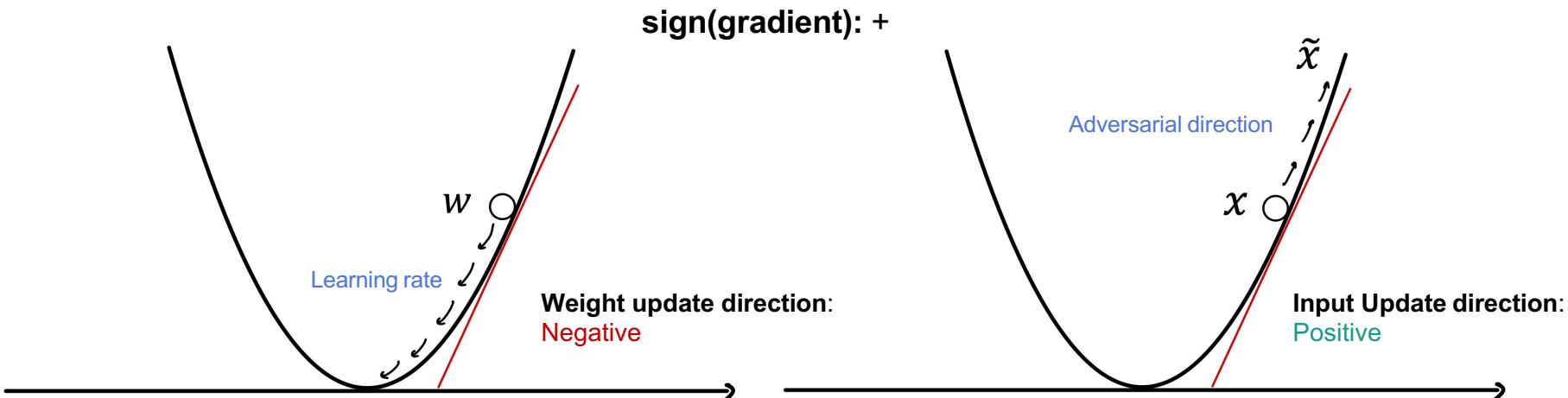
# Fast Gradient Sign Method (FGSM)

- Gradient Descent Method

OPPOSITE direction of the gradient of the cost function

- Fast Gradient Sign Method (FGSM)

SAME direction of the gradient of the cost function



---

## How adversarial example is formed

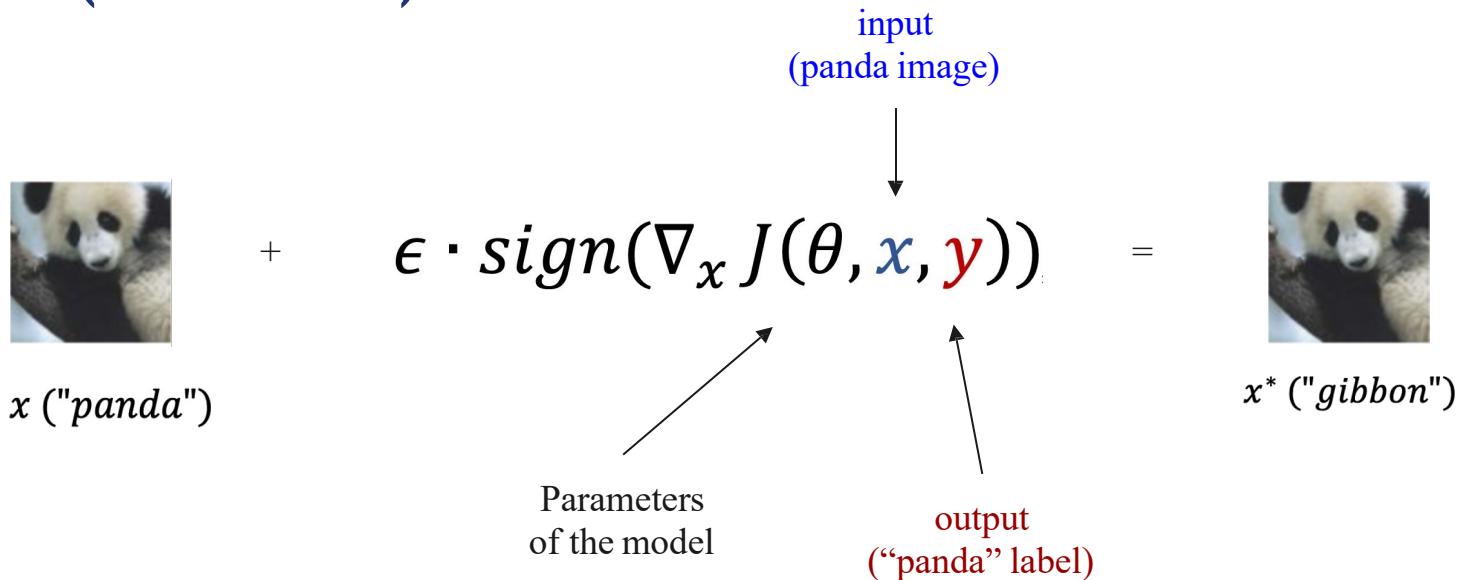
$$x + \epsilon \cdot sign(\nabla_x J(\theta, x, y))$$

Cost Function

Gradient

Adversarial Example

## FGSM (Continued)



---

## Mathematical Notation and Concepts

$$\text{activation value} \quad \longrightarrow \quad w^T \tilde{x} = w^T x + w^T \eta \quad \longleftarrow \text{activation growth}$$

↑  
original desired output

$$\text{perturbation} \longrightarrow \eta = \epsilon \cdot \text{sign}(\nabla_{x^*} J(\theta, x, y))$$

# Deciding perturbation

---

FGSM uses the “max norm constraint”:

(In all definitions,  $x = (x_1, x_2, \dots, x_n)$ )

$$L^\infty \text{ distance: } \|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

$L^\infty$  : moving as many pixels as possible but only by a small number

$$L^1 \text{ distance: } \|x\|_1 = \sum_{i=1}^n |x_i|$$

$L^1$  : summed absolute value difference between  $x$  and  $x^*$

## Example 1: 1-Dimensional Calculation

$$w^T \tilde{x} = w^T x + w^T \eta = w^T(x + \eta)$$

$$\begin{array}{c|c|c} x & w & w^T x \\ \hline \begin{pmatrix} 3 \\ -2 \\ 5 \end{pmatrix} & * \begin{pmatrix} 7 \\ 10 \\ 20 \end{pmatrix} & = \begin{pmatrix} 21 \\ -20 \\ 100 \end{pmatrix} \Rightarrow 101 \end{array}$$

activation value  
(**WITHOUT** perturbation)

$$\begin{array}{c|c|c|c|c} \eta = sign(w) & x & \eta & w & w^T \tilde{x} \\ \hline sign\left(\begin{pmatrix} 7 \\ 10 \\ 20 \end{pmatrix}\right) = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} & \begin{pmatrix} 3 \\ -2 \\ 5 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} & & * \begin{pmatrix} 7 \\ 10 \\ 20 \end{pmatrix} = \begin{pmatrix} 28 \\ -10 \\ 120 \end{pmatrix} \Rightarrow 138 & activation value \\ & & & & (\text{WITH perturbation}) \end{array}$$

## Example 2: 3-Dimensional Calculation

$$x + \epsilon \cdot sign(\nabla_x J(\theta, x, y))$$

$sign(w_x) \rightarrow \text{POSITIVE}$

$sign(w_y) \rightarrow \text{NEGATIVE} \times \epsilon_{vector} = -\epsilon + x_{vector}$

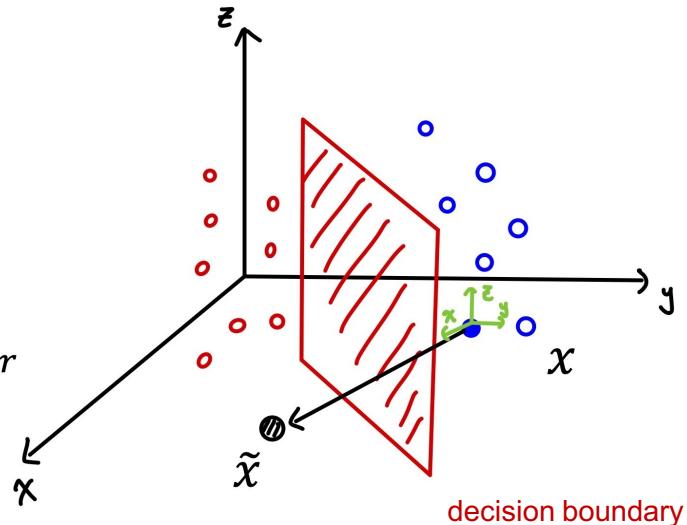
$sign(w_z) \rightarrow \text{POSITIVE}$

$+\epsilon$

$-\epsilon$

$+\epsilon$

$$= x^*_{vector}$$



## Adversarial Defense (FGSM)

$\alpha$  : proportion to use between the original data and the adversarial example

- (1)  $\tilde{J}(\theta, x, y)$  : cost function of the original data
  - (2)  $J(\theta, \tilde{x}, y)$  : cost function of the adversarial example
  - (3)  $J(\theta, x, y)$  : cost function of both original data AND adversarial example

---

# Adversarial Attack in Natural Language Processing



# Hardship of natural language adversarial attack

## Image domain (CONTINUOUS values)

Adding a minimal noise to the pixels is not noticeable through naked eyes

## Text domain (DISCRETE values)

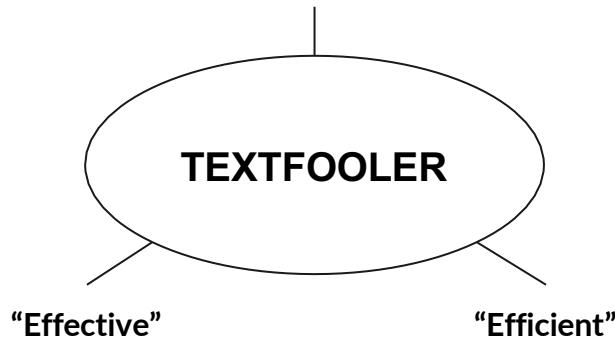
The difference between the original text and the adversarial example is easily recognizable

# Introduction

---

- Proposing *TextFooler*

“Utility-preserving”

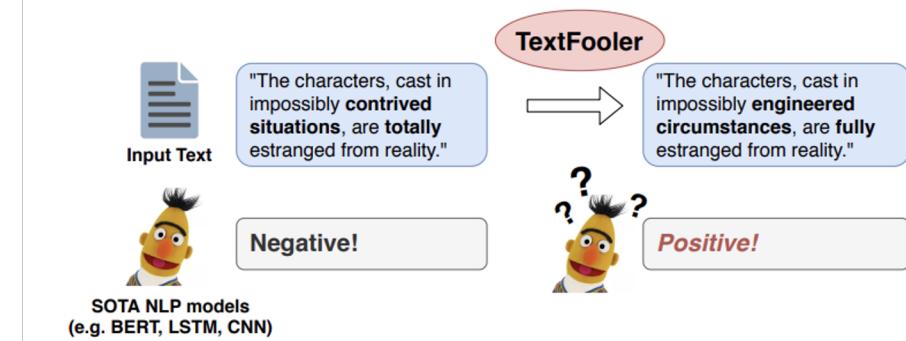


- How to test such the robustness?

Models: 1. WordLSTM 2. WordCNN 3. BERT

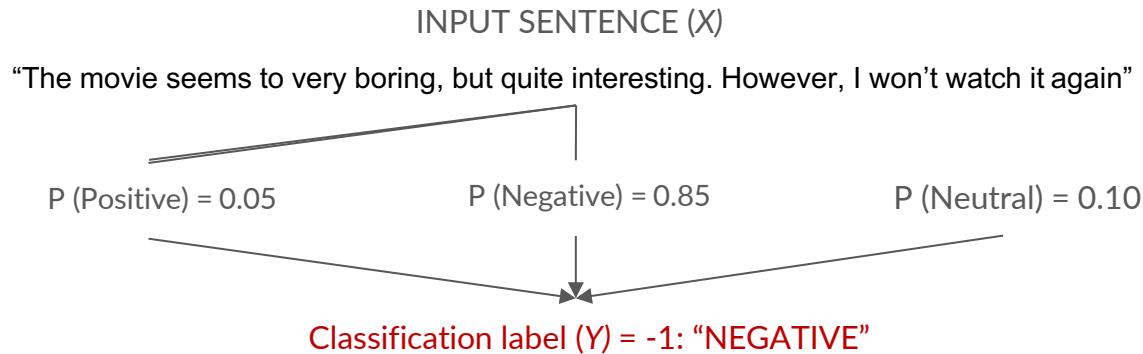
Task: 5 classification tasks and 2 textual entailment tasks

Classification Task: Is this a *positive* or *negative* review?



# Classification & Recognizing Text Entailment (NLP tasks)

---



---

## HYPOTHESIS SENTENCE

Input: “It was a very touching movie”

Entailment

“The movie made me cry”

Contradiction

“The movie made me laugh”

Neutral

“There was no discount for the movie ticket”

# Forming text adversarial example

---

- Text adversarial example need to meet the following requirement:

$$F(X_{\text{adv}}) \neq F(X), \text{ and } \text{Sim}(X_{\text{adv}}, X) \geq \epsilon,$$

Classification result  
of adversarial text

Classification result  
of original data

Semantic similarity  
between  $X_{\text{adv}}$  and  $X$

Minimum  
similarity

# TEXTFOOLER Attack Algorithm

---

Steps:

1. Word Importance Score
2. Word Transformer
  - a. Adversarial example candidates POS (Part of Speech) checking
  - b. Semantic Similarity Filter
  - c. Finalizing Adversarial Example

---

## Algorithm 1 Adversarial Attack by TEXTFOOLER

---

**Input:** Sentence example  $X = \{w_1, w_2, \dots, w_n\}$ , the corresponding ground truth label  $Y$ , target model  $F$ , sentence similarity function  $\text{Sim}(\cdot)$ , sentence similarity threshold  $\epsilon$ , word embeddings  $\text{Emb}$  over the vocabulary  $\text{Vocab}$ .

**Output:** Adversarial example  $X_{\text{adv}}$

```
1: Initialization:  $X_{\text{adv}} \leftarrow X$ 
2: for each word  $w_i$  in  $X$  do
3:   Compute the importance score  $I_{w_i}$  via Eq. (2)
4: end for
5:
6: Create a set  $W$  of all words  $w_i \in X$  sorted by the descending order of their importance score  $I_{w_i}$ .
7: Filter out the stop words in  $W$ .
8: for each word  $w_j$  in  $W$  do
9:   Initiate the set of candidates  $\text{CANDIDATES}$  by extracting the top  $N$  synonyms using  $\text{CosSim}(\text{Emb}_{w_j}, \text{Emb}_{\text{word}})$  for each word in  $\text{Vocab}$ .
10:   $\text{CANDIDATES} \leftarrow \text{POSFilter}(\text{CANDIDATES})$ 
11:   $\text{FINCANDIDATES} \leftarrow \{\}$ 
12:  for  $c_k$  in  $\text{CANDIDATES}$  do
13:     $X' \leftarrow \text{Replace } w_j \text{ with } c_k \text{ in } X_{\text{adv}}$ 
14:    if  $\text{Sim}(X', X_{\text{adv}}) > \epsilon$  then
15:      Add  $c_k$  to the set  $\text{FINCANDIDATES}$ 
16:       $Y_k \leftarrow F(X')$ 
17:       $P_k \leftarrow F_{Y_k}(X')$ 
18:    end if
19:  end for
20:  if there exists  $c_k$  whose prediction result  $Y_k \neq Y$  then
21:    In  $\text{FINCANDIDATES}$ , only keep the candidates  $c_k$  whose prediction result  $Y_k \neq Y$ 
22:     $c^* \leftarrow \underset{c \in \text{FINCANDIDATES}}{\operatorname{argmax}} \text{Sim}(X, X'_{w_j \rightarrow c})$ 
23:     $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
24:  return  $X_{\text{adv}}$ 
25: else if  $P_{Y_k}(X_{\text{adv}}) > \min_{c_k \in \text{FINCANDIDATES}} P_k$  then
26:    $c^* \leftarrow \underset{c_k \in \text{FINCANDIDATES}}{\operatorname{argmin}} P_k$ 
27:    $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
28: end if
29: end for
30: return None
```

---

# 1. Word Importance Ranking

---

“Measuring the influence the word,  $\underline{w_i}$  ”

$$I_{w_i} = \begin{cases} \underline{F_Y(X) - F_Y(X_{\setminus w_i})}, & \text{if } F(X) = F(\underline{X_{\setminus w_i}}) = Y \\ (\underline{F_Y(X) - F_Y(X_{\setminus w_i})}) + (\underline{F_{\bar{Y}}(X_{\setminus w_i})} - \underline{F_{\bar{Y}}(X)}), \\ & \text{if } F(X) = Y, F(\underline{X_{\setminus w_i}}) = \bar{Y}, \text{ and } Y \neq \bar{Y}. \end{cases}$$

Classification output      Input without the word,  $\underline{w_i}$   
Two different labels

“Prediction change before, and after the word,  $\underline{w_i}$  ”

## a. Candidates and POS Checking

---

**Output:** Adversarial example  $X_{\text{adv}}$

- 1: Initialization:  $X_{\text{adv}} \leftarrow X$
- 2: **for** each word  $w_i$  in  $X$  **do**
- 3: (1) Compute the importance score  $I_{w_i}$  via Eq. (2)
- 4: **end for**
  
- 8: **for** each word  $w_j$  in  $W$  **do** (3)
- 9:   Initiate the set of candidates CANDIDATES by extracting  
     the top  $N$  synonyms using CosSim(Emb<sub>w<sub>j</sub></sub>, Emb<sub>word</sub>) for  
     each word in Vocab. (2)
- 10:   CANDIDATES  $\leftarrow$  POSFilter(CANDIDATES) (4)

(1): Process importance score  
for every word in the  
sentence example

(2): Cosine Similarity Score  
between the Embedding(deleting  
word) and Embedding(Vocab)

(3): Extract top N synonyms and  
append to CANDIDATES list

(4): Check POS (Part of Speech) for  
every candidate word and filter

## b. Semantic Similarity Filter

---

(3) Cosine Similarity between  $X$  (original sentence) and  $X_{adv}$  (Adversarial Example)

```
11: FINCANDIDATES  $\leftarrow \{ \}$ 
12: for  $c_k$  in CANDIDATES do (2)
13:   (1)  $X' \leftarrow$  Replace  $w_j$  with  $c_k$  in  $X_{adv}$ 
14:   if  $Sim(X', X_{adv}) > \epsilon$  then (4)
15:     Add  $c_k$  to the set FINCANDIDATES
16:     (3)  $Y_k \leftarrow F(X')$ 
17:      $P_k \leftarrow F_{Y_k}(X')$ 
18:   end if
19: end for
```

(1) Substitute that specific word in the sentence with each of the words in the CANDIDATES

(2) Such sentence becomes  $X_{adv}$  (Adversarial Example)

(4) Words with similarity score  $> \epsilon$  (defined by the programmer) will be stored in FINCANDIDATES list

## c. Finalizing the Adversarial Example

In the descending order of similarity scores, replace the word:

```
(1)
20: if there exists  $c_k$  whose prediction result  $Y_k \neq Y$  then
21:   In FINCANDIDATES, only keep the candidates  $c_k$  whose
22:   prediction result  $Y_k \neq Y$ 
23:    $c^* \leftarrow \operatorname{argmax}_{c \in \text{FINCANDIDATES}} \text{Sim}(X, X'_{w_j \rightarrow c})$ 
24:    $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
25:   return  $X_{\text{adv}}$ 
26: else if  $P_{Y_k}(X_{\text{adv}}) > \min_{c_k \in \text{FINCANDIDATES}} P_k$  then
27:    $c^* \leftarrow \operatorname{argmin}_{c_k \in \text{FINCANDIDATES}} P_k$ 
28:    $X_{\text{adv}} \leftarrow \text{Replace } w_j \text{ with } c^* \text{ in } X_{\text{adv}}$ 
29: end if
30: end for
31: return None (2)
```

(1) IF the prediction of the target model changes:

- Within those candidates that changed the output of the target model
- Select the word that had the highest similarity score between X and  $X_{\text{adv}}$ .

(2) ELSE IF choose the word with the least confidence level

(word that is most likely to change the prediction of the model)

- Prediction changed → Attack Success!

Run through this process in the descending order of importance score of each word.

# Summary

---

Adversarial **ATTACK**



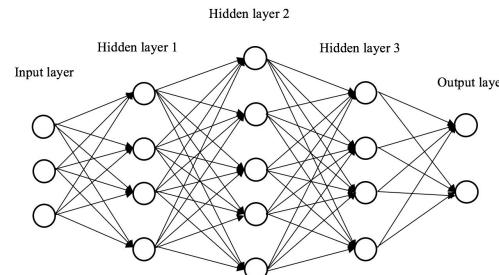
Adversarial **TRAINING**



**ROBUST**  
deep learning model



$x^*$  ("*gibbon*")



$x$  ("*panda*")

# Reference

---

- Attack & Defense (1): Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *arXiv preprint arXiv:1412.6572* (2014).
- Attack & Defense (2): Jin, Di, et al. "Is bert really robust? a strong baseline for natural language attack on text classification and entailment." *Proceedings of the AAAI conference on artificial intelligence*. Vol. 34. No. 05. 2020.
- Kariya, Mahendra. "Paper Discussion: Explaining and Harnessing Adversarial Examples." *Medium*, Medium, 16 Nov. 2018, <https://medium.com/@mahendrakariya/paper-discussion-explaining-and-harnessing-adversarial-examples-908a1b7123b5>.
- Dong, Yinpeng, et al. "Boosting adversarial attacks with momentum." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
- **1-dimensional example:** "Explaining And Harnessing Adversarial Examples 논문-리뷰." *Explaining And Harnessing Adversarial Examples* 논문-리뷰, 10 July 2020, [velog.io/@miae0112/Explaining-And-Harnessing-Adversarial-Examples-%EB%85%BC%EB%AC%B8-%EB%A6%AC%EB%B7%BO](https://velog.io/@miae0112/Explaining-And-Harnessing-Adversarial-Examples-%EB%85%BC%EB%AC%B8-%EB%A6%AC%EB%B7%BO).
- "[텍스트 분류 모델 공격 기법] TextFooler: Is BERT Really Robust?" YouTube, 22 Nov. 2020, [www.youtube.com/watch?v=EF-IYFTKZiE&t=992s](https://www.youtube.com/watch?v=EF-IYFTKZiE&t=992s).
- Pan, Zhixin, and Prabhat Mishra. "Fast Approximate Spectral Normalization for Robust Deep Neural Networks." *arXiv preprint arXiv:2103.13815* (2021).
- Cho, Yoon Sang. "Adversarial Attacks and Defenses in Deep Learing." *Adversarial Attacks and Defenses in Deep Learing*, 2020, pp. 5–32, dmqm.korea.ac.kr/activity/seminar/289.

---



# Thank you!

Jaechul Roh (ID: 20473590)