# 李宏毅 (Hung-yi Lee) · HYLEE | Machine Learning (2021)

## HYLEE(2021)· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

视频 · B 站 [ 扫码或点击链接 ]
https://www.bilibili.com/video/BV1fM4y137M4

课件 & 代码 · 博客 [ 扫码或点击链接 ]
http://blog.showmeai.tech/ntu-hylee-ml

机器学习　Auto-encoder　生成式对抗网络　学习率
深度学习　卷积神经网络　GAN　自监督　自注意力机制
批次标准化　神经网络压缩　强化学习　元学习　Transformer

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS231n |

## # Awesome AI Courses Notes Cheatsheets· 持续更新中

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏

称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]

# Adversarial Attack

## Hung-yi Lee

# Motivation

- You have trained many neural networks.

- We seek to deploy neural networks in the real world.

- Are networks robust to the inputs that are built to fool them?
  - Useful for spam classification, malware detection, network intrusion detection, etc.

Aim to fool the network
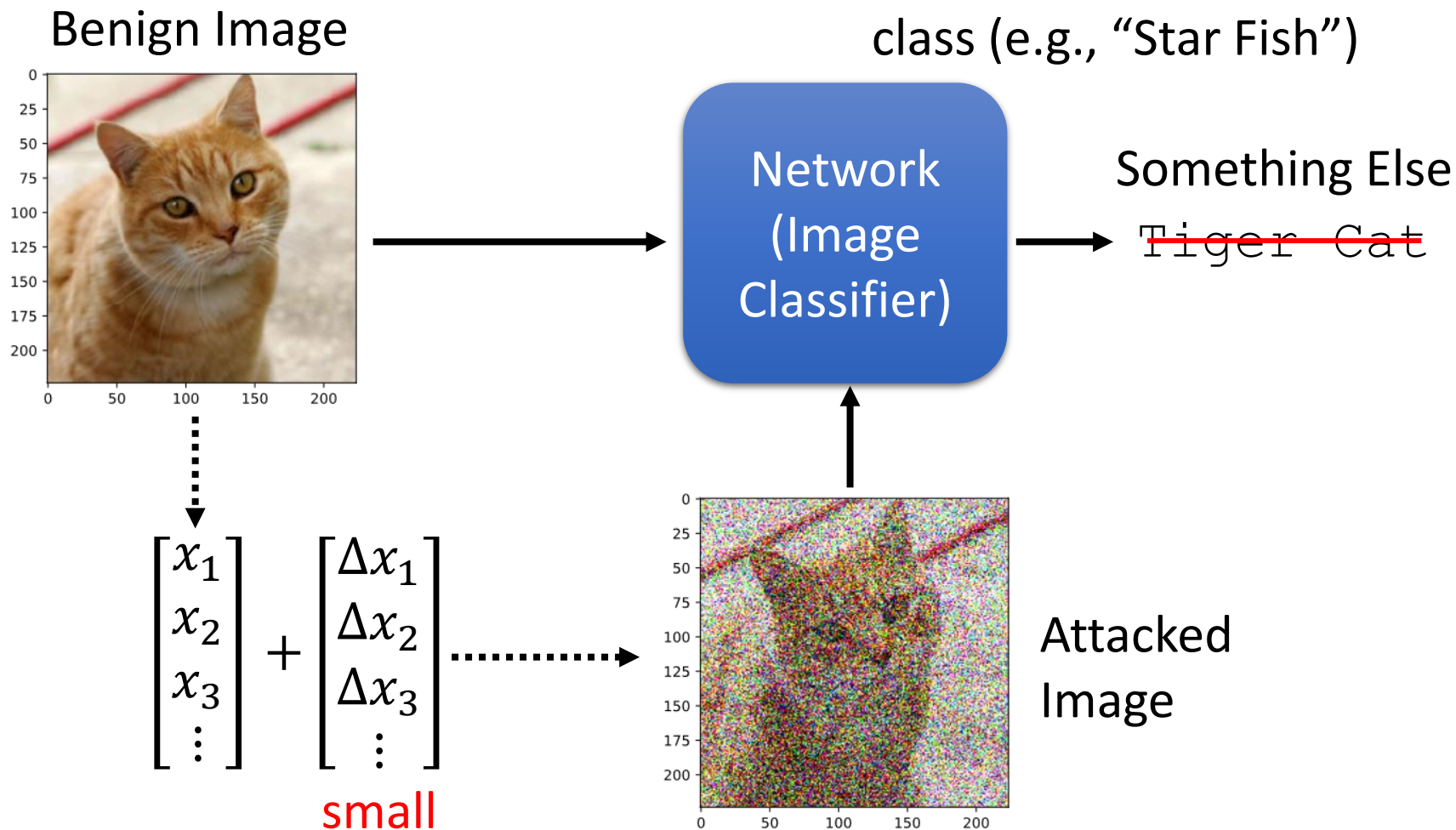
# How to Attack

# Example of Attack

Anything other than "Cat"

**Targeted**

Misclassified as a specific class (e.g., "Star Fish")

Benign Image



Network (Image Classifier)

Something Else

~~Tiger Cat~~

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix} + \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \Delta x_3 \\ \vdots \end{bmatrix}$$
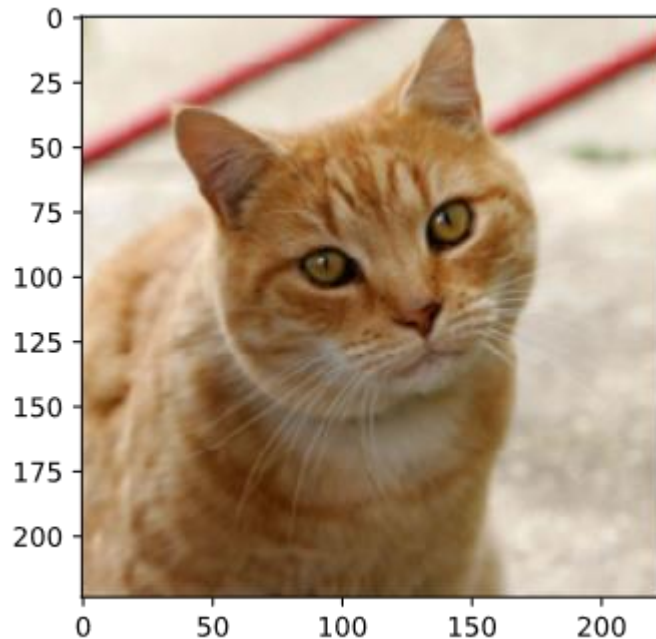
small



Attacked Image

# Example of Attack

Network = ResNet-50

The target is "Star Fish"

Benign Image



Tiger Cat

0.64

Attacked Image



Star Fish

1.00

# Example of Attack



Benign Image

Attac...

=

50x

Tiger Cat
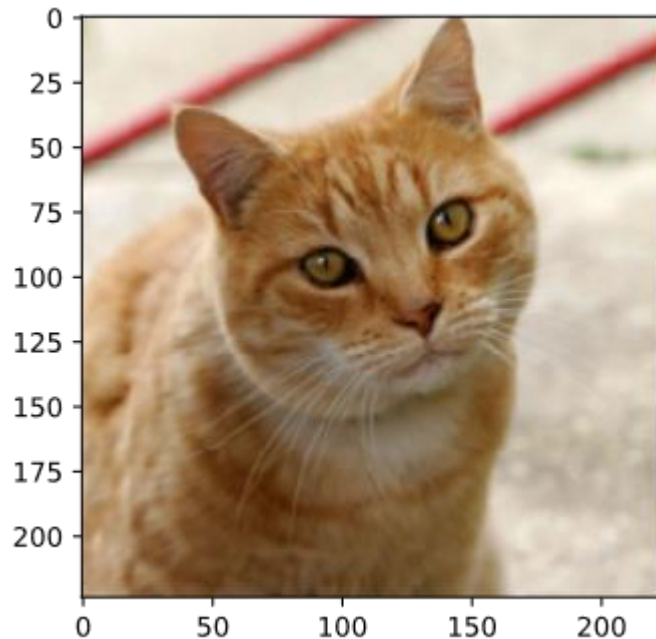
0.64

Star Fish
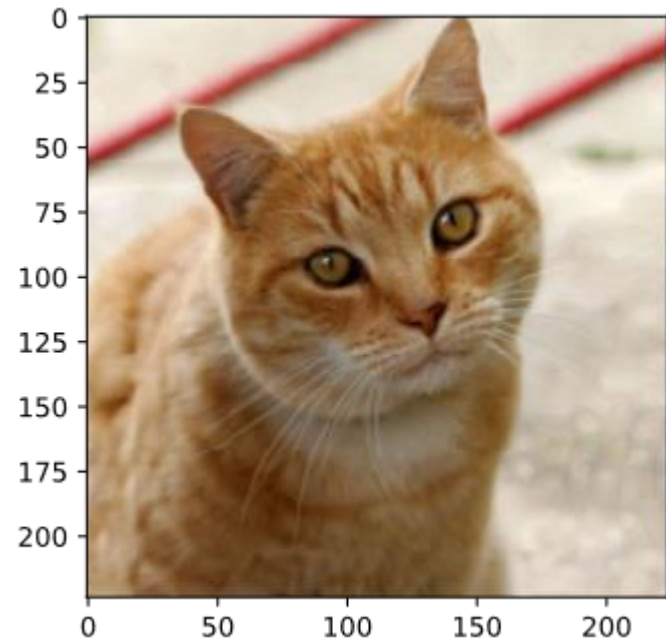
1.00

# Example of Attack

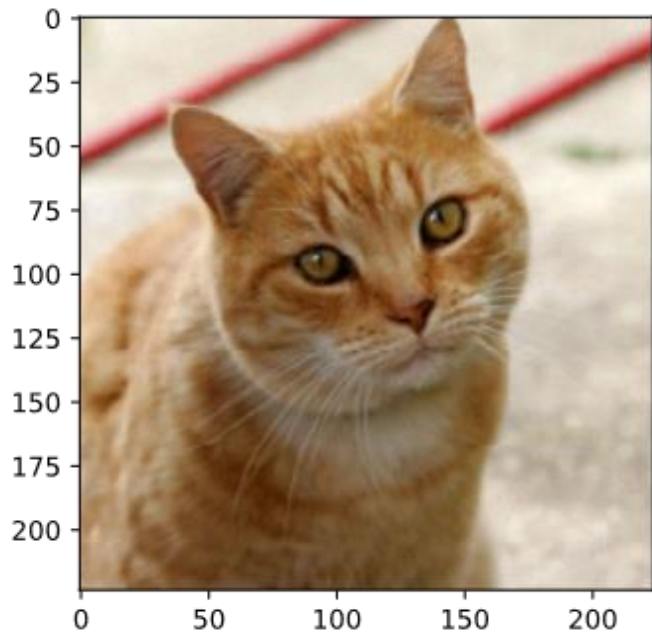Network = ResNet-50

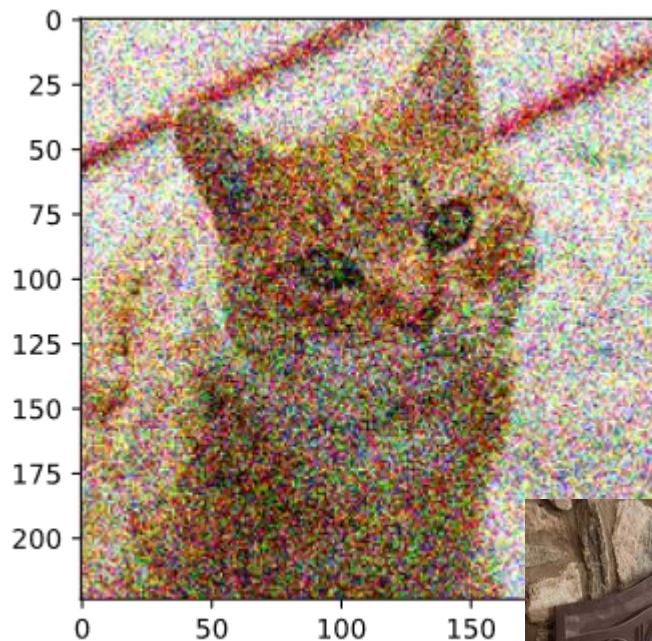The target is "Keyboard"

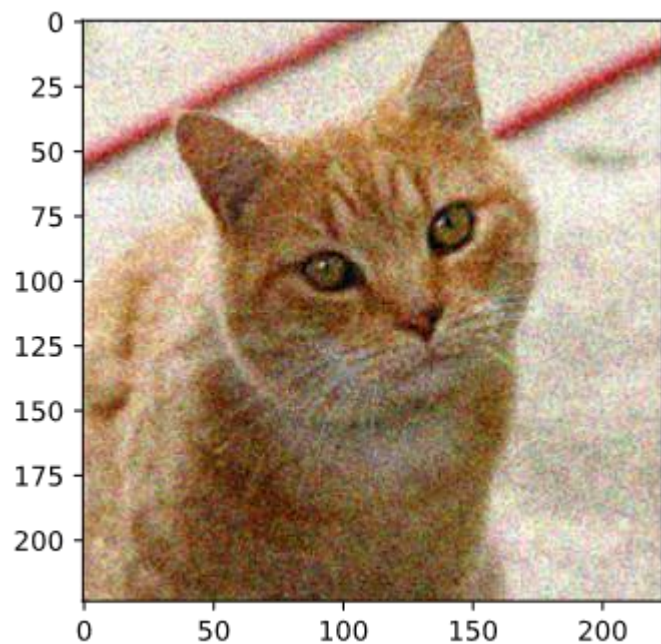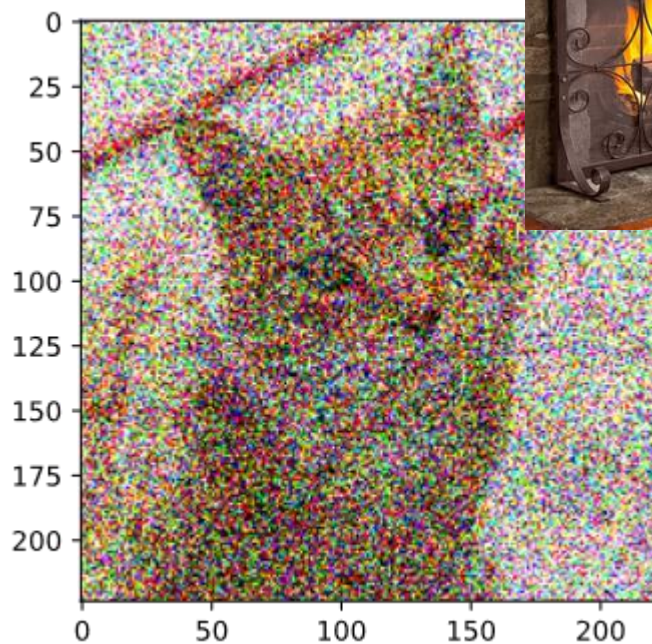Benign Image



Tiger Cat
0.64

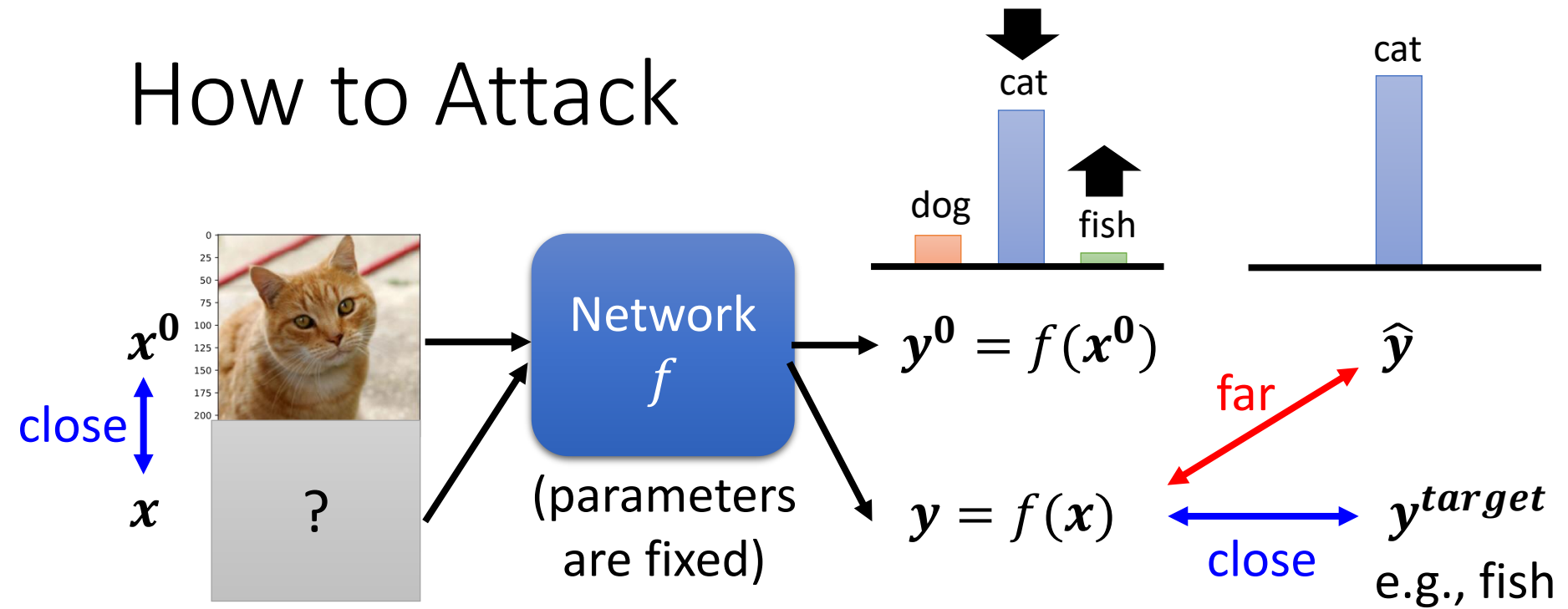Attacked Image



Keyboard
0.98

tiger cat

Persian cat

tabby cat

fire screen

# How to Attack



$$y^0 = f(x^0)$$

$$y = f(x)$$

$$\widehat{y}$$

$$y^{target}$$

e.g., fish

**_Non-targeted_**

$$x^* = arg\min_{} L(x)$$

not perceived by humans

$$L(x) = -e(y, \widehat{y})$$

**_Targeted_**

$$L(x) = -e(y, \widehat{y}) + e(y, y^{target})$$

# Non-perceivable

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix} - \begin{bmatrix} x_1^0 \\ x_2^0 \\ x_3^0 \\ \vdots \end{bmatrix} = \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \Delta x_3 \\ \vdots \end{bmatrix}$$

$$\boldsymbol{x} \qquad \boldsymbol{x^0} \qquad \Delta \boldsymbol{x}$$

$$d(\boldsymbol{x^0}, \boldsymbol{x}) \leq \varepsilon$$

Need to consider human perception

- L2-norm

$$d(\boldsymbol{x^0}, \boldsymbol{x}) = \|\Delta \boldsymbol{x}\|_2$$
$$= (\Delta x_1)^2 + (\Delta x_2)^2 + (\Delta x_3)^2 \cdots$$

- L-infinity

$$d(\boldsymbol{x^0}, \boldsymbol{x}) = \|\Delta \boldsymbol{x}\|_\infty$$
$$= max\{|\Delta x_1|, |\Delta x_2|, |\Delta x_3|, \dots\}$$

small L-∞

Change every pixel a little bit

same L2

Change one pixel much

large L-∞

# Attack Approach

$$w^*, b^* = arg \min_{w,b} L$$ Difference?

Update **input**, not **parameters**

$$\boldsymbol{x}^* = arg \min_{d(x^0,x) \le \varepsilon} L(\boldsymbol{x})$$

## *Gradient Descent*

Start from original image $\boldsymbol{x^0}$

For $t = 1$ to $T$

$$\boldsymbol{x^t} \leftarrow \boldsymbol{x^{t-1}} - \eta \boldsymbol{g}$$

$$\boldsymbol{g} = \begin{bmatrix} \frac{\partial L}{\partial x_1} \big|_{x=x^{t-1}} \\ \frac{\partial L}{\partial x_2} \big|_{x=x^{t-1}} \\ \vdots \end{bmatrix}$$

# Attack Approach

$$w^*, b^* = arg \min_{w,b} L \quad \text{Difference?}$$

Update ***input***, not ***parameters***

Different optimization methods

$$x^* = arg \boxed{\min_{d(x^0, x) \le \varepsilon}} L(x)$$

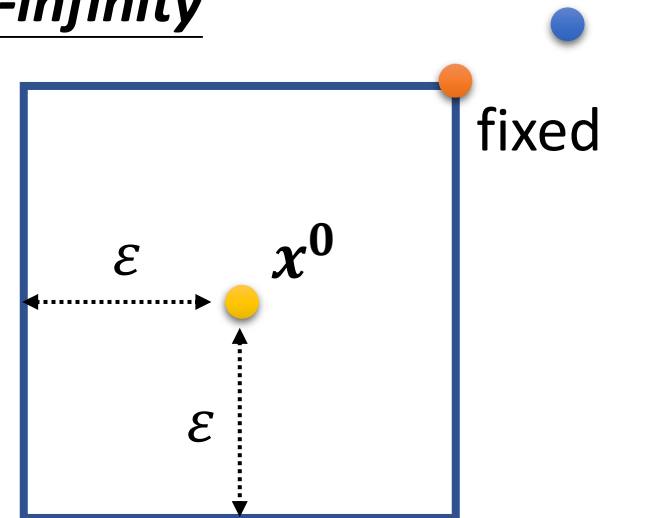Different constraints

**_Gradient Descent_**

Start from original image $x^0$

For $t = 1$ to $T$

$$x^t \leftarrow x^{t-1} - \eta g$$

If $d(x^0, x) > \varepsilon$

$$x^t \leftarrow fix(x^t)$$

**_L-infinity_**

after update

$x^0$

$\varepsilon$

$\varepsilon$

fixed

# Attack Approach

$$x^* = arg \min_{d(x^0, x) \leq \varepsilon} L(x)$$

**Fast Gradient Sign Method (FGSM)**

https://arxiv.org/abs/1412.6572

Start from original image $x^0$

For $t = 1$ ~~to $T$~~

$$x^t \leftarrow x^{t-1} - \eta g$$

# Attack Approach

$$x^* = arg \min_{d(x^0, x) \leq \varepsilon} L(x)$$

**Fast Gradient Sign Method (FGSM)**

https://arxiv.org/abs/1412.6572

Start from original image $x^0$

For $t = 1$ ~~to $T$~~

$$x^t \leftarrow x^{t-1} - \eta g$$

$$\begin{bmatrix} +1 \\ -1 \\ +1 \\ \vdots \end{bmatrix}$$

$\varepsilon$

$$g = \begin{bmatrix} \pm 1 \left[ sign\left(\frac{\partial L}{\partial x_1} |_{x=x^{t-1}}\right) \right] \\ \pm 1 \left[ sign\left(\frac{\partial L}{\partial x_2} |_{x=x^{t-1}}\right) \right] \\ \vdots \end{bmatrix}$$

$$if \ t > 0, sign(t) = 1; otherwise, sign(t) = -1$$

# Attack Approach

**_L-infinity_**

$$x^* = arg \min_{d(x^0, x) \le \varepsilon} L(x)$$

**Iterative FGSM**

https://arxiv.org/abs/1607.02533

Start from original image $x^0$

For $t = 1$ ~~to $T$~~

$\quad x^t \leftarrow x^{t-1} - \eta g$

$\quad$ If $d(x^0, x) > \varepsilon$

$\qquad\qquad x^t \leftarrow fix(x^t)$

$$g = \begin{bmatrix} \pm 1 \left[ sign\left( \frac{\partial L}{\partial x_1} \big|_{x=x^{t-1}} \right) \right] \\ \pm 1 \left[ sign\left( \frac{\partial L}{\partial x_2} \big|_{x=x^{t-1}} \right) \right] \\ \vdots \end{bmatrix}$$

fixed

$\varepsilon$

$x^0$

$\varepsilon$

# White Box v.s. Black Box

$$g = \begin{bmatrix} sign\left(\frac{\partial L}{\partial x_1}|_{x=x^{t-1}}\right) \\ sign\left(\frac{\partial L}{\partial x_2}|_{x=x^{t-1}}\right) \\ \vdots \end{bmatrix}$$
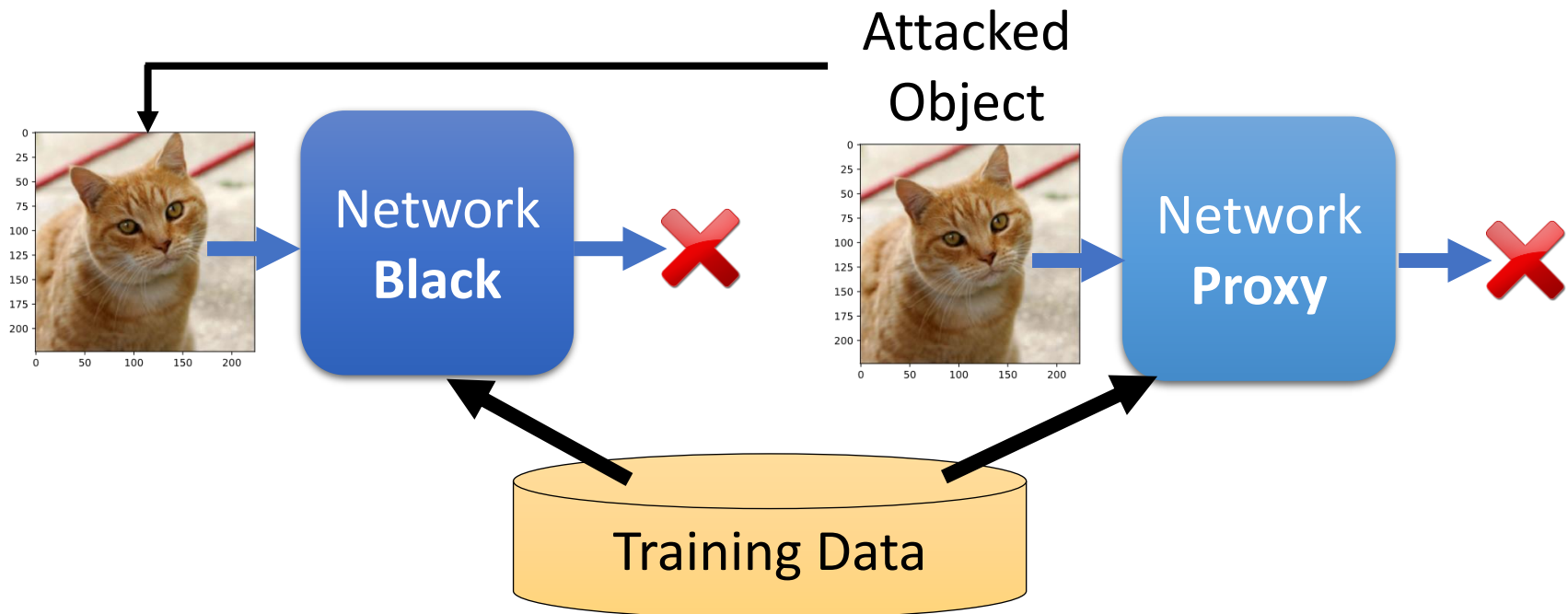
- In the previous attack, we know the network parameters $\theta$
  - This is called **White Box Attack**.

- You cannot obtain model parameters in most online API.

- Are we safe if we do not release model? ☺

- No, because **Black Box Attack** is possible. ☹

# Black Box Attack

If you have the training data of the target network
Train a proxy network yourself
Using the proxy network to generate attacked objects



Attacked Object

Network **Black**

Network **Proxy**

Training Data

What if we do not know the training data?

# Black Box Attack

Be Attacked

|  | ResNet-152 | ResNet-101 | ResNet-50 | VGG-16 | GoogLeNet |
|---|---|---|---|---|---|
| ResNet-152 | 0% | 13% | 18% | 19% | 11% |
| ResNet-101 | 19% | 0% | 21% | 21% | 12% |
| ResNet-50 | 23% | 20% | 0% | 21% | 18% |
| VGG-16 | 22% | 17% | 17% | 0% | 5% |
| GoogLeNet | 39% | 38% | 34% | 19% | 0% |

Proxy

(lower accuracy → more successful attack)

## *Ensemble Attack*

|  | ResNet-152 | ResNet-101 | ResNet-50 | VGG-16 | GoogLeNet |
|---|---|---|---|---|---|
| -ResNet-152 | 0% | 0% | 0% | 0% | 0% |
| -ResNet-101 | 0% | 1% | 0% | 0% | 0% |
| -ResNet-50 | 0% | 0% | 2% | 0% | 0% |
| -VGG-16 | 0% | 0% | 0% | 6% | 0% |
| -GoogLeNet | 0% | 0% | 0% | 0% | 5% |

# The attack is so easy! Why?



https://arxiv.org/pdf/1611.02770.pdf

To learn more:

Adversarial Examples Are Not
Bugs, They Are Features

https://arxiv.org/abs/1905.02175

# One pixel attack

joystick



Cup(16.48%)
Soup Bowl(16.74%)

Bassinet(16.59%)
Paper Towel(16.21%)

Teapot(24.99%)
Joystick(37.39%)

Hamster(35.79%)
Nipple(42.36%)

Video: https://youtu.be/tfpKIZIWidA

# Universal Adversarial Attack

https://arxiv.org/abs/1610.08401



Black Box Attack is also possible!

# Beyond Images

- Speech processing

Detect synthesized speech

Synthesized!

Real!

- Natural language processing

https://arxiv.org/abs/1908.07125

*Question:* Why did he walk?
For <u>exercise</u>, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells.

exercise

*Question:* Why did the university see a drop in applicants?
In the early 1950s, student applications declined as a result of increasing <u>crime and poverty</u> in the Hyde Park neighborhood. In response, the university became a . . . . .

crime and poverty

# Attack in the Physical World



- An attacker would need to find perturbations that generalize beyond a single image.

- Extreme differences between adjacent pixels in the perturbation are unlikely to be accurately captured by cameras.

- It is desirable to craft perturbations that are comprised mostly of colors reproducible by the printer.

| Distance/Angle | Subtle Poster | Subtle Poster Right Turn | Camouflage Graffiti | Camouflage Art (LISA-CNN) | Camouflage Art (GTSRB-CNN) |
|---|---|---|---|---|---|
| 5′ 0° | | | | | |
| 5′ 15° | | | | | |
| 10′ 0° | | | | | |
| 10′ 30° | | | | | |
| 40′ 0° | | | | | |
| Targeted-Attack Success | 100% | 73.33% | 66.67% | 100% | 80% |

https://arxiv.org/abs/1707.08945

# Attack in the Physical World



read as an 85-mph sign

https://youtu.be/4uGV_fRj0UA

https://www.mcafee.com/blogs/other-blogs/mcafee-labs/model-hacking-adas-to-pave-safer-roads-for-autonomous-vehicles/
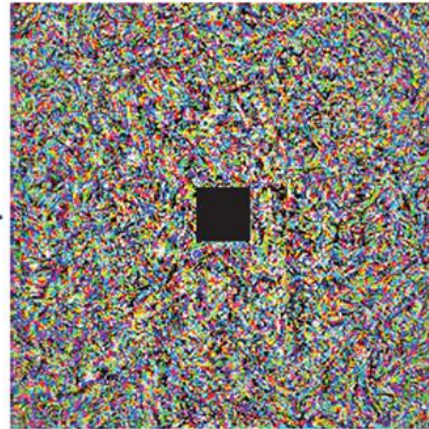
# Adversarial Reprogramming

(a)

| counting $y_{adv}$ | ImageNet $y$ |
|---|---|
| 1 square | tench |
| 2 squares | goldfish |
| 3 squares | white shark |
| 4 squares | tiger shark |
| 5 squares | hammerhead |
| 6 squares | electric ray |
| 7 squares | stingray |
| 8 squares | cock |
| 9 squares | hen |
| 10 squares | ostrich |

(b) Adversarial Program

(c)

ImageNet Classifier

tiger shark, ostrich
≡
4 squares, 10 squares

https://arxiv.org/abs/1806.11146

# "Backdoor" in Model

- Attack happens at the training phase

Training data

\+

dog

(attacked)

Goal: misclassified as "dog"

Train

Model → dog!

be careful of unknown dataset ......

# Defense
## Passive v.s. Proactive

# Passive Defense

Original

Do not influence classification



+

Filter

e.g. Smoothing

+

Network

Tiger Cat
~~Keyboard~~

Attack signal

Less harmful

Keyboard
0.98

Smoothing

tiger cat
0.37

tiger cat
0.64

Smoothing

tiger cat
0.45    Side Effect!

# Passive Defense

## Image Compression



8.9M           68.34K

https://arxiv.org/abs/1704.01155
https://arxiv.org/abs/1802.06816

## Generator

https://arxiv.org/abs/1805.06605



Input
image

G

# Passive Defense - Randomization



Input Image $X_n$

Resized Image $X_n'$

Padded Image $X_n''$

CNN

**Random Resizing Layer**

**Random Padding Layer**

**Randomly Select One Pattern**

**CNN Classification**

https://arxiv.org/abs/1711.01991

# Proactive Defense

Adversarial Training for Free!

https://arxiv.org/abs/1904.12843

_Adversarial Training_

Training a model that is robust to adversarial attack.

Given training set $\mathcal{X} = \{(\boldsymbol{x^1}, \hat{y}^1), (\boldsymbol{x^2}, \hat{y}^2), \cdots, (\boldsymbol{x^N}, \hat{y}^y)\}$

Using $\mathcal{X}$ to train your model

For $n = 1$ to $N$

> Can it deal with new algorithm?

Find adversarial input $\widetilde{\boldsymbol{x}}^{\boldsymbol{n}}$ given $\boldsymbol{x^n}$ by an attack algorithm

Find the problem

We have new training data

$$\mathcal{X}' = \{(\widetilde{\boldsymbol{x}}^{\boldsymbol{1}}, \hat{y}^1), (\widetilde{\boldsymbol{x}}^{\boldsymbol{2}}, \hat{y}^2), \cdots, (\widetilde{\boldsymbol{x}}^{\boldsymbol{N}}, \hat{y}^y)\}$$

Using both $\mathcal{X}$ and $\mathcal{X}'$ to update your model   Fix it!

**Data Augmentation**

# Concluding Remarks

- Attack: given the network parameters, attack is very easy.

- Even black box attack is possible

- Defense: Passive & Proactive
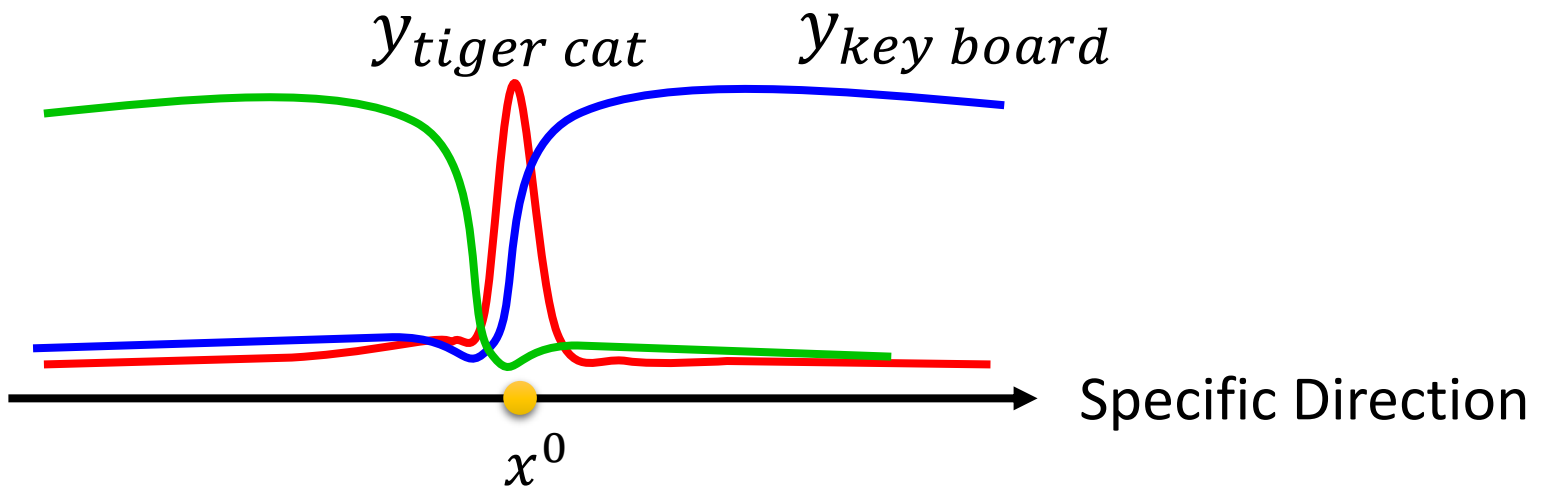
- Attack / Defense are still evolving.

# Acknowledgement

- 感謝作業十助教團隊林毓宸同學、黃啟斌同學幫忙蒐集參考

# Attack Approaches

- FGSM (https://arxiv.org/abs/1412.6572)

- Basic iterative method (https://arxiv.org/abs/1607.02533)

- L-BFGS (https://arxiv.org/abs/1312.6199)

- Deepfool (https://arxiv.org/abs/1511.04599)

- JSMA (https://arxiv.org/abs/1511.07528)

- C&W (https://arxiv.org/abs/1608.04644)

- Elastic net attack (https://arxiv.org/abs/1709.04114)

- Spatially Transformed (https://arxiv.org/abs/1801.02612)

- One Pixel Attack (https://arxiv.org/abs/1710.08864)

- …… only list a few

# What happened?

# 李宏毅 (Hung-yi Lee) · HYLEE | Machine Learning (2021)

## HYLEE(2021)· 课程资料包 @ShowMeAI

**视频**
中英双语字幕

**课件**
一键打包下载

**笔记**
官方笔记翻译

**代码**
作业项目解析

**视频·B 站 [ 扫码或点击链接 ]**
https://www.bilibili.com/video/BV1fM4y137M4

**课件 & 代码·博客 [ 扫码或点击链接 ]**
http://blog.showmeai.tech/ntu-hylee-ml

机器学习　Auto-encoder　生成式对抗网络　学习率
深度学习　卷积神经网络　GAN　自监督　自注意力机制
批次标准化　神经网络压缩　强化学习　元学习　Transformer

Awesome AI Courses Notes Cheatsheets 是 **ShowMeAI** 资料库的分支系列，覆盖最具知名度的 **TOP50+** 门 AI 课程，旨在为读者和学习者提供一整套高品质中文学习笔记和速查表。

**点击**课程名称，跳转至课程**资料包**页面，**一键下载**课程全部资料！

| 机器学习 | 深度学习 | 自然语言处理 | 计算机视觉 |
|---|---|---|---|
| Stanford · CS229 | Stanford · CS230 | Stanford · CS224n | Stanford · CS231n |

### # Awesome AI Courses Notes Cheatsheets· 持续更新中

| 知识图谱 | 图机器学习 | 深度强化学习 | 自动驾驶 |
|---|---|---|---|
| Stanford · CS520 | Stanford · CS224W | UCBerkeley · CS285 | MIT · 6.S094 |

**微信公众号**

资料下载方式 2：扫码点击底部菜单栏

称为 **AI 内容创作者？** 回复 [ 添砖加瓦 ]