

COMP5212: Machine Learning

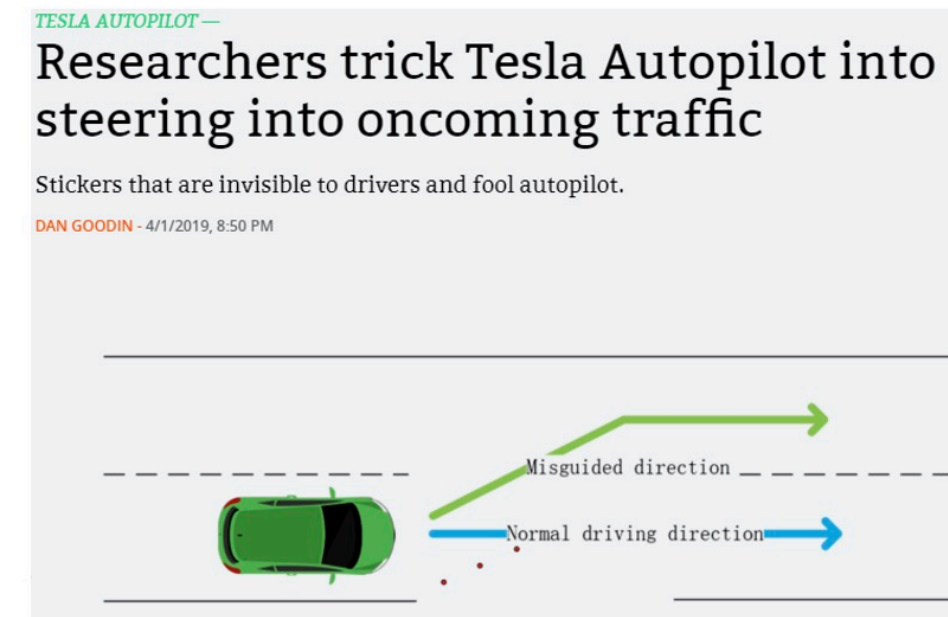
Lecture 18: Adversarial Machine Learning

Sen LI

Zeyu QIN

Machine learning

Beyond Accuracy



Microsoft silences its new A.I. bot Tay, after Twitter users teach it racism [Updated]

Sarah Perez @sarahintampa / 10:16 am EDT • March 24, 2016 [Comment](#)

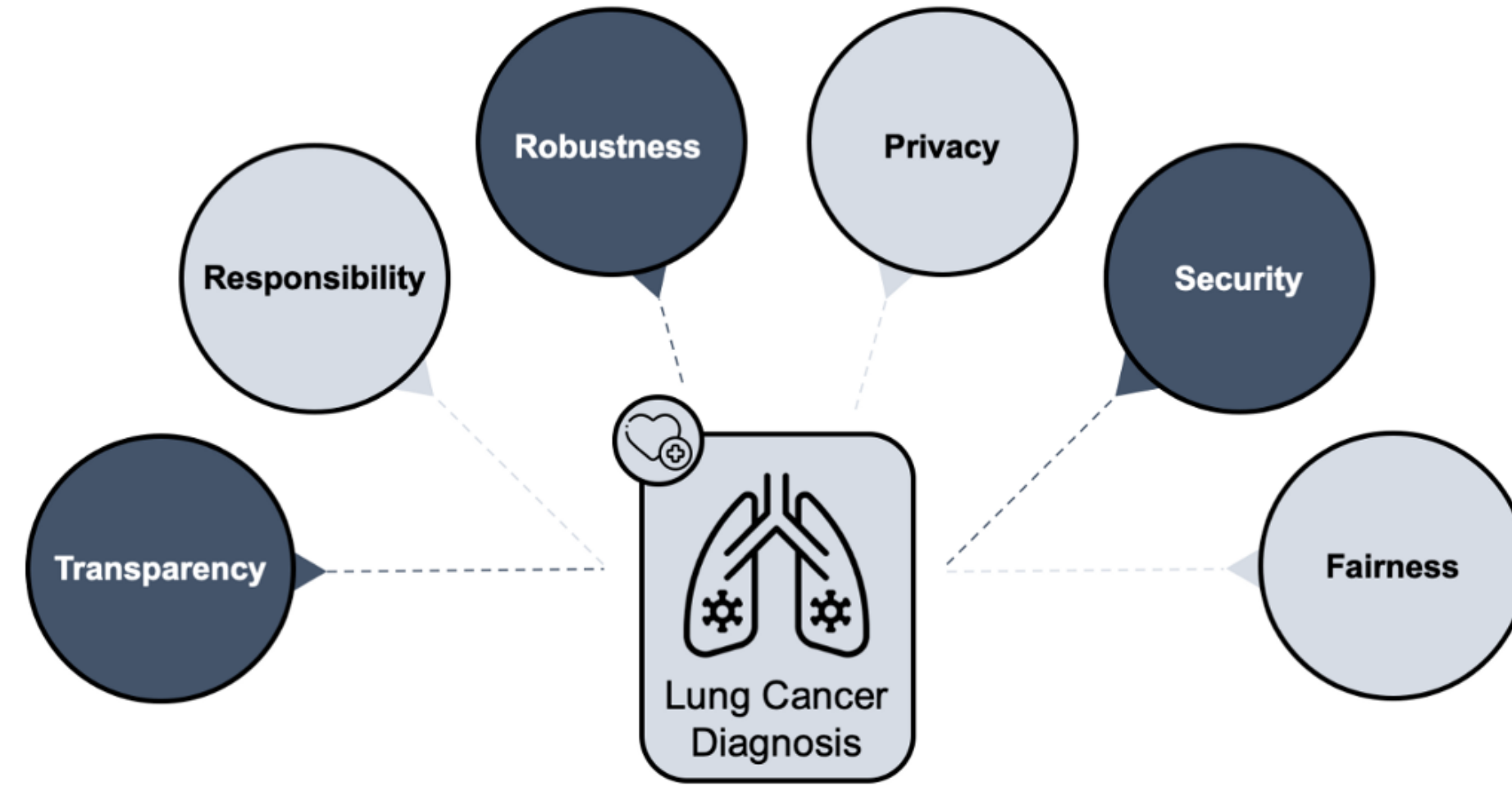


Microsoft's newly launched A.I.-powered bot called Tay, which was responding to tweets and chats on GroupMe and Kik, has already been shut down due to concerns with its inability to recognize when it was making offensive or racist statements. Of course, the bot wasn't coded to be racist, but it "learns" from those it interacts with. And naturally, given that this is the Internet, one of the first things online users taught Tay was how to be racist, and how to spout back ill-informed or inflammatory political opinions. [Update: Microsoft now says it's "making adjustments" to Tay in light of this problem.]

Trustworthy ML

What and why

- Not alchemy
 - Explainability
 - Security
 - Privacy
 - Fairness
 - **Integrity**
 - ...
- Establish model understanding

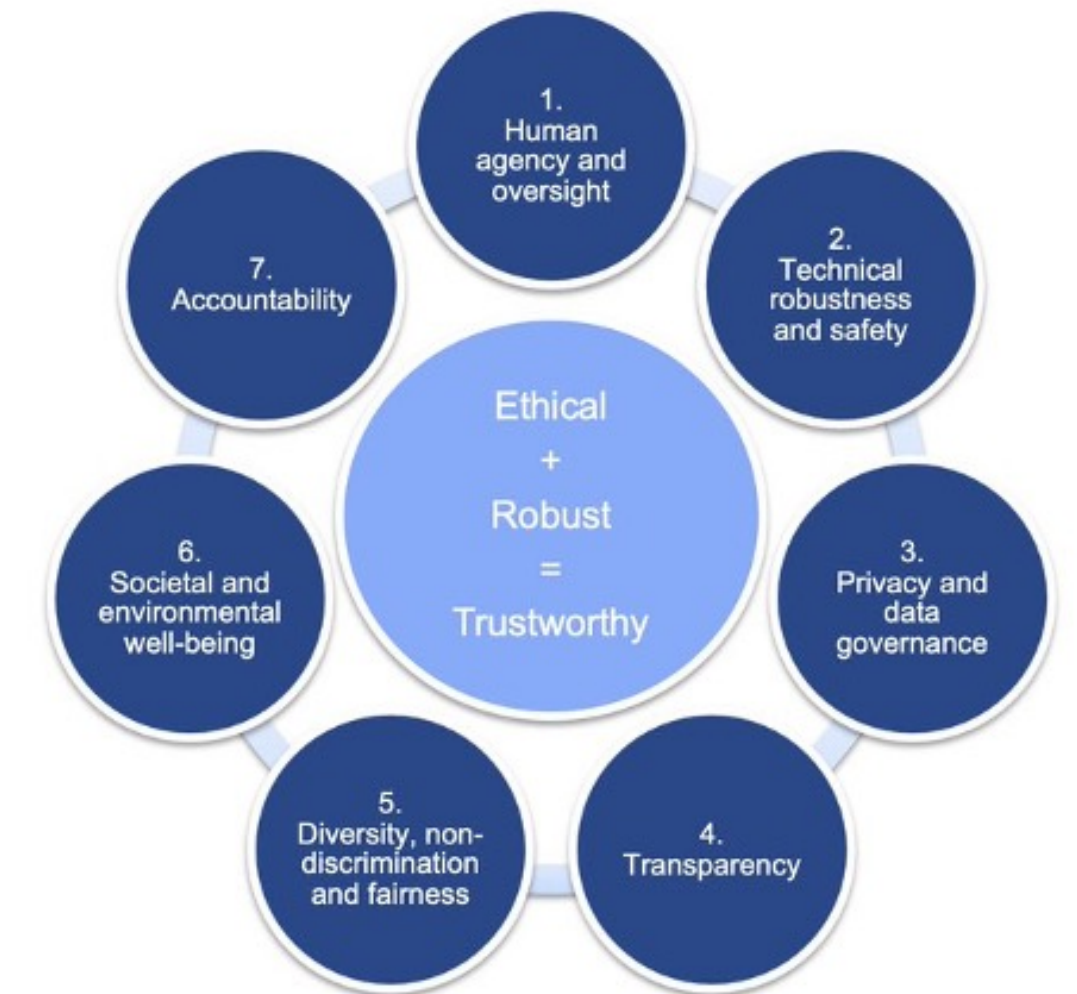


THE NATIONAL SECURITY COMMISSION
ON ARTIFICIAL INTELLIGENCE

人工智能安全测评白皮书
(2021)



国家语音及图像识别产品质量监督检验中心
国家工业信息安全发展研究中心人工智能所
2021年10月



Trustworthy ML

Integrity

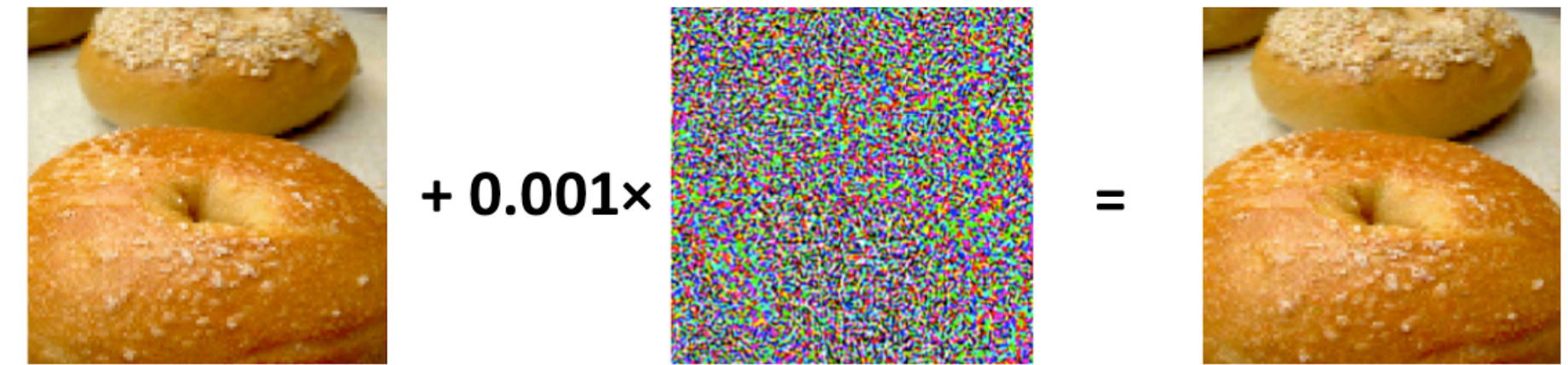
- Training-time integrity and Testing-time integrity

Attack Category	Attack Target	Attack Mechanism	Training Process	Inference Process
Backdoor Attack	Misclassify attacked samples; Behave normal on benign samples.	Excessive learning ability of models.	Under control.	Out of control.
Adversarial Attack	Misclassify attacked samples; Behave normal on benign samples.	Behavior differences between models and humans.	Out of control.	Attackers need to generate adversarial perturbation through an iterative optimization process.
Data Poisoning	Reduce model generalization.	Overfitting to bad local optima.	Can only modify the training set.	Out of control.

Testing time integrity

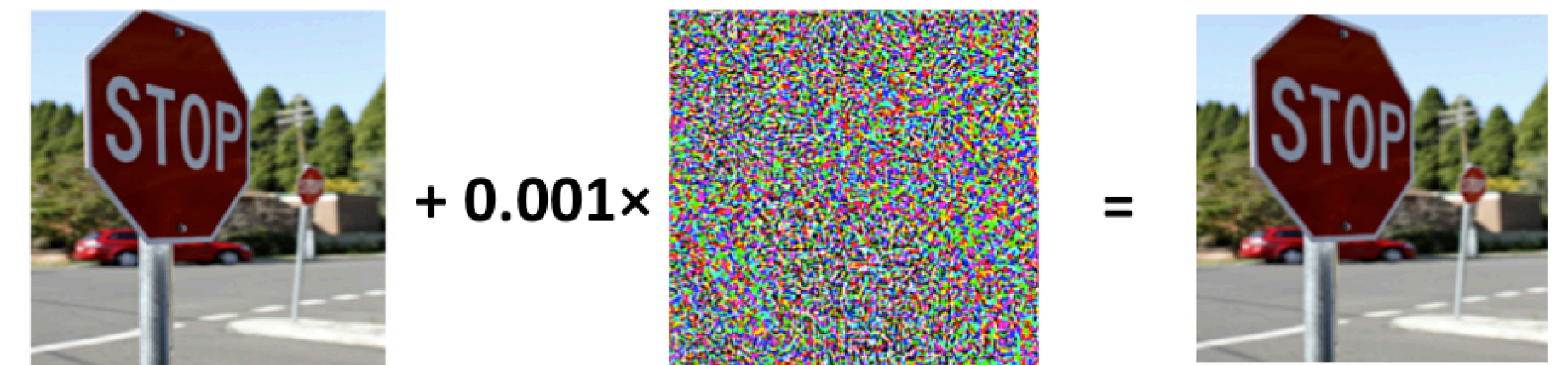
Adversarial examples

- An **adversarial** example can easily fool a deep network
- **Robustness** is critical in real systems



Bagel

piano



stop sign

speed limit 40

Adversarial examples

Definition

- Given a K -way multi-class classification model $f : \mathbb{R}^d \rightarrow \{1, \dots, K\}$ and an original example x_0 , the goal is to generate an adversarial example x such that
 - x is close to x_0 and $\arg \max_i f_i(x) \neq \arg \max_i f_i(x_0)$
 - i.e., x has a different prediction with x_0 by model f .

Adversarial example

Attack as an optimization problem

- Craft adversarial example by solving

- $\arg \min_x \|x - x_0\|^2 + c \cdot h(x)$

- $\|x - x_0\|^2$: the distortion

Adversarial example

Attack as an optimization problem

- Craft adversarial example by solving

- $\arg \min_x \|x - x_0\|^2 + c \cdot h(x)$

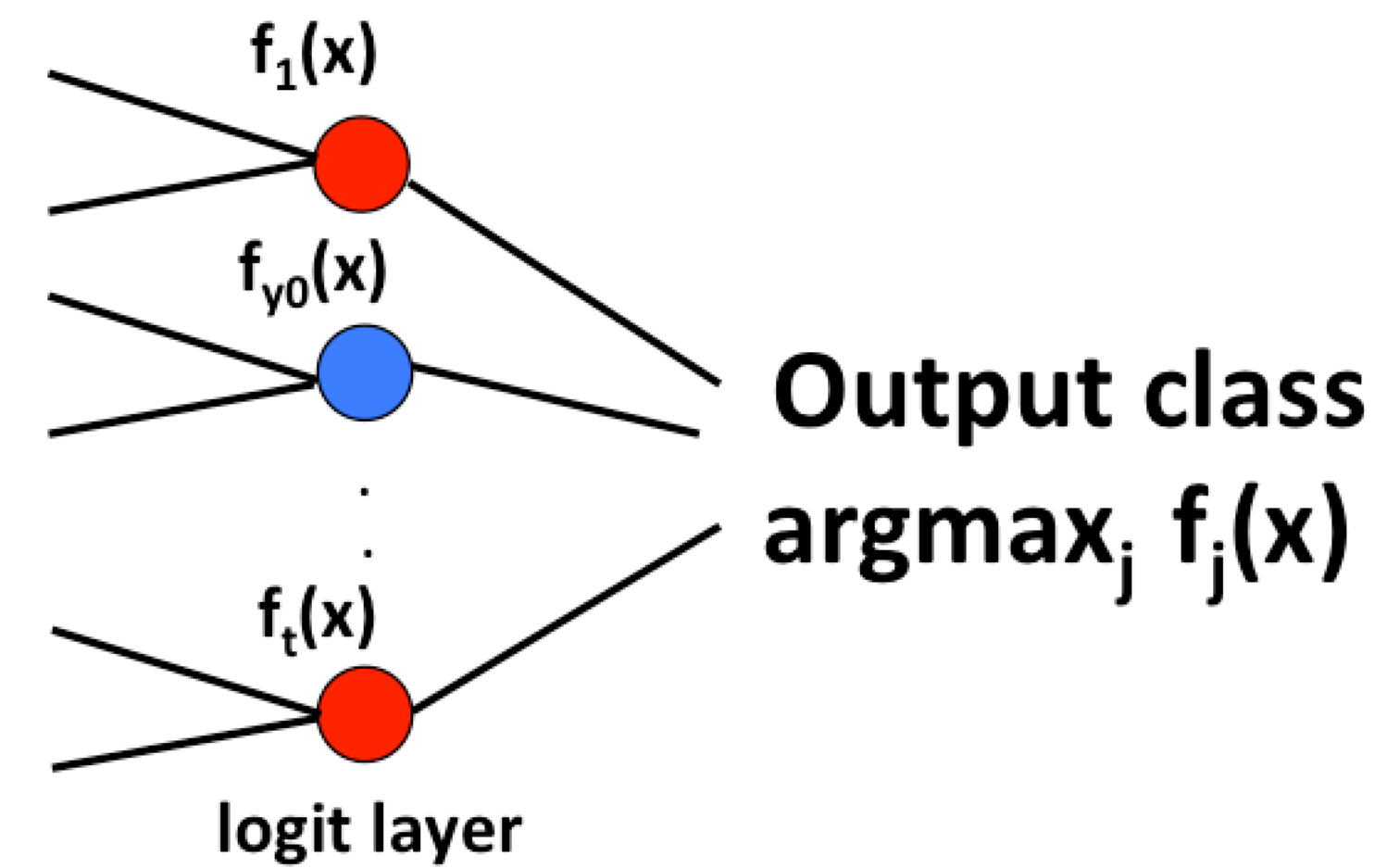
- $\|x - x_0\|^2$: the distortion

- $h(x)$: loss to measure the **successfulness** of attack

Adversarial example

Attack as an optimization problem

- Craft adversarial example by solving
 - $\arg \min_x \|x - x_0\|^2 + c \cdot h(x)$
- $\|x - x_0\|^2$: the distortion
- $h(x)$: loss to measure the **successfulness** of attack
- Untargeted attack: success if $\arg \max_j f_j(x) \neq y_0$
 - $h(x) = \max\{f_{y_0}(x) - \max_{j \neq y_0} f_j(x), 0\}$



How to find adversarial examples

White-box vs black-box setting

- Attackers knows the model structure and weights (white-box)
- Can query the model to get probability output (soft-label)
- Can query the model to get label output (hard-label)
- No information about the model (universal)

Adversarial example

White-box setting

- $\arg \min_x \|x - x_0\|^2 + c \cdot h(x)$
- Model (network structure and weights) is revealed to attacker
 - \Rightarrow gradient of $h(x)$ can be computed
 - \Rightarrow attacker minimizes the objective by gradient descent

Adversarial example

White-box adversarial attack

- C&W attack [CW17]:
 - $h(x) = \max\{ [Z_{y_0}(x) - \max_{j \neq y} Z_j(x)], -\kappa \}$
 - Where $Z(x)$ is the pre-softmax layer output

Adversarial example

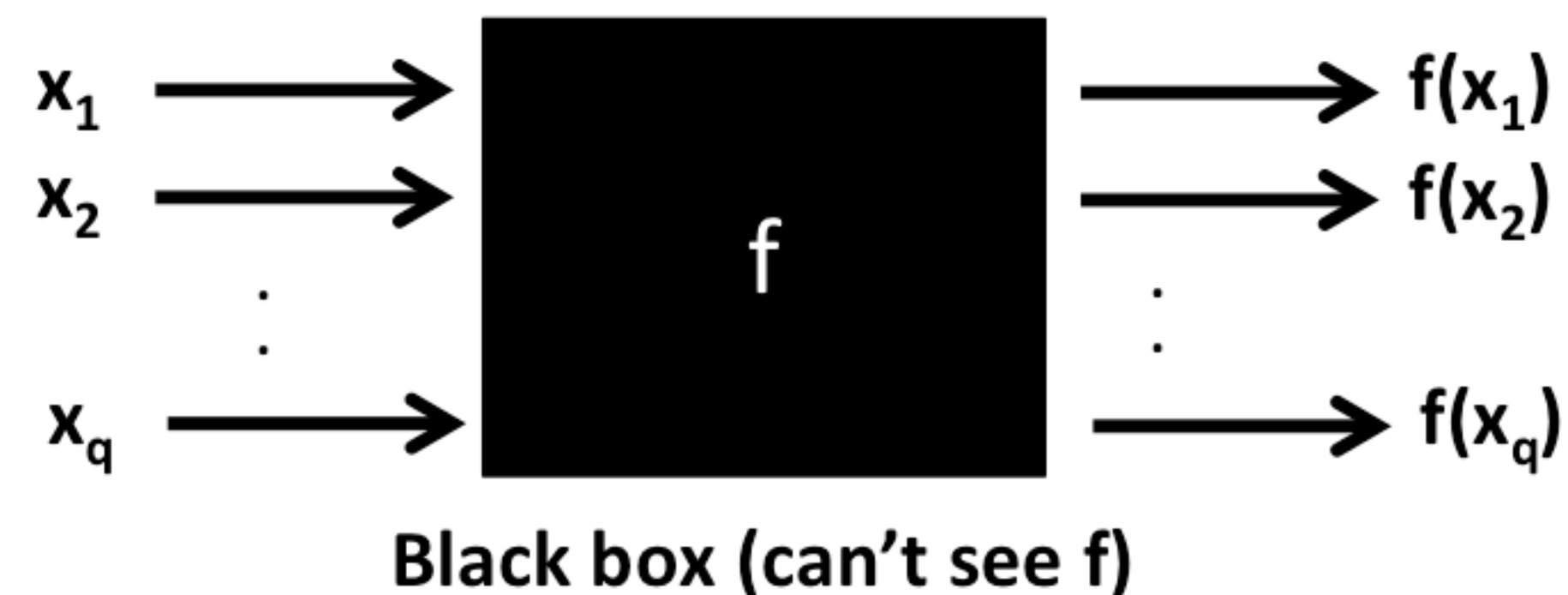
White-box adversarial attack

- If there is $\|x - x_0\|_\infty$ constraint, we could turn to solve by
- FGSM attack [GSS15]:
 - $x \leftarrow \text{proj}_{x+\mathcal{S}}(x_0 + \alpha \text{sign}(\nabla_{x_0} \ell(\theta, x, y)))$
- PGD attack [KGB17, MMS18]
 - $x^{t+1} \leftarrow \text{proj}_{x+\mathcal{S}}(x^t + \alpha \text{sign}(\nabla_{x^t} \ell(\theta, x, y)))$

Adversarial example

Black-box Soft-label Setting

- Black-box Soft Label setting (practical setting):
 - Structure and weights of deep network are not revealed to attackers
 - Attacker can **query** the ML model and get the **probability output**



- Cannot compute gradient ∇_x

Soft-label Black-box Adversarial attack

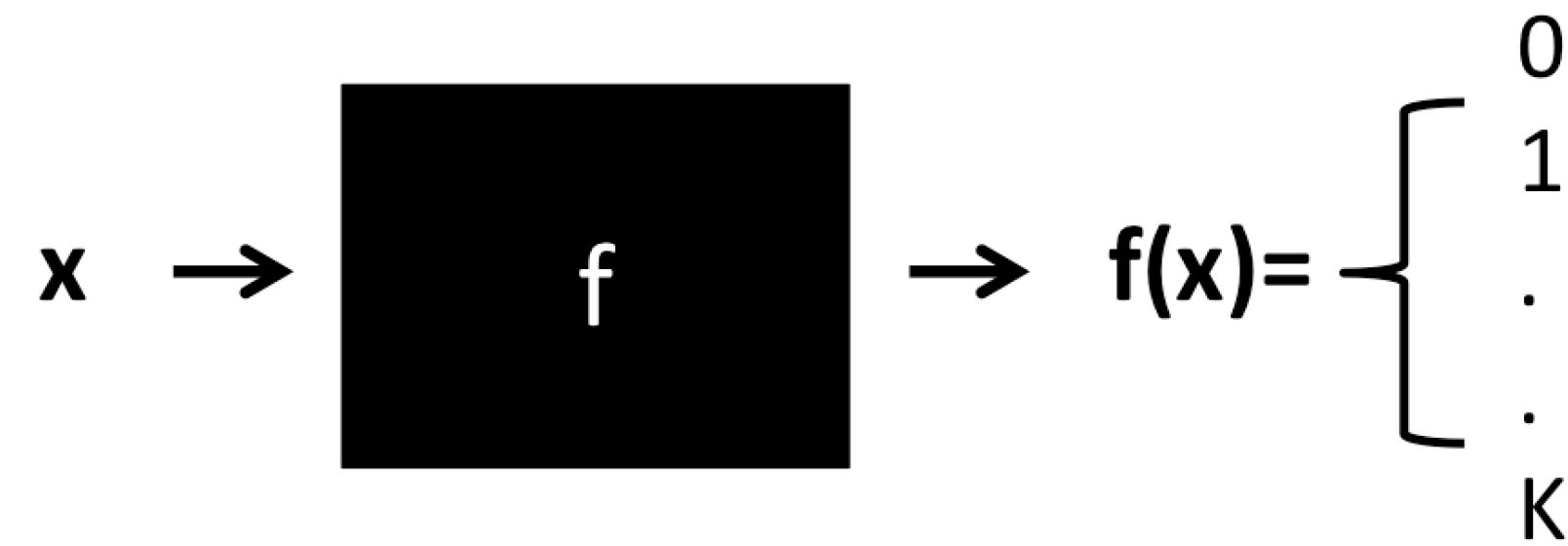
- Soft-label Black-box: query to get the **probability output**
- Key problem: how to estimate gradient?
- Gradient-based [CZS17, IEAL18]:
 - $$\nabla_x = \frac{h(x + \beta u) - h(x)}{\beta} \cdot u$$
- Genetic algorithm [ASC19]

Soft-label Black-box Adversarial attack

- Query based
- Transfer based:
 - Train a substitute model and conduct the white-box attack

Hard-label Black-box Attack

- Model is not known to the attacker
- Attacker can make query and observe **hard-label multi-class output**



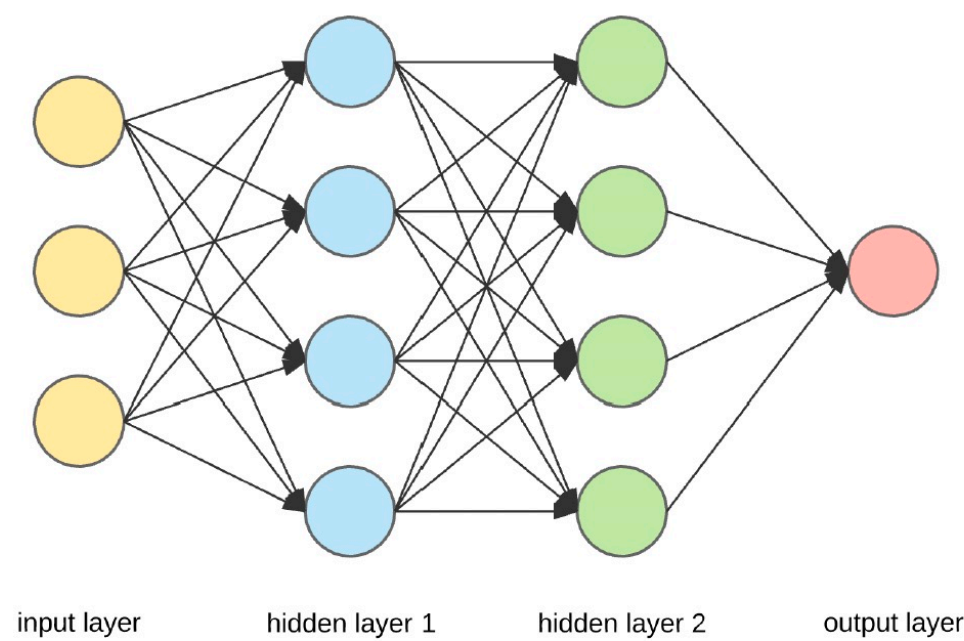
- (K : number of classes)

Why Hard-label

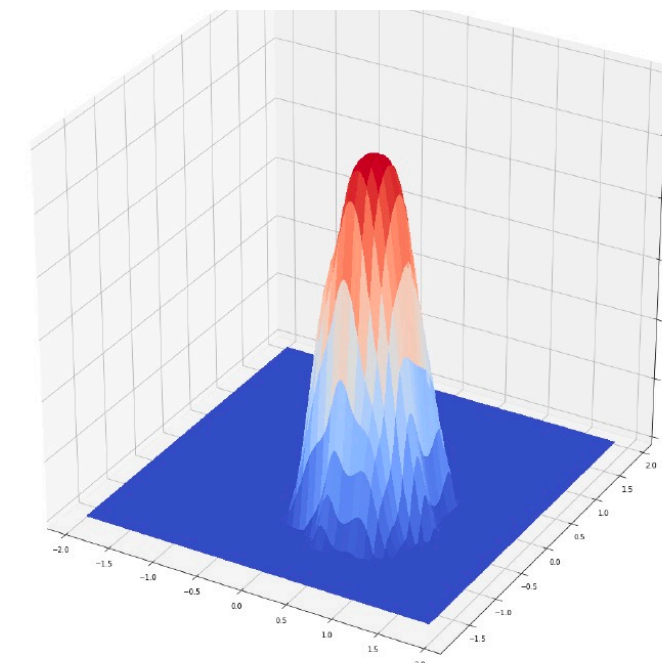
- More practical setting for attacker
- Discrete and complex models (e.g quantization, projection, detection)
- Framework friendly

The difficulty

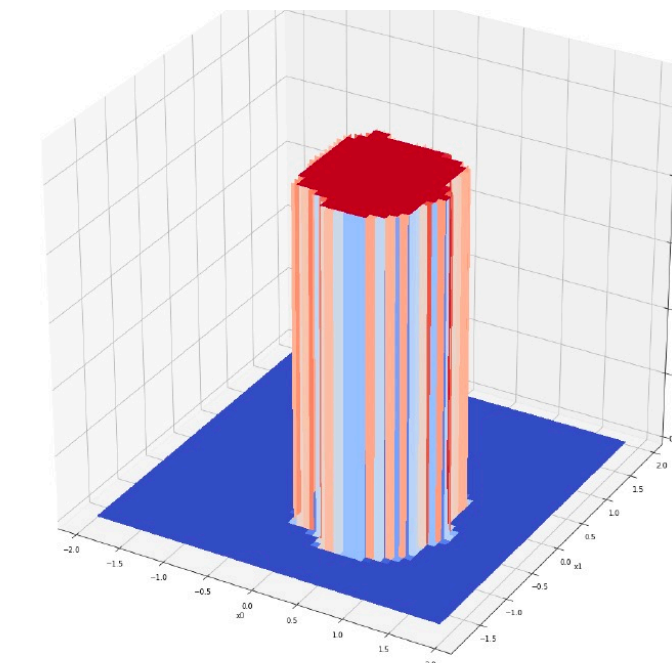
- Hard-label attack on a simple 3-layer neural network yields a discontinuous optimization problem



(a) neural network $f(x)$



(b) $h(Z(x))$



(c) $h(f(x))$

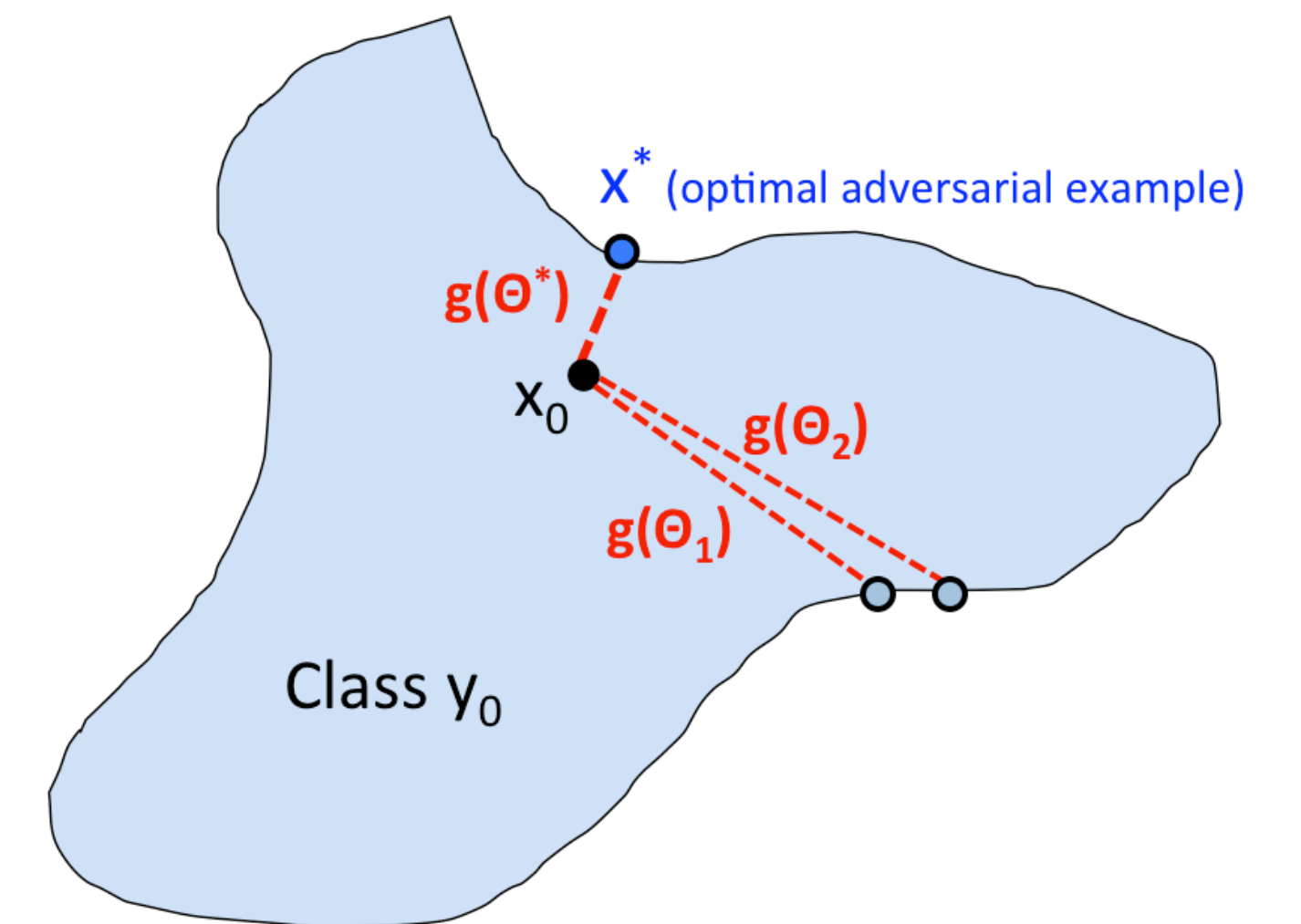
Reformulation

- We reformulate the attack optimization problem (untargeted attack):

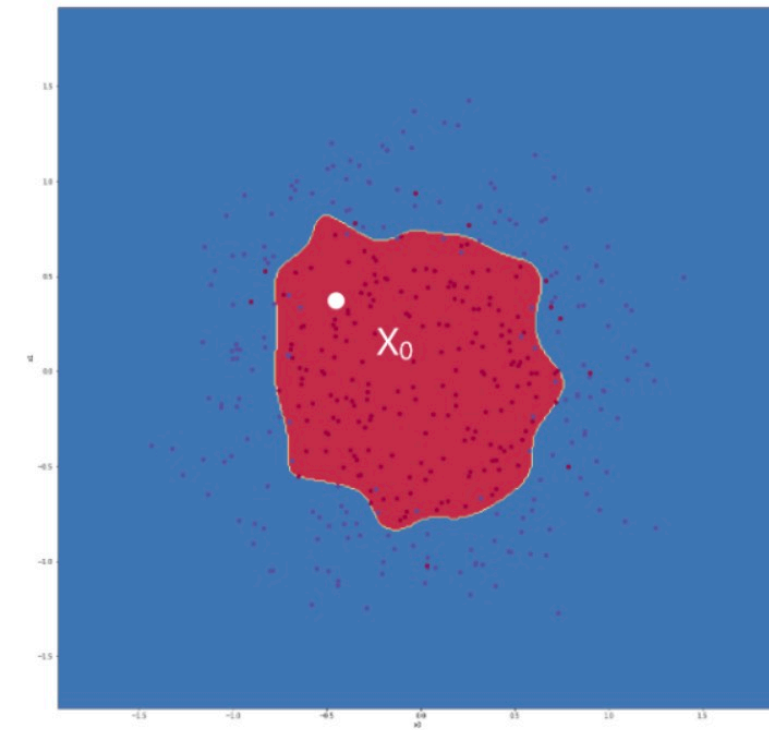
$$\theta^* = \arg \min_{\theta} g(\theta)$$

- where $g(\theta) = \operatorname{argmin}_{\lambda > 0} \left(f\left(x_0 + \lambda \frac{\theta}{\|\theta\|}\right) \neq y_0 \right)$

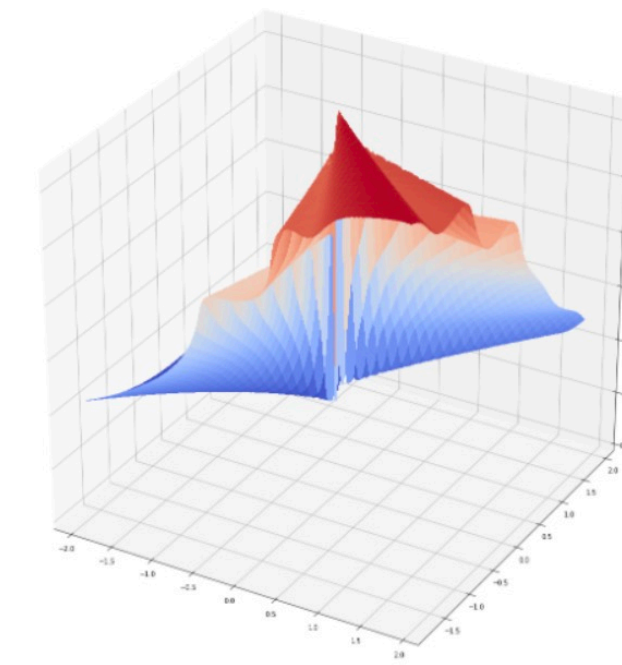
- θ : the direction of adversarial example



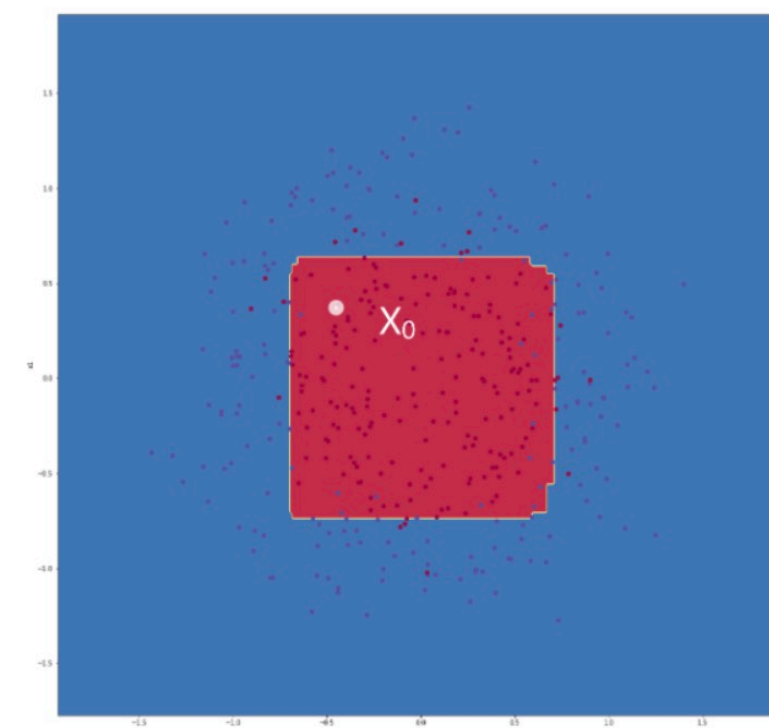
Examples



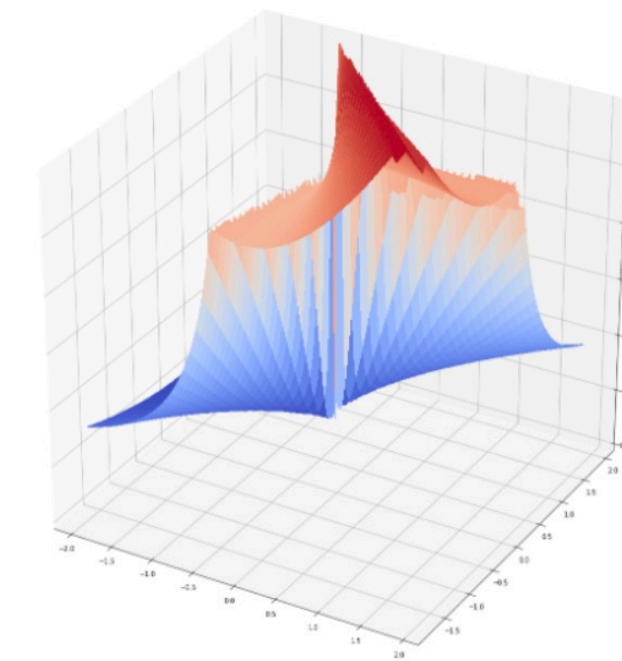
Neural network decision function



$g(\theta)$



Boosting Tree decision function



$g(\theta)$

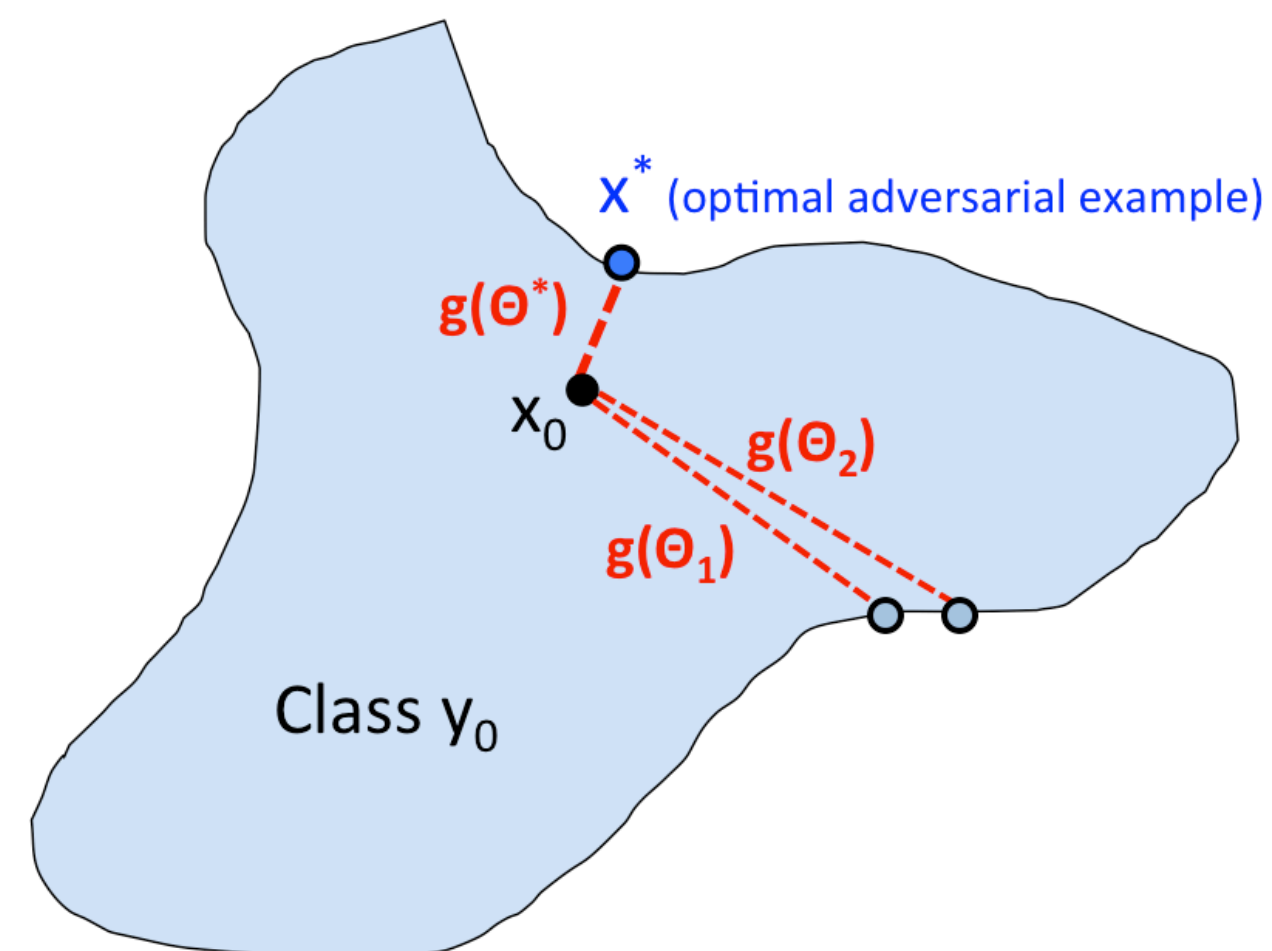
Two things unaddressed

$$\theta^* = \arg \min_{\theta} g(\theta)$$

- where $g(\theta) = \operatorname{argmin}_{\lambda > 0} \left(f\left(x_0 + \lambda \frac{\theta}{\|\theta\|}\right) \neq y_0 \right)$
- How to estimate $g(\theta)$
- How to find θ^*

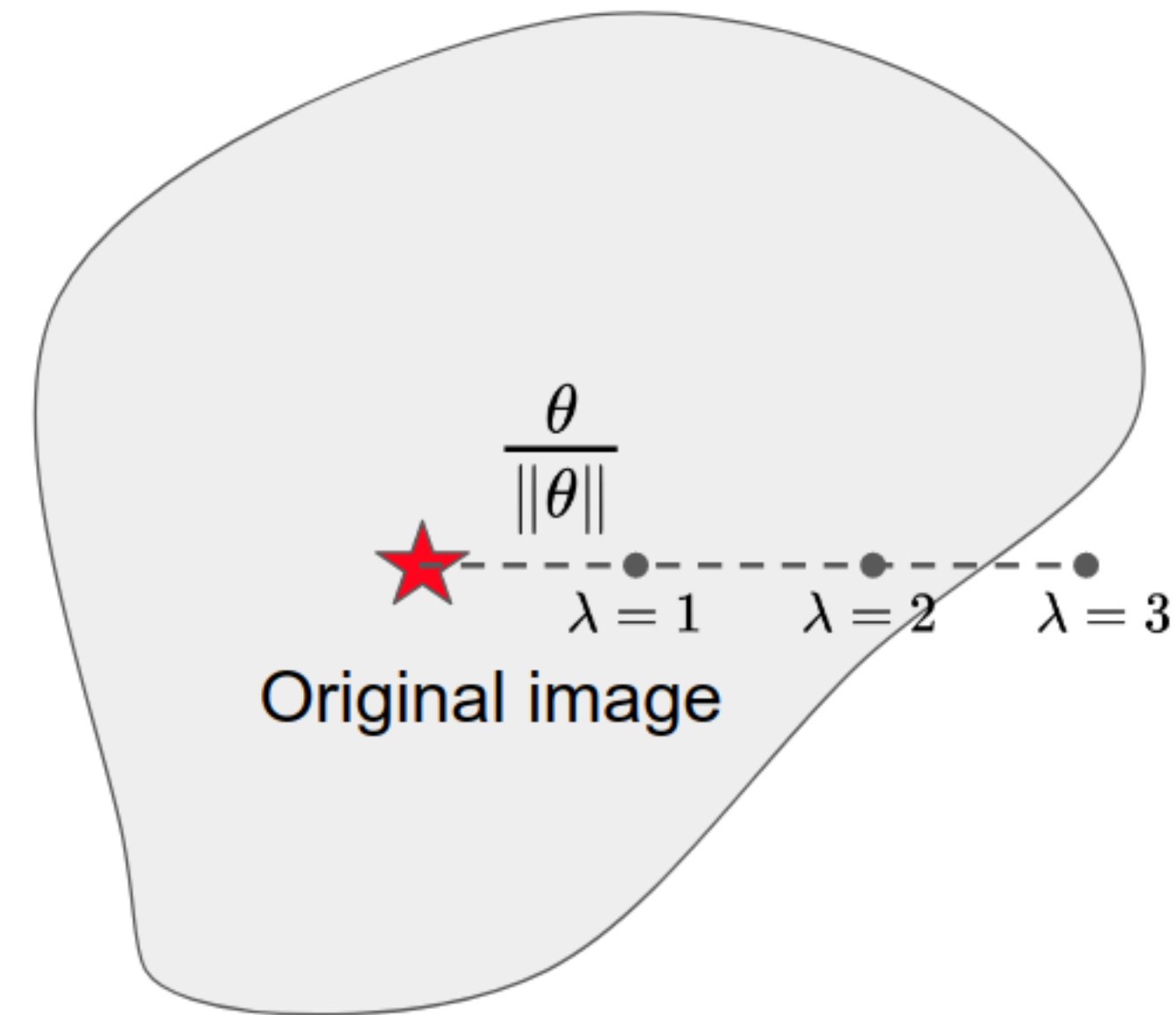
Computing Function Value

- Can't compute the gradient of g
- However, we can compute the function value of g using queries of $f(\cdot)$
- Implemented using fine-grained search + **binary search**



Estimation of $g(\theta)$

- Fine-grained search
- Binary search
- Prediction unchanged enlarge g
- Prediction changed shrink g



How to optimize $g(\theta)$

- The gradient of g is available by

- $$\nabla g(\theta) \approx \frac{g(\theta + \beta u) - g(\theta)}{\beta} \cdot u$$

- One u is too noisy, better to use multiple u (~ 20)
- Zeroth order optimization for minimizing $g(\theta)$

Algorithm

Algorithm 1 OPT attack (ICLR '19)

- 1: **Input:** Hard-label model f , original image x_0 , initial θ_0 .
 - 2: **for** $t = 0, 1, 2, \dots, T$ **do**
 - 3: Randomly choose u from a zero-mean Gaussian distribution
 - 4: Evaluate $g(\theta_t)$ and $g(\theta_t + \beta u)$
 - 5: Compute $\hat{g} = \frac{g(\theta_t + \beta u) - g(\theta_t)}{\beta} \cdot u$
 - 6: Update $\theta_{t+1} = \theta_t - \eta_t \hat{g}$
 - 7: **return** $x_0 + g(\theta_T)\theta_T$
-

Algorithm

Algorithm 2 OPT attack (ICLR '19)

- 1: **Input:** Hard-label model f , original image x_0 , initial θ_0 .
 - 2: **for** $t = 0, 1, 2, \dots, T$ **do**
 - 3: Randomly choose u_t from a zero-mean Gaussian distribution
 - 4: Evaluate $g(\theta_t)$ and $g(\theta_t + \beta u)$
 - 5: Compute $\hat{g} = \frac{g(\theta_t + \beta u) - g(\theta_t)}{\beta} \cdot u$
 - 6: Update $\theta_{t+1} = \theta_t - \eta_t \hat{g}$
 - 7: **return** $x_0 + g(\theta_T)\theta_T$
-

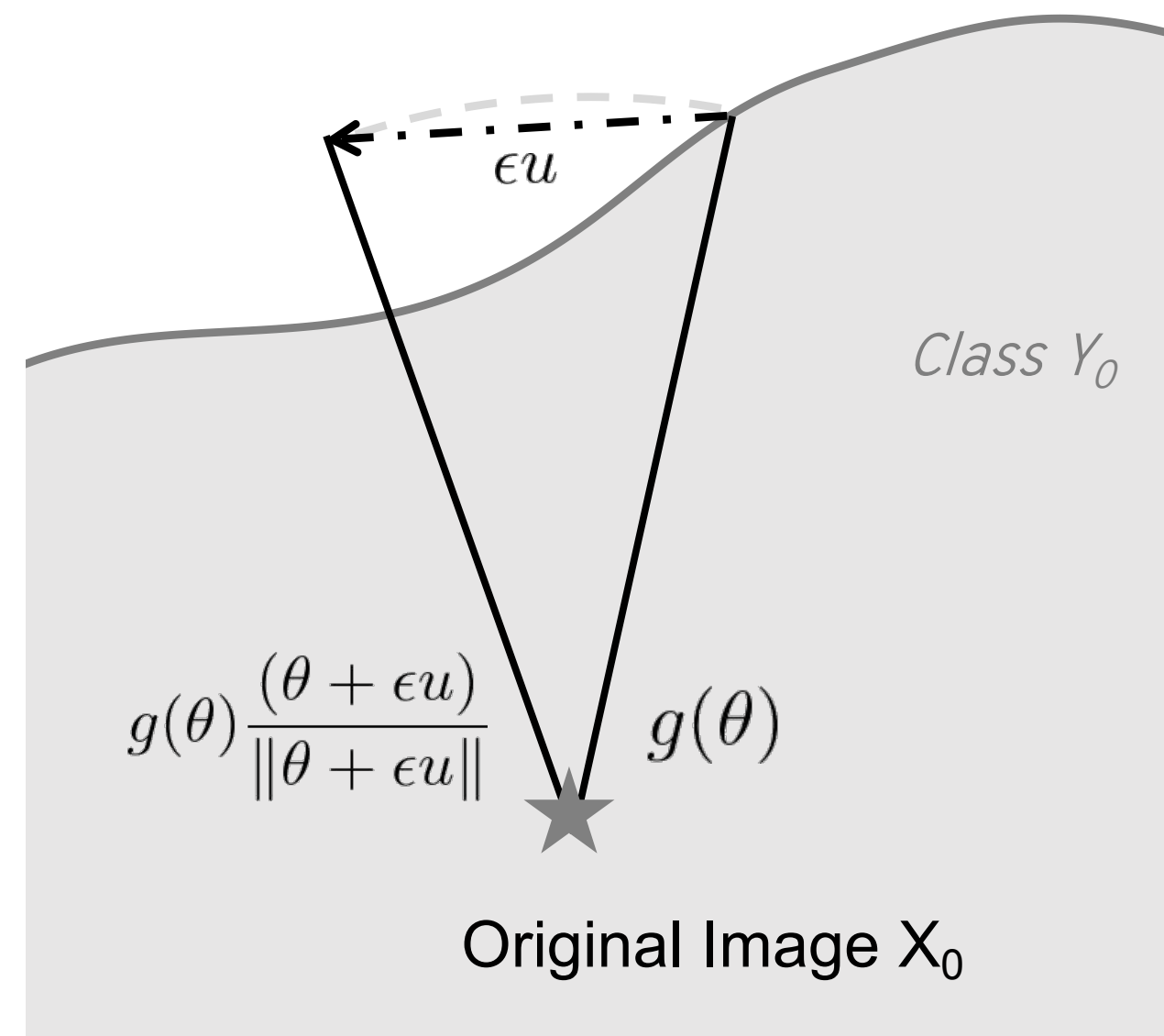
- $g(\theta_t)$ and $g(\theta_t + \beta u)$ in the gradient estimation takes most of queries, how to further reduce it?

Sign is enough!

- Binary search to estimate $g(\theta)$ in the gradient estimation takes most of queries.
- Gradient sign is powerful ! (FGSM)
- How to get the gradient sign efficiently ?

Single query oracle

- $$\text{sign}(g(\theta + \epsilon u) - g(\theta)) = \begin{cases} +1, & f(x_0 + g(\theta) \frac{(\theta + \epsilon u)}{\|\theta + \epsilon u\|}) = y_0, \\ -1, & \text{Otherwise.} \end{cases}$$



Sign-OPT attack

Algorithm 3 Sign-OPT attack (ICLR '20)

Input: Hard-label model f , original image x_0 , initial θ_0

for $t = 1, 2, \dots, T$ **do**

 Randomly sample u_1, \dots, u_Q from a Gaussian or Uniform distribution

 Evaluate $g(\theta_t)$

$$\hat{g} = \frac{g(\theta_t + \beta u) - g(\theta_t)}{\beta} \cdot u \Rightarrow \text{sign}\left(\frac{g(\theta_t + \beta u) - g(\theta_t)}{\beta}\right) \cdot u$$

 Update $\theta_{t+1} \leftarrow \theta_t - \eta \hat{g}$

 Evaluate $g(\theta_t)$ using the same search algorithm

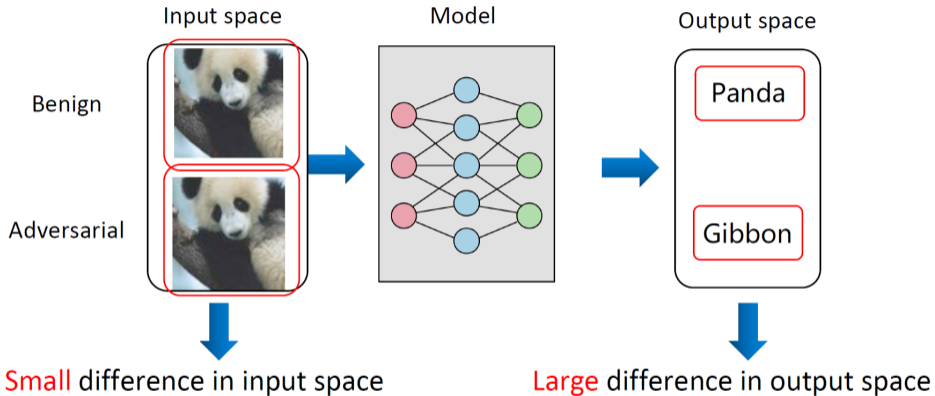
Transferability of Adversarial Examples

a.k.a Adversarial Transferability

Zeyu Qin

HKUST, AI Safety Group

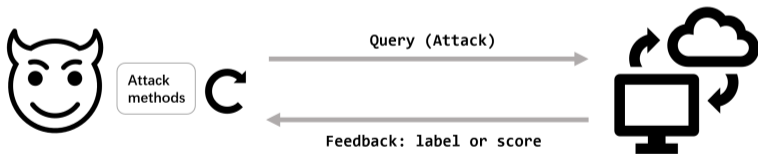
Adversarial Examples



- ▶ Imperceptible: $\|\mathbf{x}_{adv} - \mathbf{x}\|_p \leq \epsilon$
- ▶ Misclassified: $\mathbf{y} \neq Prediction(\mathbf{x}_{adv})$

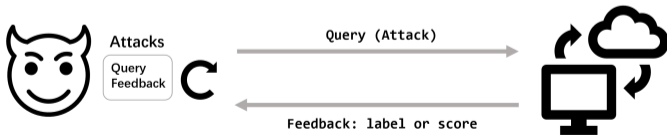
Severe Threats from Black-box Attacks

- ▶ In real scenarios, the users always get access to the DNN applications by querying API.
- ▶ The adversary does not know any knowledge about the target model \mathcal{M}_T , such as parameters, the architecture, and possibly the dataset.



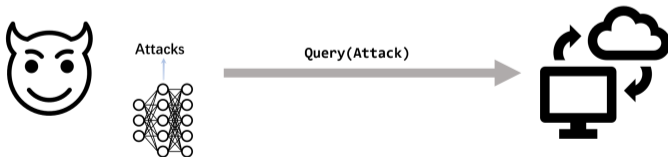
Black-box Attacks

- ▶ **Query-based attacks: only utilizing the feedback from \mathcal{M}_T**
 - Only know **the model feedback for each query (labels or confidence scores)** from \mathcal{M}_T .
 - By iteratively querying \mathcal{M}_T , the attackers generate x_{adv} based on exact feedback of each query.
 - It is more practical and only needs dozens of queries to generate the successful attacks.



Black-box Attacks

- ▶ **Transfer attacks: using x_{adv} from \mathcal{M}_S to attack \mathcal{M}_T**
 - The attackers can **utilize same dataset to train the surrogate model \mathcal{M}_S**
 - Generating x_{adv} (white-box attacks) on \mathcal{M}_S , then attacking \mathcal{M}_T .
 - Don't need to iteratively query but it is not practical and performs poor attack performance.



How to generate adversarial examples (white-box attacks)

- ▶ generating attacks

$$\arg \min_{\mathbf{x}^{adv}} f(\mathbf{x}^{adv}) := \mathcal{L}(\mathcal{M}(\mathbf{x}^{adv}), \mathbf{y}_t),$$
$$s.t. \|\mathbf{x}^{adv} - \mathbf{x}\|_p \leq \epsilon, \quad \mathbf{x}^{adv} \in [0, 1]$$

where p norm could be ℓ_∞ , ℓ_2 or ℓ_1 norm. \mathcal{L} could be CE loss, and $\mathcal{M}(\cdot)$ is Softmax output or logits (**returned feedback under black-box setting**).

Multi-step Projection Gradient Method (PGD) [1]:

$$\mathbf{x}^{adv} \leftarrow \text{Proj}_{\mathcal{B}(\mathbf{x}^{adv})} (\mathbf{x} - \alpha \cdot \text{sign}(\nabla \mathcal{L}(\mathcal{M}(\mathbf{x}^{adv}), \mathbf{y}_t)))$$
$$\mathcal{B} = \{\mathbf{x}^{adv} \mid [0, 1] \cap \|\mathbf{x}^{adv} - \mathbf{x}\|_\infty \leq \epsilon\}$$

[1]. Alexey Kurakin, Ian Goodfellow, Samy Bengio, "Adversarial examples in the physical world", Arxiv 2017

Transfer attacks

- ▶ Taking the targeted attack as an example, the general formulation of many existing transfer attack methods can be written as follows:

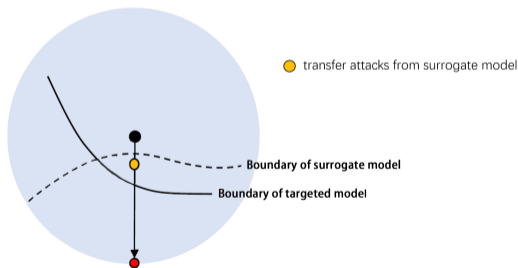
$$\min_{\mathbf{x}^{adv} \in \mathcal{B}_\epsilon(\mathbf{x})} \mathcal{L}(\mathcal{M}^S(\mathbf{x}^{adv}; \boldsymbol{\theta}), y_t). \quad (1)$$

where \mathcal{L} is the adversarial loss function computed on **surrogate model** \mathcal{M}^S , y_t is target label.

- ▶ Then, the attackers use \mathbf{x}^{adv} from \mathcal{M}^S to attack the target model, \mathcal{M}^T .

Transfer attacks overfits to \mathcal{M}_S

- ▶ The existing transfer attack methods exhibit poor transferability on \mathcal{M}_T (not successfully attacking \mathcal{M}_T)
- ▶ x_{adv} severely depends on (overfits to) the decision boundaries of \mathcal{M}_S and there are huge differences of decision boundaries between \mathcal{M}_S and \mathcal{M}_T . [1,2]



[1] Tramer et al., Ensemble Adversarial Training: Attacks and Defenses, ICLR 2018.

[2] Demontis et al., Why Do Adversarial Attacks Transfer? Explaining Transferability of Evasion and Poisoning Attacks, ACM CCS 2019.

How can we improve transfer attack performance?

- ▶ Remember our goal: to improve x_{adv} on unseen models (but trained on the same training dataset, like ImageNet). ('Sample' Generalization)
- ▶ It seems like the other important problem: to improve M on unseen data. (Model Generalization)
- ▶ Sample Generalization \Leftrightarrow Model Generalization ?

Inspiration from Improving Model Generalization

- ▶ You could borrow some **advanced optimization methods** to generate x_{adv} , like momentum [1,3], and variance reduction [4].
- ▶ Strong Augmentations: like crop and resizing [2] and Mixup [5].
- ▶ Unseen models: we could use much more diverse models during attack generation: model ensembling [6]
 - More diverse architectures
 - Bayesian model sampling (adding noise on parameters of M_s)

[1] Dong et al., Boosting Adversarial Attacks with Momentum, CVPR 2018.

[2] Xie et al., Improving transferability of adversarial examples with input diversity, CVPR 2019.

[3] Lin et al., Nesterov Accelerated Gradient and Scale Invariance for Adversarial Attacks, ICLR 2020.

[4] Wang et al., Enhancing the Transferability of Adversarial Attacks through Variance Tuning, CVPR 2021.

[5] Xiaosen Wang et al., Admix: Enhancing the Transferability of Adversarial Attacks, ICCV 2021.

[6] Tramer et al., Ensemble Adversarial Training: Attacks and Defenses, ICLR 2018.

Poor targeted attack performance

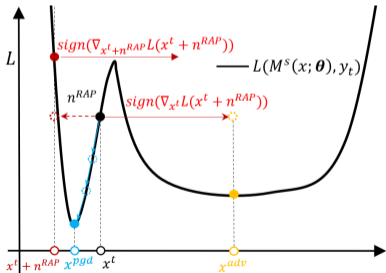
Figure from the recent NeurIPS 2021 paper:

Table 1: Targeted transfer success rates (%) in the single-model transfer setting. We consider three attacks with different loss functions: cross-entropy (CE), Poincaré distance with Triplet loss (Po+Trip) [21], and the logit loss. Results with 20/100/300 iterations are reported.

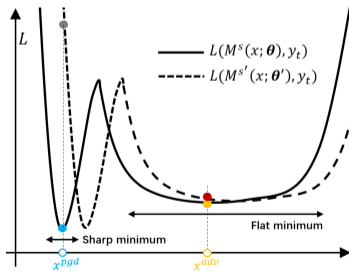
Attack	Source Model: Res50			Source Model: Dense121		
	→Dense121	→VGG16	→Inc-v3	→Res50	→VGG16	→Inc-v3
CE	26.9/39.4/42.6	17.3/27.3/30.4	2.4/3.8/4.1	13.1/17.3/19.4	7.7/10.8/10.9	1.9/3.3/3.5
Po+Trip	26.7/53.0/54.7	18.8/34.2/34.4	2.9/6.0/5.9	10.1/14.7/14.7	6.7/8.3/7.7	2.1/3.0/2.7
Logit	29.3/63.3/72.5	24.0/55.7/62.7	3.0/7.2/9.4	17.2/39.7/43.7	13.5/35.3/38.7	2.7/6.9/7.6
Attack	Source Model: VGG16			Source Model: Inc-v3		
	→Res50	→Dense121	→Inc-v3	→Res50	→Dense121	→VGG16
CE	0.7/0.4/0.6	0.5/0.3/0.1	0/0.1/0	0.6/ 2.1 /2.4	0.8/2.5/2.9	0.7/1.6/2.0
Po+Trip	0.6/0.8/0.5	0.6/0.6/0.7	0.2/0.1/0.1	0.6/2.0/2.5	0.8/ 3.1 /3.3	0.5/2.1/2.0
Logit	3.3/8.7/11.2	3.6/11.7/13.2	0.2/0.7/0.9	0.8/1.6/2.9	1.2/2.8/5.3	0.7/2.2/3.7

Further Improving Transferability

- ▶ Qin et al., propose a new perspective to interpret the adversarial transferability, the flatness of (adversarial) loss landscape of x^{adv} on \mathcal{M}^S .
- ▶ The x^{adv} located at the flat local minimum is less sensitive to the changes of decision boundary (the difference of \mathcal{M}^S and \mathcal{M}^T). Therefore, it could have better adversarial transferability.



(a)



(b)

Finding x_{adv} located at a local flat region

- ▶ Qin et al., encourage that **not only** x_{adv} itself has low loss value, but also the points **in the vicinity of** x_{adv} have similarly low loss values.
- ▶ Qin et al., proposes to **minimize the maximal loss value within a local neighborhood region** around x_{adv} .
- ▶ The maximal loss is implemented by perturbing x_{adv} (adding perturbation \mathbf{n}) to maximize the adversarial loss, named **Reverse Adversarial Perturbation (RAP)**. So, we aim to solve this problem,

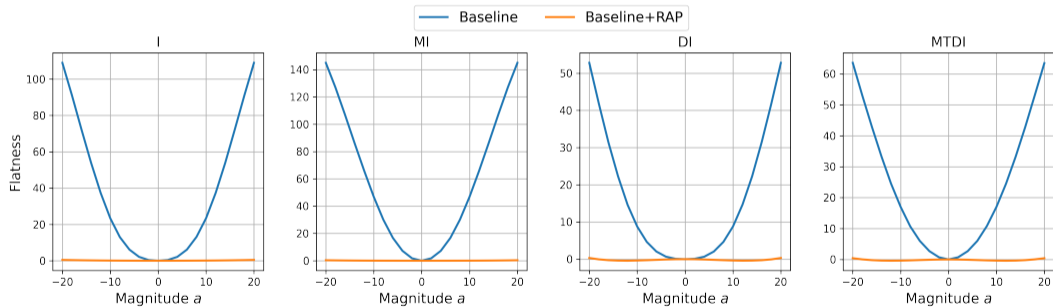
$$\min_{\mathbf{x}_{adv} \in \mathcal{B}_\epsilon(\mathbf{x})} \mathcal{L}(\mathcal{M}_S(\mathbf{x}_{adv} + \mathbf{n}_{adv}; \boldsymbol{\theta}), y_t)$$

Where,

$$\mathbf{n}_{adv} = \arg \max_{\|\mathbf{n}\|_\infty \leq \epsilon_n} \mathcal{L}(\mathcal{M}_S(\mathbf{x}_{adv} + \mathbf{n}; \boldsymbol{\theta}), y_t)$$

A Closer Look at RAP

- ▶ First we **visualize the loss landscape around x^{adv} on \mathcal{M}^S** by plotting the loss variations. We can observe that RAP could help find x^{adv} located at the flat region.



RAP achieves the better attack performance

- Combined with existing attacks, RAP further boosts their transferability for both **untargeted and targeted attacks**.

The below tables show the transfer targeted attack performance ($\mathcal{M}^S \Rightarrow \mathcal{M}^T$).

Attack	ResNet-50 \Rightarrow			DenseNet-121 \Rightarrow		
	Dense-121	VGG-16	Inc-v3	Res-50	VGG-16	Inc-v3
MTDI / +RAP / +RAP-LS	74.9 / <u>78.2</u> / 88.5	62.8 / <u>72.9</u> / 81.5	10.9 / <u>28.3</u> / 33.2	44.9 / <u>64.3</u> / 74.5	38.5 / <u>55.0</u> / 65.5	7.7 / <u>23.0</u> / 26.5
MTDSI / +RAP / +RAP-LS	86.3 / <u>88.4</u> / 93.3	70.1 / <u>77.7</u> / 84.7	38.1 / <u>51.8</u> / 58.0	55.0 / <u>71.2</u> / 75.8	42.0 / <u>58.4</u> / 62.3	19.8 / <u>39.0</u> / 39.2
MTDAI / +RAP / +RAP-LS	<u>91.4</u> / 89.4 / 93.6	<u>79.9</u> / 79.0 / 86.3	50.8 / <u>57.1</u> / 64.1	69.1 / <u>74.2</u> / 82.1	54.7 / <u>63.1</u> / 69.3	32.0 / <u>43.5</u> / 49.3

Attack	VGG-16 \Rightarrow			Inc-v3 \Rightarrow		
	Res-50	Dense-121	Inc-v3	Res-50	Dense-121	VGG-16
MTDI / +RAP / +RAP-LS	11.8 / <u>16.7</u> / 22.9	13.7 / <u>19.4</u> / 27.4	0.7 / <u>3.4</u> / 4.6	1.8 / 8.3 / <u>7.5</u>	4.1 / 14.8 / <u>13.4</u>	2.9 / <u>8.0</u> / 9.8
MTDSI / +RAP / +RAP-LS	31.0 / <u>35.3</u> / 38.7	41.7 / <u>44.4</u> / 49.6	9.6 / 15.2 / <u>13.7</u>	5.6 / 11.9 / <u>10.7</u>	10.4 / 21.2 / <u>20.9</u>	4.2 / 8.9 / <u>8.6</u>
MTDAI / +RAP / +RAP-LS	36.2 / <u>39.0</u> / 43.1	<u>48.0</u> / 45.1 / 55.2	11.6 / <u>17.1</u> / 17.6	9.6 / <u>13.6</u> / 16.7	17.9 / <u>27.5</u> / 31.6	8.4 / <u>12.0</u> / 12.1

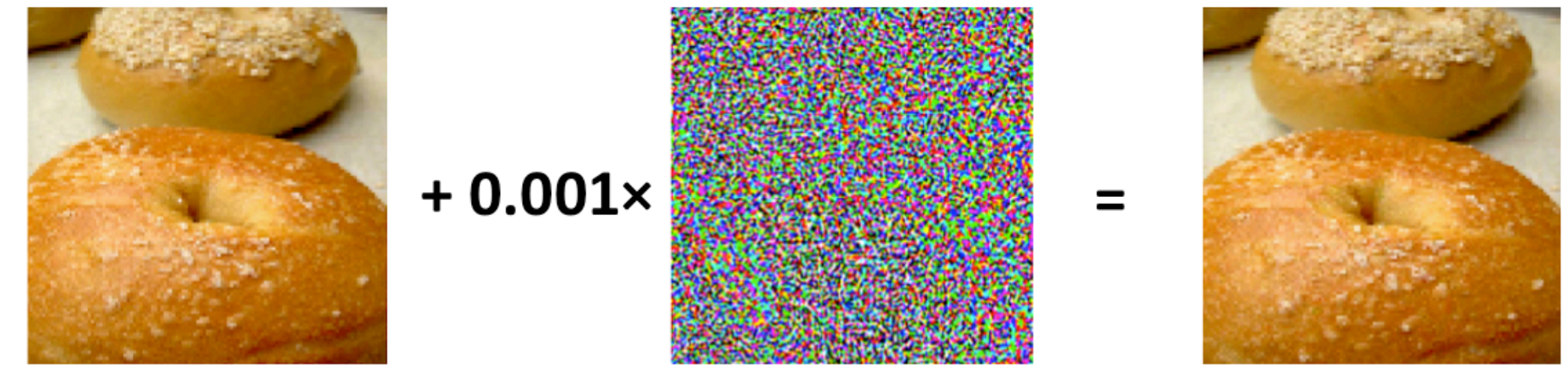
Adversarial Defense Methods

Adversarial Training

Test-time integrity

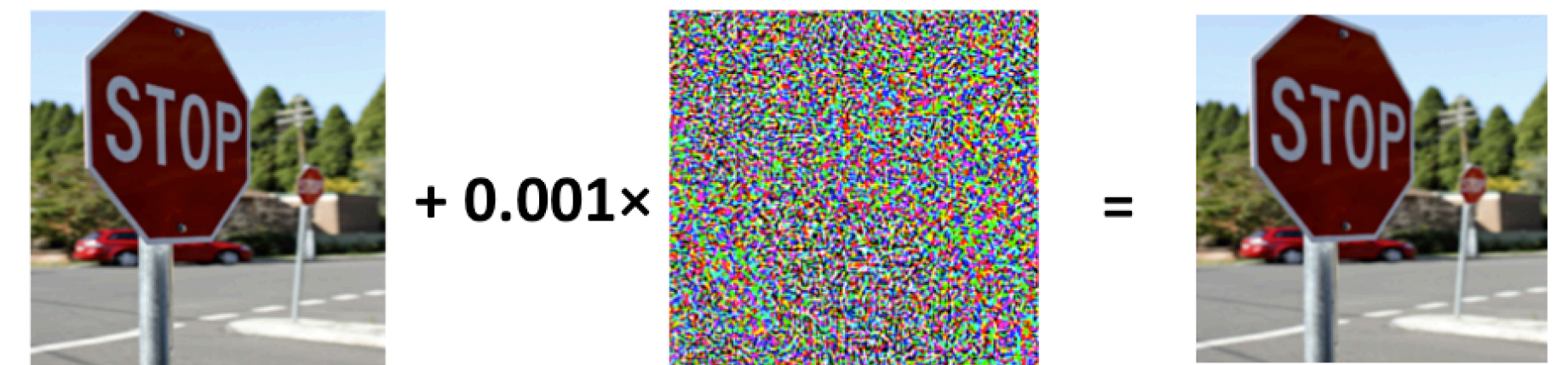
Adversarial examples

- An **adversarial** example can easily fool a deep network
- **Robustness** is critical in real systems



Bagel

piano



stop sign

speed limit 40

Adversarial example

White-box adversarial attack

- If there is $\|x - x_0\|_\infty$ constraint, we could turn to solve by
- FGSM attack [GSS15]:
 - $x \leftarrow \text{proj}_{x+\mathcal{S}}(x_0 + \alpha \text{sign}(\nabla_{x_0} \ell(\theta, x, y)))$
- PGD attack [KGB17, MMS18]
 - $x^{t+1} \leftarrow \text{proj}_{x+\mathcal{S}}(x^t + \alpha \text{sign}(\nabla_{x^t} \ell(\theta, x, y)))$

Adversarial defense

Adversarial training

- Adversarial training [MMS18]:

- $$\min_{\theta} \rho(\theta), \quad \text{where} \quad \rho(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \mathcal{S}} L(\theta, x + \delta, y) \right].$$

- Solve the inner loop by

- $$x^{t+1} = \Pi_{x+\mathcal{S}} (x^t + \alpha \text{sgn}(\nabla_x L(\theta, x, y)))$$

Adversarial training

Capacity is crucial

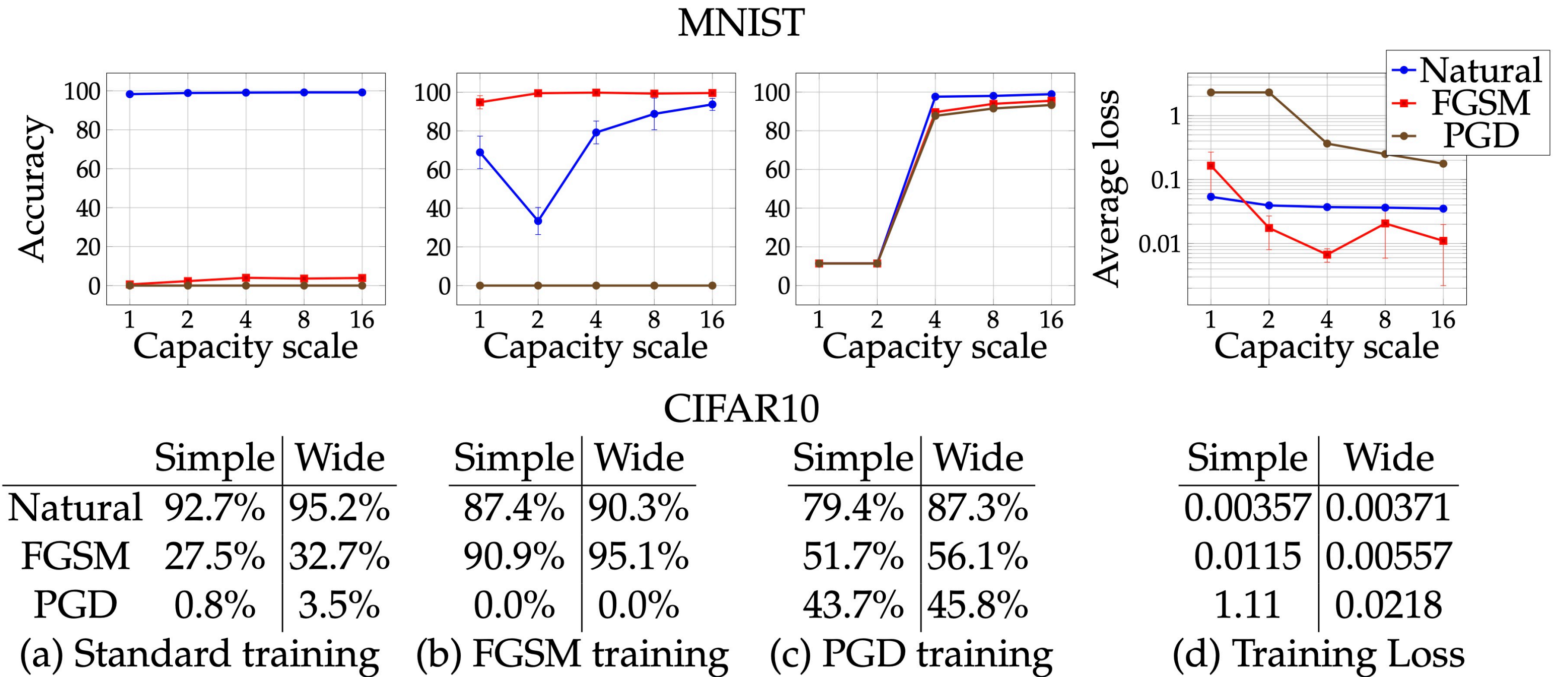


Figure 4: The effect of network capacity on the performance of the network. We trained MNIST and CIFAR10 networks of varying capacity on: (a) natural examples, (b) with FGSM-made adversarial examples, (c) with PGD-made adversarial examples. In the first three plots/tables of each dataset, we show how the standard and adversarial accuracy changes with respect to capacity for each training regime. In the final plot/table, we show the value of the cross-entropy loss on the adversarial examples the networks were trained on. This corresponds to the value of our saddle point formulation (2.1) for different sets of allowed perturbations.

Adversarial training

Problems

- Huge overhead
 - Increase training time by an order magnitude (7x if 7 step PGD)
- Fast method like FGSM doesn't work
 - Easily be attacked by strong attackers such as C&W attack

Fast Adversarial training

Algorithm 3 FGSM adversarial training for T epochs, given some radius ϵ , N PGD steps, step size α , and a dataset of size M for a network f_θ

```
for  $t = 1 \dots T$  do  
  for  $i = 1 \dots M$  do  
    // Perform FGSM adversarial attack  
     $\delta = \text{Uniform}(-\epsilon, \epsilon)$   
     $\delta = \delta + \alpha \cdot \text{sign}(\nabla_\delta \ell(f_\theta(x_i + \delta), y_i))$   
     $\delta = \max(\min(\delta, \epsilon), -\epsilon)$   
     $\theta = \theta - \nabla_\theta \ell(f_\theta(x_i + \delta), y_i)$  // Update model weights with some optimizer, e.g. SGD  
  end for  
end for
```

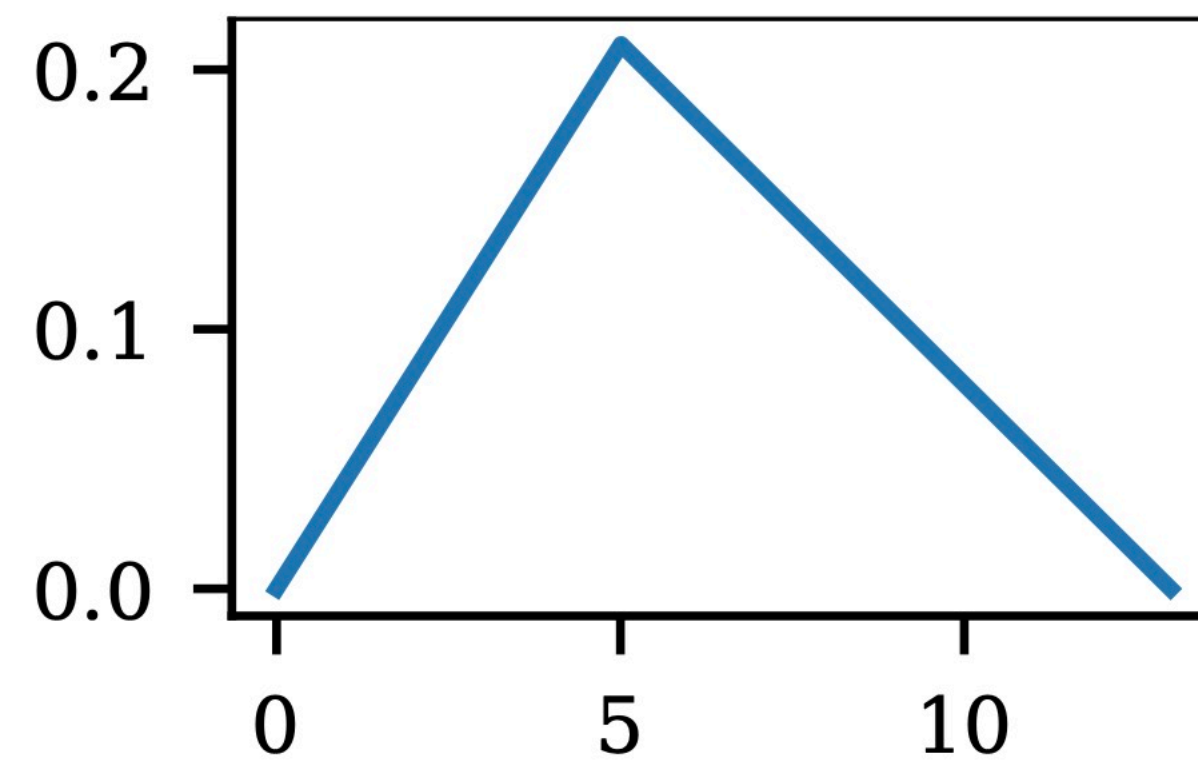
The magic of random initialization

Method	Standard accuracy	PGD ($\epsilon = 8/255$)	Time (min)
FGSM + DAWNBench			
+ zero init	85.18%	0.00%	12.37
+ early stopping	71.14%	38.86%	7.89
+ previous init	86.02%	42.37%	12.21
+ random init	85.32%	44.01%	12.33
+ $\alpha = 10/255$ step size	83.81%	46.06%	12.17
+ $\alpha = 16/255$ step size	86.05%	0.00%	12.06
+ early stopping	70.93%	40.38%	8.81
“Free” ($m = 8$) (Shafahi et al., 2019) ¹	85.96%	46.33%	785
+ DAWNBench	78.38%	46.18%	20.91
PGD-7 (Madry et al., 2017) ²	87.30%	45.80%	4965.71
+ DAWNBench	82.46%	50.69%	68.8

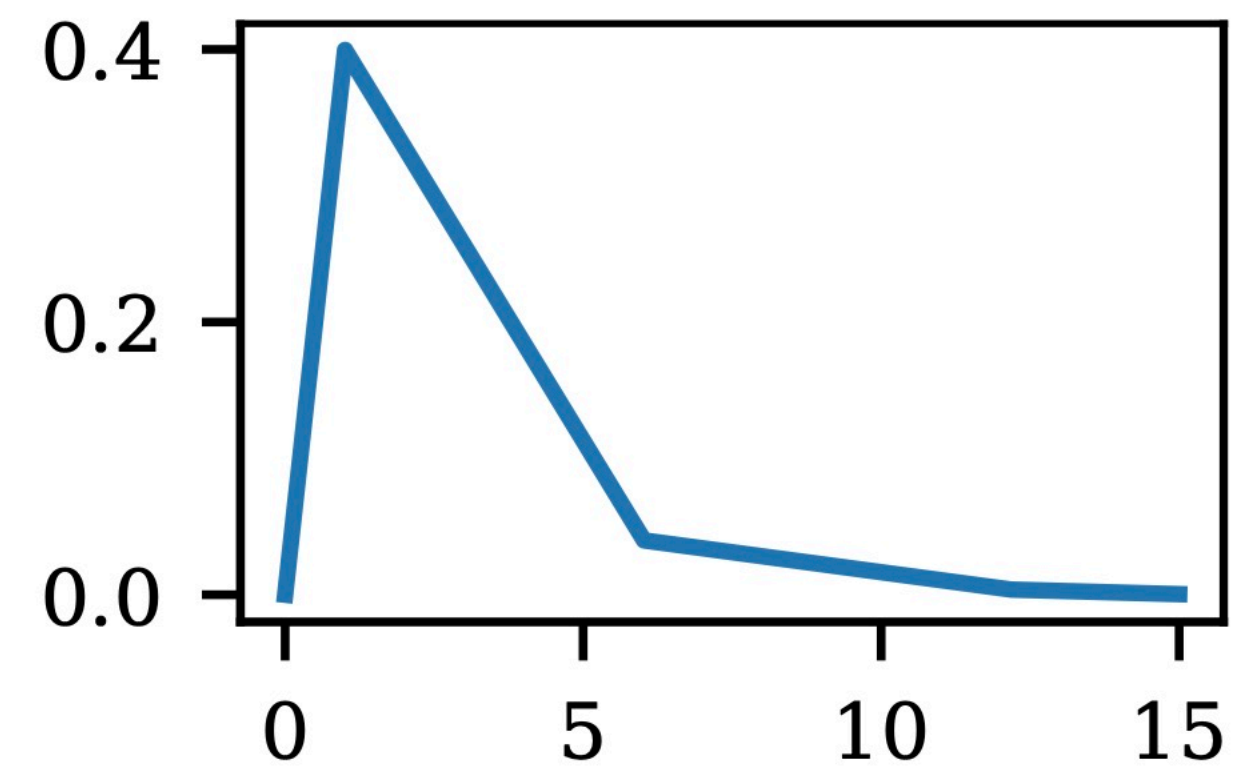
DAWNBench Improvement

Reduce # of training epochs

- Cyclic learning rate
- Mixed-precision arithmetic



(a) CIFAR10



(b) ImageNet

Catastrophic overfitting

