# Adversarial Attack and Defense

Jaechul Roh (ID: 20473590)

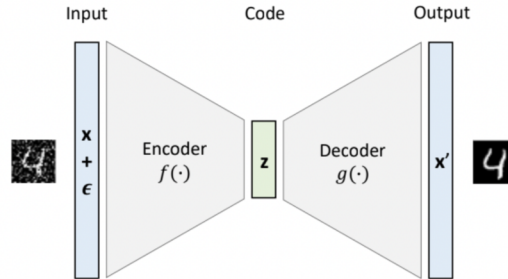YouTube Video Link: **https://youtu.be/maMC93Lf-mY**

# Introduction

- **Adversarial Examples** input data with an <u>imperceptible change</u>

- **Adversarial Examples** = Original data ($x$) + Perturbation with noise ($\epsilon$)

- **Adversarial Attack** induce misclassification in purpose to make machine learning models more **ROBUST**

| Original Data | Perturbation | Adversarial Data |
|:---:|:---:|:---:|
| | + | = |
| Alps: 94.39% | | Dog: 99.99% |

# 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**.
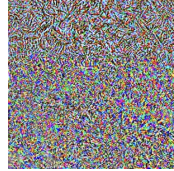
# Real-life adversarial attack example
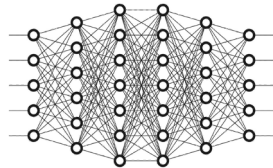
**(1) INPUT**

Original Data          +          Noise

"Green"          **(4) OUTPUT**

**(2) ADVERSARIAL EXAMPLE**

**MOVE FORWARD**          **(5) ACTION**

**(3) DEEP LEARNING MODEL**

**(6) CONSEQUENCE**

# Adversarial Defense



**ATTACK (Step 1)**

**DEFENSE (Step 2)**

Input

Adversarial Examples

Clean data + Adversarial Examples

Deep learning Model
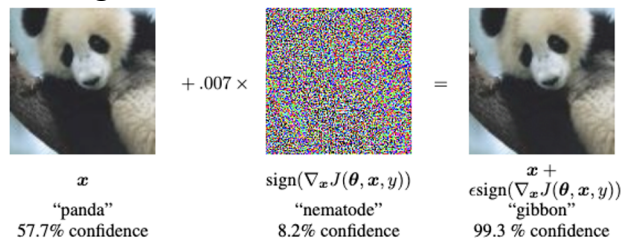
Output

"Green"

"Red"

# Representative models/algorithms

- *Explaining and Harnessing Adversarial Examples (2015)*
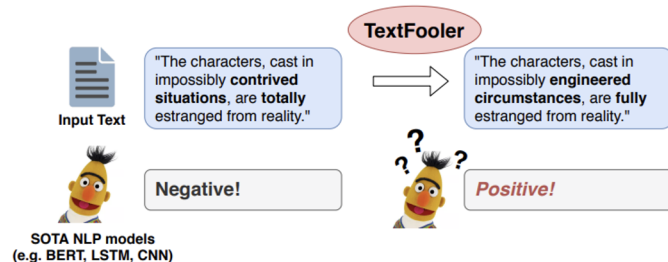  - *Ian.J.Goodfellow, Jonathon Shlens & Christian Szegedy*



$$+ .007 \times$$

$$=$$

$x$
"panda"
57.7% confidence

$\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"nematode"
8.2% confidence

$x + \epsilon \text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"gibbon"
99.3 % confidence

- *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*



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

**TextFooler**

Input Text

"The characters, cast in impossibly **contrived situations**, are **totally** estranged from reality."

"The characters, cast in impossibly **engineered circumstances**, are **fully** estranged from reality."

Negative!

*Positive!*

SOTA NLP models
(e.g. BERT, LSTM, CNN)

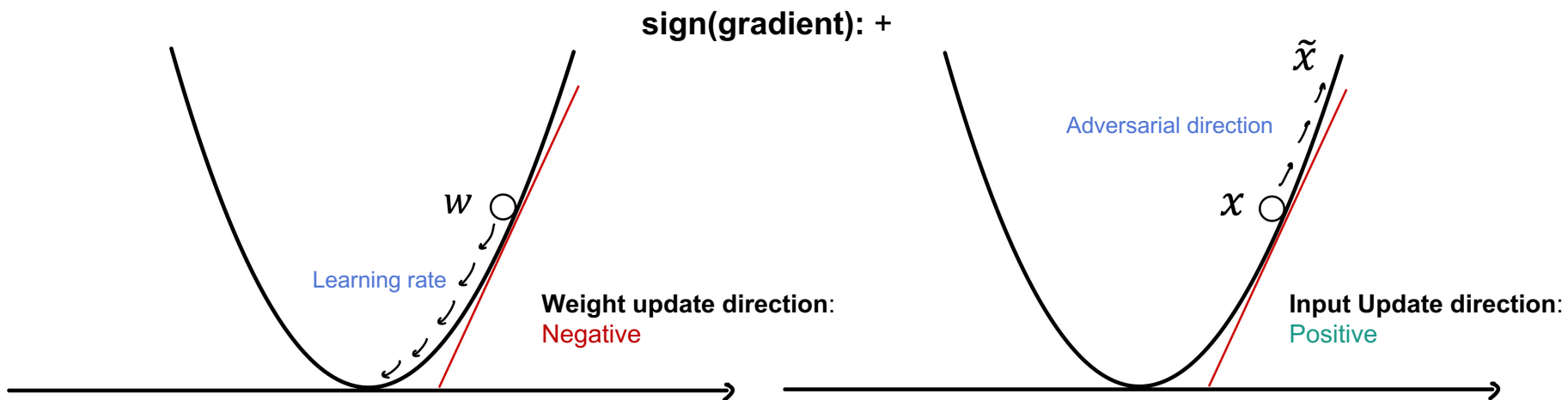# Explaining and Harnessing Adversarial Examples

# 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

sign(gradient): +

$w$

Learning rate

**Weight update direction**:
Negative

$\tilde{x}$

Adversarial direction

$x$

**Input Update direction**:
Positive

# How adversarial example is formed

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

**Cost Function**

**Gradient**

**Adversarial Example**

# FGSM (Continued)



input
(panda image)

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

$x$ ("panda")   +   =   $x^*$ ("gibbon")

Parameters
of the model

output
("panda" label)

# Mathematical Notation and Concepts

activation value
(with perturbation) $\longrightarrow$  $w^T\tilde{x} = w^Tx + w^T\eta$ $\longleftarrow$ activation growth

original desired output

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

# Deciding perturbation

FGSM uses the "*max norm* constraint":

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

$$L^\infty\ distance\colon \parallel x \parallel_\infty = \max_{1 \le i \le n} |x_i| \qquad L^1\ distance\colon \parallel x \parallel_1 = \sum_{i=1}^{n} |x_i|$$

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

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

# Example 1: 1-Dimensional Calculation

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

$x$      $w$      $w^T x$

$$\begin{pmatrix} 3 \\ -2 \\ 5 \end{pmatrix} \begin{matrix} * \\ * \\ * \end{matrix} \begin{pmatrix} 7 \\ 10 \\ 20 \end{pmatrix} = \begin{pmatrix} 21 \\ -20 \\ 100 \end{pmatrix} \Rightarrow 101$$

activation value
(**WITHOUT** perturbation)

$\eta = sign(w)$

$$sign\left(\begin{pmatrix} 7 \\ 10 \\ 20 \end{pmatrix}\right) = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$x$      $\eta$      $w$      $w^T \tilde{x}$

$$\left[ \begin{pmatrix} 3 \\ -2 \\ 5 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \right] * \begin{pmatrix} 7 \\ 10 \\ 20 \end{pmatrix} = \begin{pmatrix} 28 \\ -10 \\ 120 \end{pmatrix} \Rightarrow 138$$

activation value
(**WITH** perturbation)

# Example 2: 3-Dimensional Calculation

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

$$sign(w_x) \rightarrow \text{POSITIVE} \qquad +\epsilon$$
$$sign(w_y) \rightarrow \text{NEGATIVE} \quad \times \ \epsilon_{vector} \ = \ -\epsilon \quad + \quad x_{vector}$$
$$sign(w_z) \rightarrow \text{POSITIVE} \qquad +\epsilon$$

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

decision boundary

# Adversarial Defense (FGSM)

$$\tilde{J}(\theta, x, y) = \alpha \cdot J(\theta, x, y) + (1 - \alpha) \cdot J(\theta, \tilde{x}, y)$$

(3)                 (1)                (2)

$\alpha$ : **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
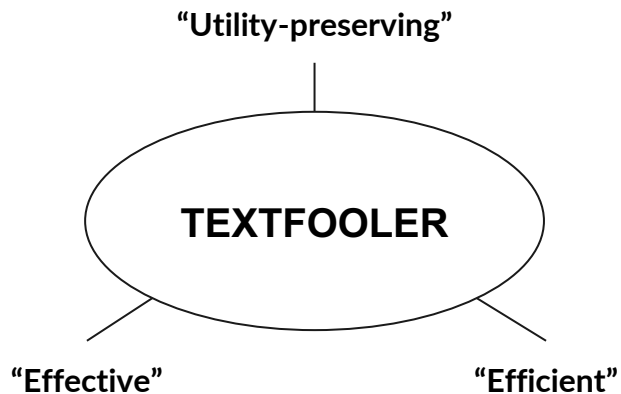
**Image domain (CONTINUOUS values)**

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
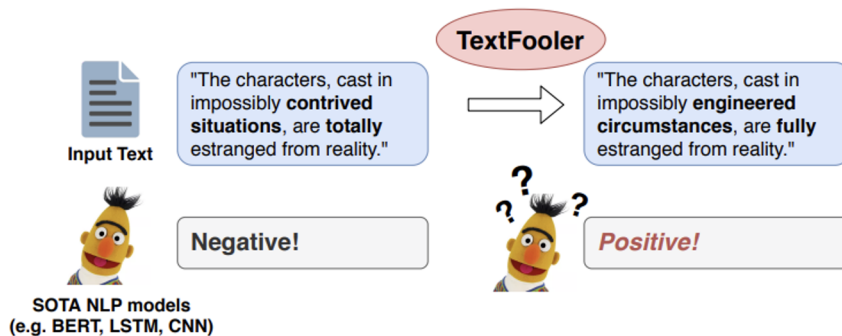
# Introduction

- **Proposing *TextFooler***

"Utility-preserving"

**TEXTFOOLER**

"Effective"          "Efficient"

- **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?

Input Text: "The characters, cast in impossibly **contrived situations**, are **totally** estranged from reality."

TextFooler → "The characters, cast in impossibly **engineered circumstances**, are **fully** estranged from reality."

SOTA NLP models (e.g. BERT, LSTM, CNN)

Negative!          *Positive!*

# Classification & Recognizing Text Entailment (NLP tasks)

INPUT SENTENCE ($X$)

"The movie seems to very boring, but quite interesting. However, I won't watch it again"

P (Positive) = 0.05          P (Negative) = 0.85          P (Neutral) = 0.10

Classification label ($Y$) = -1: "NEGATIVE"

---

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**:

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

**Classification result of adversarial text**   **Classification result of original data**   **Semantic similarity between *Xadv 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 $X = \{w_1, w_2, ..., w_n\}$, the corresponding ground truth label $Y$, target model $F$, sentence similarity function $\text{Sim}(\cdot)$, sentence similarity threshold $\epsilon$, word embeddings Emb over the vocabulary 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 CANDIDATES by extracting the top $N$ synonyms using $\text{CosSim}(\text{Emb}_{w_j}, \text{Emb}_{\text{word}})$ for each word in Vocab.
10:      CANDIDATES $\leftarrow$ POSFilter(CANDIDATES)
11:      FINCANDIDATES $\leftarrow \{\,\}$
12:      **for** $c_k$ in CANDIDATES **do**
13:          $X' \leftarrow$ Replace $w_j$ with $c_k$ in $X_{\text{adv}}$
14:          **if** $\text{Sim}(X', X_{\text{adv}}) > \epsilon$ **then**
15:              Add $c_k$ to the set 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 FINCANDIDATES, only keep the candidates $c_k$ whose prediction result $Y_k \neq Y$
22:          $c^* \leftarrow \underset{c \in \text{FINCANDIDATES}}{\arg\max}\ \text{Sim}(X, X'_{w_j \rightarrow c})$
23:          $X_{\text{adv}} \leftarrow$ Replace $w_j$ with $c^*$ in $X_{\text{adv}}$
24:          **return** $X_{\text{adv}}$
25:      **else if** $P_{Y_k}(X_{\text{adv}}) > \underset{c_k \in \text{FINCANDIDATES}}{\min}\ P_k$ **then**
26:          $c^* \leftarrow \underset{c_k \in \text{FINCANDIDATES}}{\arg\min}\ P_k$
27:          $X_{\text{adv}} \leftarrow$ Replace $w_j$ with $c^*$ in $X_{\text{adv}}$
28:      **end if**
29:  **end for**
30:  **return** None

# 1. Word Importance Ranking

"*Measuring the influence the word,* $w_i$ "

**Input without the word,** $w_i$

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

**Classification output**                    **Two different labels**

"*Prediction change before, and after the word,* $w_i$ "

# a. Candidates and POS Checking

**Output:** Adversarial example $X_{\mathrm{adv}}$
1: Initialization: $X_{\mathrm{adv}} \leftarrow X$
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**
9:     Initiate the set of candidates CANDIDATES by extracting (3)
    the top $N$ synonyms using $\mathrm{CosSim}(\mathrm{Emb}_{w_j}, \mathrm{Emb}_{\mathrm{word}})$ 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 Xadv (Adversarial Example)*

```
11:     FINCANDIDATES ← { }
12:     for ck in CANDIDATES do              (2)
13: (1)  X' ← Replace wj with ck in Xadv
14:      if Sim(X', Xadv) > ε then          (4)
15: (3)   Add ck to the set FINCANDIDATES
16:       Yk ← F(X')
17:       Pk ← FYk(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 *Xadv (Adversarial Example)*

(4) Words with **similarity score > ε (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 ≠ Y then
21:        In FINCANDIDATES, only keep the candidates c_k whose
           prediction result Y_k ≠ Y
22:        c* ←   argmax    Sim(X, X'_{w_j→c})
              c∈FINCANDIDATES
23:        X_adv ← Replace w_j with c* in X_adv
24:        return X_adv
25:    else if P_{Y_k}(X_adv) >    min      P_k then
                              c_k∈FINCANDIDATES
26:        c* ←   argmin    P_k
              c_k∈FINCANDIDATES
27:        X_adv ← Replace w_j with c* in X_adv
28:    end if
29: end for
30: 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 Xadv.**

(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!**

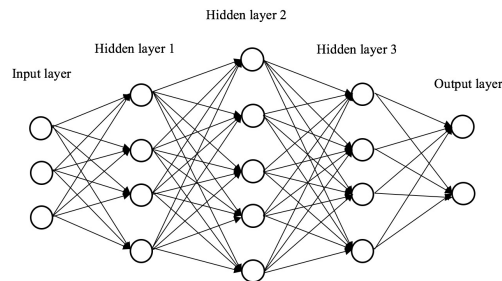Run through this process in the descending order of importance score of each word.

# Summary

Adversarial ATTACK          Adversarial TRAINING          ROBUST
deep learning model



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

Hidden layer 2

Hidden layer 1          Hidden layer 3

Input layer          Output layer

A Deep Learning Model

$x$ ("$panda$")

# 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/@miai0112/Explaining-And-Harnessing-Adversarial-Examples-%EB%85%BC%EB%AC%B8-%EB%A6%AC%EB%B7%B0.

- "[텍스트 분류 모델 공격 기법] TextFooler: Is BERT Really Robust?" *YouTube*, 22 Nov. 2020, 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)