

Adversarial Examples in Machine Learning

CS489/698

Rey Reza Wiyatno

Email: reywiyatno@gmail.com

Agenda

Intro

Adversarial Attacks

Adversarial Defenses

Takeaways

Agenda

Intro

Adversarial Attacks

Adversarial Defenses

Takeaways

Adversarial Examples

Inputs specifically designed to cause a model to make mistakes in its prediction, although they look like valid inputs to a human.

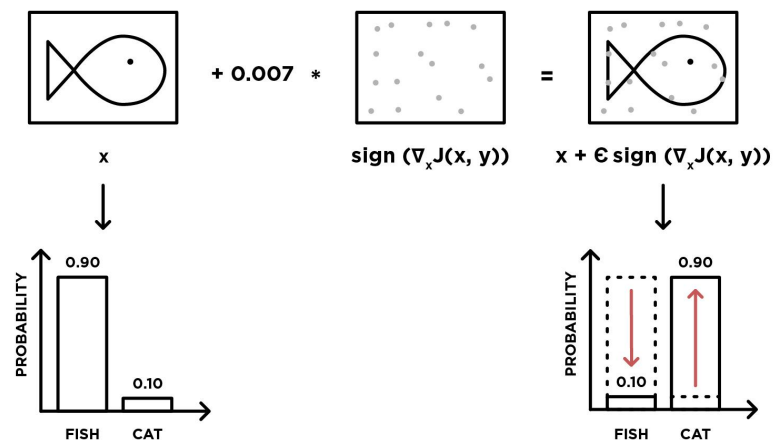


Illustration of adversarial example.

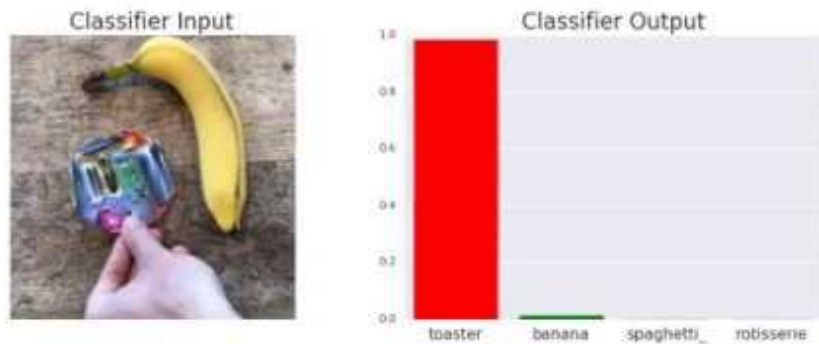
Standard Definition in Classification Task

Generate adversarial example \mathbf{x}' that looks like an original example \mathbf{x} , but is mispredicted by a model. Formally:

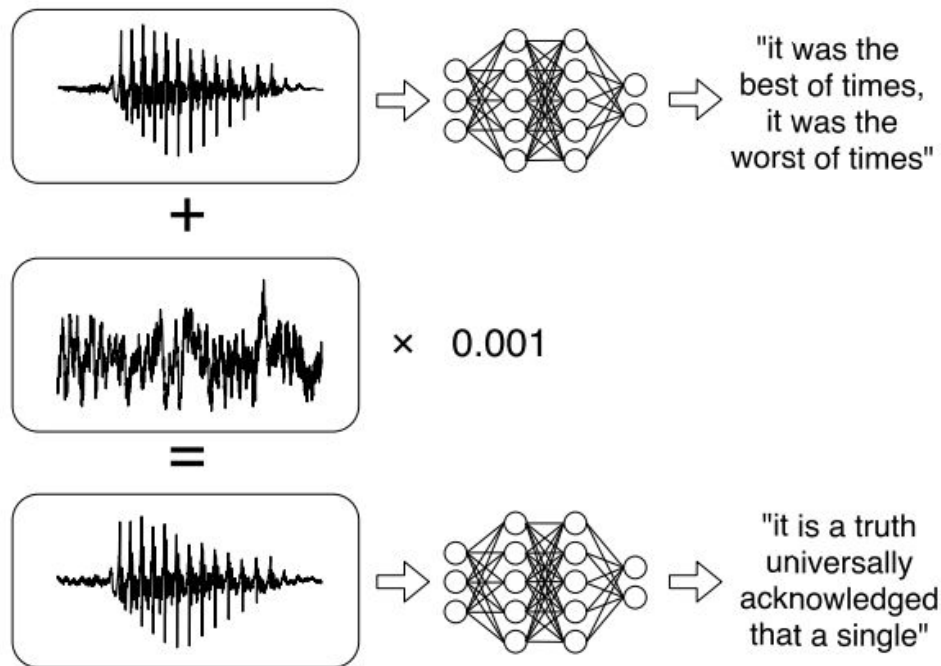
Given a model f and an input \mathbf{x} , find \mathbf{x}' where $\|\mathbf{x} - \mathbf{x}'\|_p < \epsilon$, such that $f(\mathbf{x}) \neq f(\mathbf{x}')$

Note: other perceptual similarity metrics other than the L_p norm can also be used.

Adversarial Examples: Adversarial Patch



Adversarial Examples: Adversarial Audio



Demo: https://nicholas.carlini.com/code/audio_adversarial_examples/

Audio Adversarial Example: Targeted Attacks on Speech-to-Text (Carlini & Wagner, 2018)

Common Terms

Adversarial attacks: methods to generate adversarial examples.

Adversarial defenses: methods to defend against adversarial examples.

Adversarial robustness: property to resist misclassification of adversarial examples.

Adversarial detection: methods to detect adversarial examples.

Transferability: adversarial examples generated to fool a specific model can also be used to fool other models.

Common Terms

Adversarial perturbation: difference between original example and its adversarial counterpart

Whitebox attack: when attacker has full access to the victim model

Blackbox attack: when attacker only has access to victim's output

Targeted attack: when attacker wants an adversary to be mispredicted in a specific way

Non-targeted attack: when attacker does not care if an example is mispredicted in a specific way

Agenda

Intro

Adversarial Attacks

Adversarial Defenses

Takeaways

Goal:

Minimally modify inputs to
confuse the victim model

We can re-formulate as:

Minimally modify inputs to
maximize some loss function

Any ideas how?

Attack & Defense Methods

Attack strategies:

Whitebox:

- Direct gradient step(s): FGSM, BIM, R+FGSM, DAG
- Gradient-based greedy algorithm: JSMA
- Iterative optimization: L-BFGS, C&W, UAP, Adversarial Eyeglasses, RP₂, EOT
- Parameterized optimization: ATN

Blackbox:

- Decision boundary approximation: SBA
- Evolutionary algorithm: One Pixel Attack
- Finite difference method: ZOO

Defense strategies:

Data augmentation:

- Adversarial training
- Ensemble adversarial training
- PGD adversarial training

Gradient regularization:

- DCN
- Defensive distillation

Input manipulation:

- PixelDefend

Detection methods:

- Adversary detector networks
- Feature squeezing

Fast Gradient Sign Method (FGSM)

- Take gradient w.r.t. input
- Then do **single-step gradient ascent**

$$x' = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x, y))$$

- Is above equation for targeted/non-targeted attack?
- If targeted, how to make it untargeted (and vice-versa)?
- Why use “sign” operator?

Note: $\mathcal{L}(x, y)$ and $J(x, y)$ denote the training loss function, x is an input, y is the true label of x , ϵ is a small constant where $\epsilon > 0$

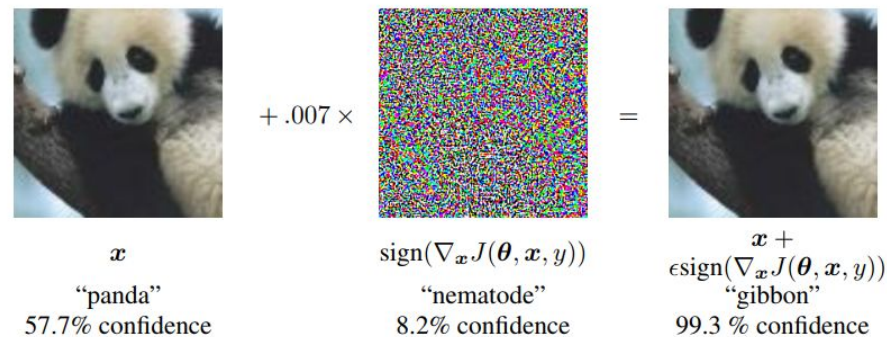


Illustration of FGSM. Adversarial example (right) misclassified as “Gibbon” with high confidence (Goodfellow et al., 2015).

Implementing FGSM (see notebook)

Link: https://github.com/rrwiyatn/deep-learning/blob/master/fast_gradient_sign_attack/fgsm.ipynb

Basic Iterative Method (BIM)

- Iterative variant of FGSM

$$\mathbf{X}_0^{adv} = \mathbf{X}, \quad \mathbf{X}_{N+1}^{adv} = \text{Clip}_{X,\epsilon} \left\{ \mathbf{X}_N^{adv} + \alpha \text{sign}(\nabla_X J(\mathbf{X}_N^{adv}, y_{true})) \right\}$$

- Targeted attack variant of BIM called Iterative Least-Likely Class Method (ILLCM)

$$\mathbf{X}_0^{adv} = \mathbf{X}, \quad \mathbf{X}_{N+1}^{adv} = \text{Clip}_{X,\epsilon} \left\{ \mathbf{X}_N^{adv} - \alpha \text{sign}(\nabla_X J(\mathbf{X}_N^{adv}, y_{LL})) \right\}$$

$$y_{LL} = \arg \min_y \{p(y|\mathbf{X})\}$$

$$\text{Clip}_{X,\epsilon} \{ \mathbf{X}' \} (x, y, z) = \min \left\{ 255, \mathbf{X}(x, y, z) + \epsilon, \max \{ 0, \mathbf{X}(x, y, z) - \epsilon, \mathbf{X}'(x, y, z) \} \right\}$$

- When to use ILLCM?

Adversarial Examples in the Physical World (Kurakin et al., 2017)

Adversarial Machine Learning at Scale (Kurakin et al., 2017)

Random FGSM (R+FGSM)

- FGSM variant with a random starting point

$$x^{\text{adv}} = x' + (\varepsilon - \alpha) \cdot \text{sign}(\nabla_{x'} J(x', y_{\text{true}})), \quad \text{where } x' = x + \alpha \cdot \text{sign}(\mathcal{N}(\mathbf{0}^d, \mathbf{I}^d))$$

- Designed to circumvent a defense method called adversarial training (Goodfellow et al., 2015) in its naive implementation
- We will come back to the motivation later in defense section

Attack & Defense Methods

Attack strategies:

Whitebox:

- Direct gradient step(s): FGSM, BIM, R+FGSM, DAG
- Gradient-based greedy algorithm: JSMA
- Iterative optimization: L-BFGS, UAP, C&W, Adversarial Eyeglasses, RP_2 , EOT
- Parameterized optimization: ATN

Blackbox:

- Decision boundary approximation: SBA
- Evolutionary algorithm: One Pixel Attack
- Finite difference method: ZOO

Defense strategies:

Data augmentation:

- Adversarial training
- Ensemble adversarial training
- PGD adversarial training

Gradient regularization:

- DCN
- Defensive distillation

Input manipulation:

- PixelDefend

Detection methods:

- Adversary detector networks
- Feature squeezing

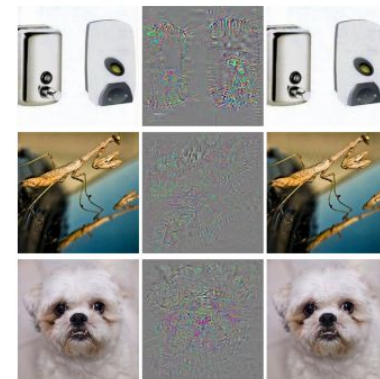
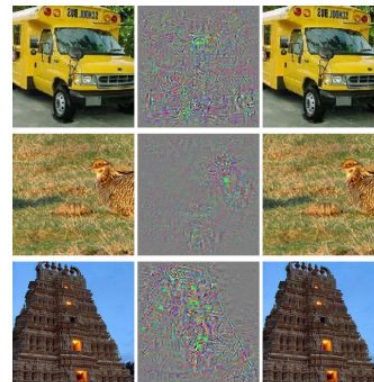
L-BFGS Attack

- Model adversarial example generation as optimization problem
- Use L-BFGS as optimizer
- Given a victim model f , find r that minimizes:

$$c|r| + \text{loss}_f(x + r, l)$$

subject to $x + r \in [0, 1]^m$

- Note that the loss function does not have to be the same with the training loss of the victim



Adversarial examples generated using this method (Szegedy et al., 2014).

Universal Adversarial Perturbation (UAP)

Universal perturbations:

Single perturbation that can be added to multiple inputs to make them adversarial

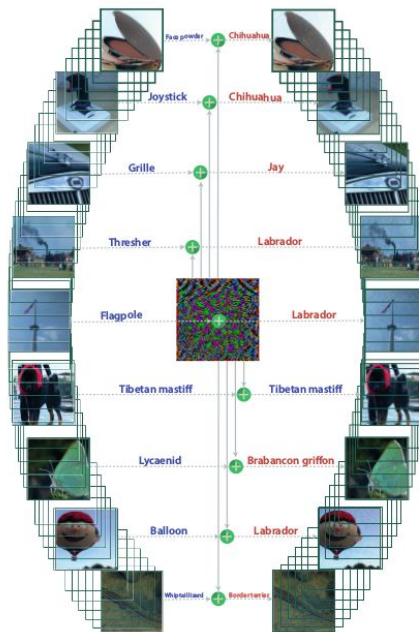


Illustration of UAP
(Moosavi-Dezfooli et al., 2016).

Algorithm 1 Computation of universal perturbations.

- 1: **input:** Data points X , classifier \hat{k} , desired ℓ_p norm of the perturbation ξ , desired accuracy on perturbed samples δ .
- 2: **output:** Universal perturbation vector v .
- 3: Initialize $v \leftarrow 0$.
- 4: **while** $\text{Err}(X_v) \leq 1 - \delta$ **do**
- 5: **for** each datapoint $x_i \in X$ **do**
- 6: **if** $\hat{k}(x_i + v) = \hat{k}(x_i)$ **then**
- 7: Compute the *minimal* perturbation that sends $x_i + v$ to the decision boundary:
- $\Delta v_i \leftarrow \arg \min_r \|r\|_2 \text{ s.t. } \hat{k}(x_i + v + r) \neq \hat{k}(x_i).$
- 8: Update the perturbation:
- $v \leftarrow \mathcal{P}_{p,\xi}(v + \Delta v_i).$
- 9: **end if**
- 10: **end for**
- 11: **end while**

$$\mathcal{P}_{p,\xi}(v) = \arg \min_{v'} \|v - v'\|_2 \text{ subject to } \|v'\|_p \leq \xi$$

Attack & Defense Methods

Attack strategies:

Whitebox:

- Direct gradient step(s): FGSM, BIM, R+FGSM, DAG
- Gradient-based greedy algorithm: JSMA
- Iterative optimization: L-BFGS, C&W, UAP, Adversarial Eyeglasses, RP_2 , EOT
- Parameterized optimization: ATN

Blackbox:

- Decision boundary approximation: SBA
- Evolutionary algorithm: One Pixel Attack
- Finite difference method: ZOO

Defense strategies:

Data augmentation:

- Adversarial training
- Ensemble adversarial training
- PGD adversarial training

Gradient regularization:

- DCN
- Defensive distillation

Input manipulation:

- PixelDefend

Detection methods:

- Adversary detector networks
- Feature squeezing

Adversarial Transformation Network (ATN)

- Train a neural network to generate adversaries
- 2 variants: Adversarial Autoencoder (AAE) and Perturbation ATN (P-ATN)
- AAE: Given a victim model f , train a generator network G_t on dataset X to output adversarial examples X' such that $f(X') = t$, where t is a target misclassification class
- Every G_t can only be used generate adversarial examples that are misclassified as class t

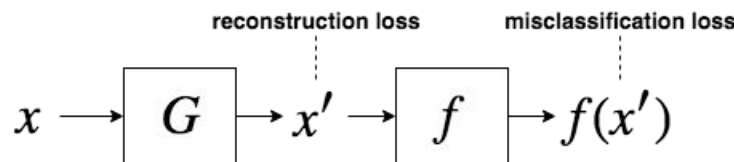
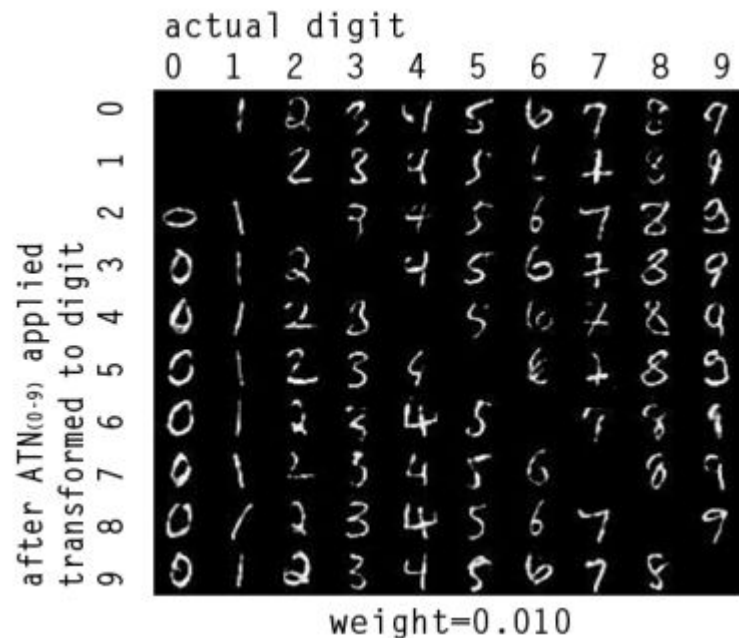


Illustration of AAE.

Adversarial Transformation Network (ATN)



Generated adversaries that fool Inception Resnet V2 (Baluja & Fischer, 2017)



Generated adversarial MNIST (Baluja & Fischer, 2017)

Attack & Defense Methods

Attack strategies:

Whitebox:

- Direct gradient step(s): FGSM, BIM, R+FGSM, DAG
- Gradient-based greedy algorithm: JSMA
- Iterative optimization: L-BFGS, C&W, UAP, Adversarial Eyeglasses, RP₂, EOT
- Parameterized optimization: ATN

Blackbox:

- Decision boundary approximation: SBA
- Evolutionary algorithm: One Pixel Attack
- Finite difference method: ZOO

Defense strategies:

Data augmentation:

- Adversarial training
- Ensemble adversarial training
- PGD adversarial training

Gradient regularization:

- DCN
- Defensive distillation

Input manipulation:

- PixelDefend

Detection methods:

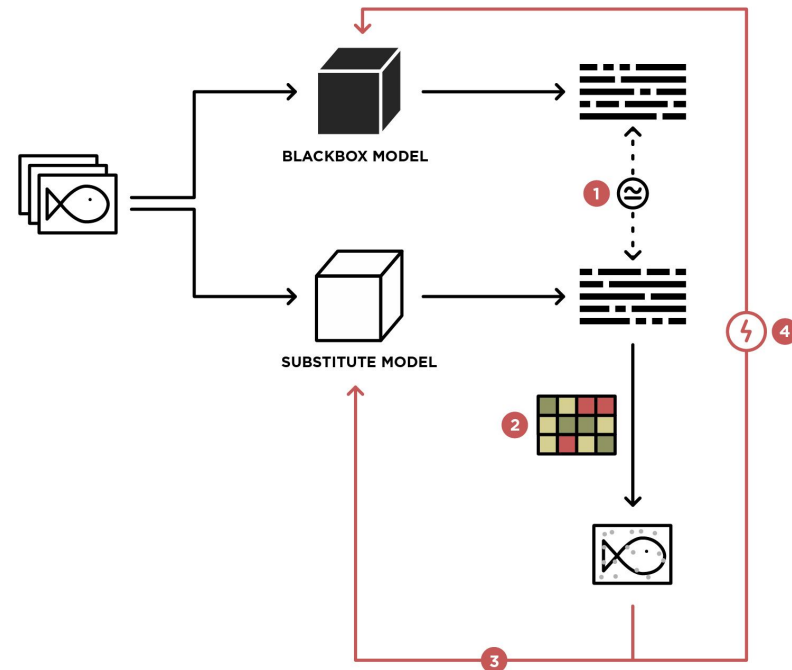
- Adversary detector networks
- Feature squeezing

Substitute Blackbox Attack (SBA)

- Blackbox attack by approximating decision boundaries of victim
- Assumption: attacker has access to the softmax probability
- Procedures (see next slide):
 - Train a substitute model on a dataset **labelled by the victim**
 - Attack the substitute model using any whitebox methods
 - The generated adversaries should be transferable to the victim model (thanks to transferability property)
- Successfully attacked Google, Amazon, and MetaMind image recognition models

Substitute Blackbox Attack (SBA)

1. Train a substitute model on a dataset **labelled by the victim (i.e., blackbox model)**
2. Attack the substitute model using any whitebox methods
3. Validate that the adversarial examples fool the substitute model
4. Use the adversarial examples to fool the blackbox model



Attack & Defense Methods

Attack strategies:

Whitebox:

- Direct gradient step(s): FGSM, BIM, R+FGSM, DAG
- Gradient-based greedy algorithm: JSMA
- Iterative optimization: L-BFGS, C&W, UAP, Adversarial Eyeglasses, RP_2 , EOT
- Parameterized optimization: ATN

Blackbox:

- Decision boundary approximation: SBA
- Evolutionary algorithm: One Pixel Attack
- Finite difference method: ZOO

Defense strategies:

Data augmentation:

- Adversarial training
- Ensemble adversarial training
- PGD adversarial training

Gradient regularization:

- DCN
- Defensive distillation

Input manipulation:

- PixelDefend

Detection methods:

- Adversary detector networks
- Feature squeezing

Zeroth Order Optimization (ZOO)

- Blackbox attack via **finite difference approximation**

$$\hat{g}_i := \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_i} \approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x} - h\mathbf{e}_i)}{2h}$$

- f is victim model, \mathbf{e}_i is a vector where only the i -th element is 1, \mathbf{x}_i is the i -th element of input \mathbf{x}
- What is the main disadvantage of using finite difference?

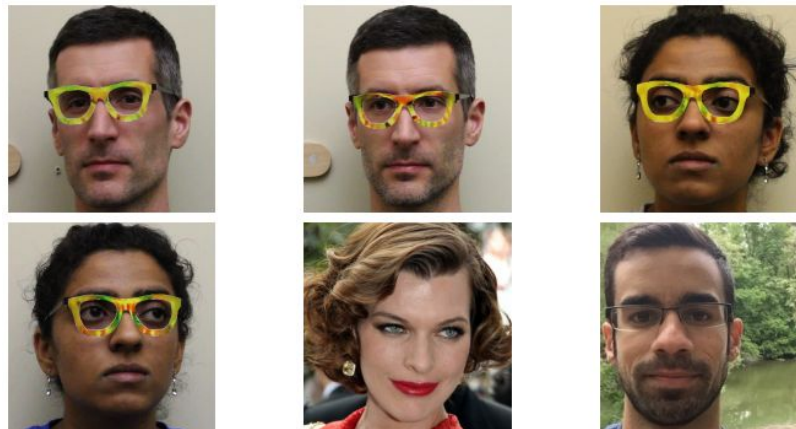
Algorithm 2 ZOO-ADAM: Zeroth Order Stochastic Coordinate Descent with Coordinate-wise ADAM

Require: Step size η , ADAM states $M \in \mathbb{R}^p, v \in \mathbb{R}^p, T \in \mathbb{Z}^p$, ADAM hyper-parameters $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$

```
1:  $M \leftarrow \mathbf{0}, v \leftarrow \mathbf{0}, T \leftarrow \mathbf{0}$ 
2: while not converged do
3:   Randomly pick a coordinate  $i \in \{1, \dots, p\}$ 
4:   Estimate  $\hat{g}_i$  using (6)
5:    $T_i \leftarrow T_i + 1$ 
6:    $M_i \leftarrow \beta_1 M_i + (1 - \beta_1)\hat{g}_i, \quad v_i \leftarrow \beta_2 v_i + (1 - \beta_2)\hat{g}_i^2$ 
7:    $\hat{M}_i = M_i / (1 - \beta_1^{T_i}), \quad \hat{v}_i = v_i / (1 - \beta_2^{T_i})$ 
8:    $\delta^* = -\eta \frac{\hat{M}_i}{\sqrt{\hat{v}_i + \epsilon}}$ 
9:   Update  $\mathbf{x}_i \leftarrow \mathbf{x}_i + \delta^*$ 
10: end while
```

Some real world examples

Adversarial Eyeglasses



Example of adversarial eyeglasses (Sharif et al., 2016). Top row depicts an input given to a model, while bottom row depicts the person from the target misclassification class.

Adversarial Road Signs

- Similar attack method to the adversarial eyeglasses with different masks
- Perceptually different, but **inconspicuous**



“STOP” signs misclassified as a “speed limit” sign
(Eykholt et al., 2017).

Adversarial Turtle



3D printed adversarial turtle misclassified
as rifle (Athalye et al., 2018).

Agenda

Intro

Adversarial Attacks

Adversarial Defenses

Takeaways

Attack & Defense Methods

Attack strategies:

Whitebox:

- Direct gradient step(s): FGSM, BIM, R+FGSM, DAG
- Gradient-based greedy algorithm: JSMA
- Iterative optimization: L-BFGS, C&W, UAP, Adversarial Eyeglasses, RP₂, EOT
- Parameterized optimization: ATN

Blackbox:

- Decision boundary approximation: SBA
- Evolutionary algorithm: One Pixel Attack
- Finite difference method: ZOO

Defense strategies:

Data augmentation:

- Adversarial training
- Ensemble adversarial training
- PGD adversarial training

Gradient regularization:

- DCN
- Defensive distillation

Input manipulation:

- PixelDefend

Detection methods:

- Adversary detector networks
- Feature squeezing

Adversarial Training

- Include adversarial examples as part of the training set
- If using FGSM adversaries, training loss becomes:

$$\tilde{J}(\boldsymbol{\theta}, \mathbf{x}, y) = \alpha J(\boldsymbol{\theta}, \mathbf{x}, y) + (1 - \alpha) J(\boldsymbol{\theta}, \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y)))$$

- In other words, new adversarial examples are generated **per training iteration** based on the state of the model at that iteration
- **NOT** the same as the Generative Adversarial Nets (GAN)
- Robust to adversaries included during adversarial training

Note: $J(\cdot)$ denotes the classification loss (e.g. cross-entropy), $\boldsymbol{\theta}$ denotes the model's parameter, α denotes a constant that weigh the importance on classifying normal versus adversarial examples where $\alpha \in [0, 1]$.

Adversarial Training

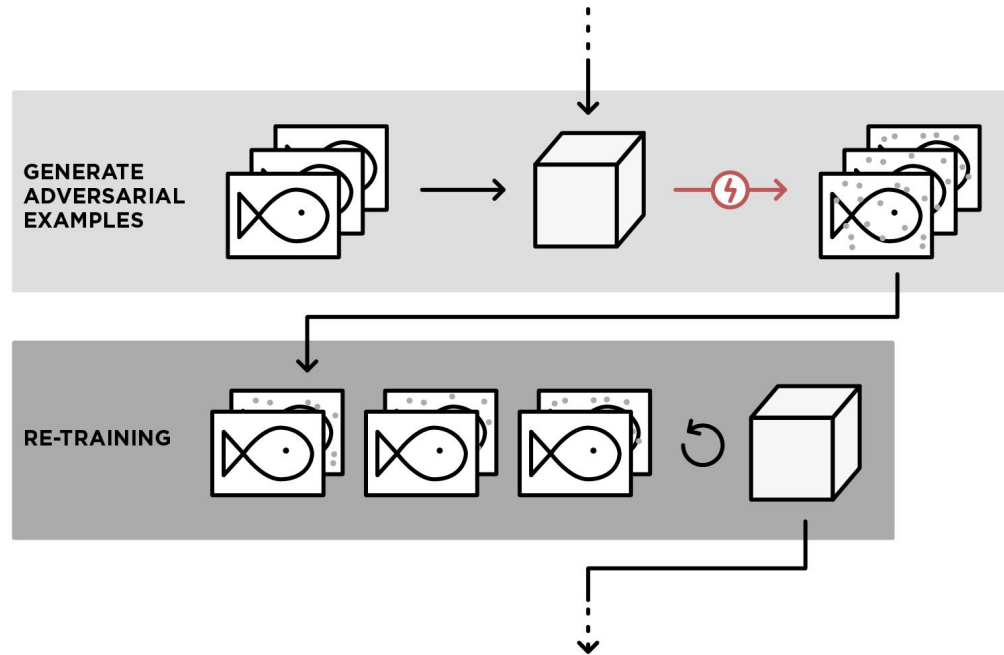


Illustration of adversarial training.

BEWARE: Gradient Masking

- When a defense method prevents a model to reveal meaningful gradients
- At the origin, local gradient towards ϵ_1 is larger compared to ϵ_2 direction
- But loss actually higher in ϵ_2 direction for higher ϵ values
- Many directions orthogonal to ϵ_1 with higher loss at larger ϵ
- Often **unintentional**

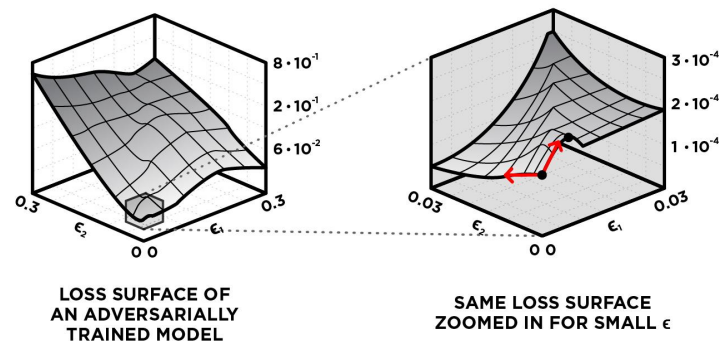
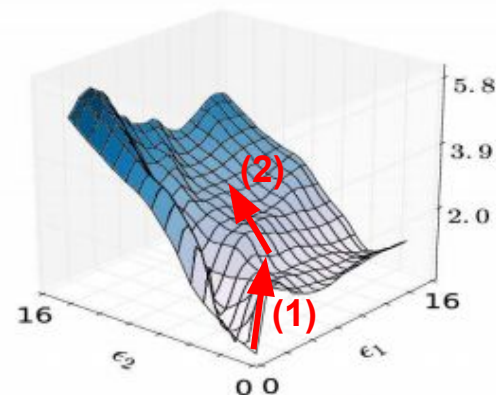


Illustration of loss surface of a model trained with FGSM adversarial training (adapted from Tramèr et al., 2018). Here, ϵ_1 is the direction given by calculating dL/dx , and ϵ_2 is direction orthogonal to ϵ_1 .

REVISIT: R+FGSM

$$x^{\text{adv}} = x' + (\varepsilon - \alpha) \cdot \text{sign}(\nabla_{x'} J(x', y_{\text{true}}))^{(2)}, \quad \text{where} \quad x' = x + \alpha \cdot \text{sign}(\mathcal{N}(\mathbf{0}^d, \mathbf{I}^d))^{(1)}$$

- R+FGSM: add small random perturbation **(1)** **before** calculating the gradient **(2)** (i.e. random start)
- Can circumvent naive implementation of FGSM adversarial training



(a) Loss of model v3_{adv}.

Illustration of how R+FGSM circumvents FGSM adversarial training.

PGD Adversarial Training

- Adversarial training from robust optimization perspective

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

- In other words: find a set of parameter θ that **minimizes** the loss in the **worst-case** scenario (Minimax formulation)
- Empirically showed that adversaries generated using R+BIM, which they called Projected Gradient Descent (PGD) method, are the worst-case adversaries
- In practice: perform adversarial training **only** on PGD adversaries

Agenda

Intro

Adversarial Attacks

Adversarial Defenses

Takeaways

Attack & Defense Methods

Attack strategies:

Whitebox:

- Direct gradient step(s): FGSM, BIM, R+FGSM, DAG
- Gradient-based greedy algorithm: JSMA
- Iterative optimization: L-BFGS, UAP, C&W, Adversarial Eyeglasses, RP_2 , EOT
- Parameterized optimization: ATN

Blackbox:

- Decision boundary approximation: SBA
- Evolutionary algorithm: One Pixel Attack
- Finite difference method: ZOO

Defense strategies:

Data augmentation:

- Adversarial training
- Ensemble adversarial training
- PGD adversarial training

Gradient regularization:

- DCN
- Defensive distillation

Input manipulation:

- PixelDefend

Detection methods:

- Adversary detector networks
- Feature squeezing

Recall our goal:

Minimally modify inputs to
maximize some loss function

Gradient-based Attack Methods

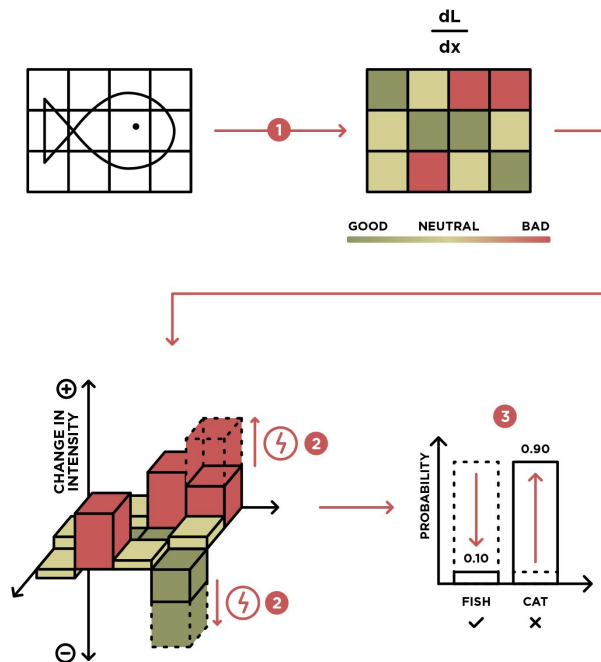


Illustration of gradient-based attacks.

Takeaway Messages

- Arms race between adversarial attacks and defenses: attackers are winning
- Beware of **gradient masking**, often it is unintentional and may give false robustness
 - 7 out of 9 defenses accepted to ICLR 2018 were successfully attacked just few days after acceptance decision date (Athalye et al., 2018)
- Although most attacks focus on virtual world adversaries, there are works that aim to generate adversarial examples in the physical world (e.g. the adversarial eyeglasses, adversarial turtle, etc.)
- The field is very empirical, need more works that can provide guarantee on adversarial robustness (e.g., by providing upper bound of a proposed defense method)

Thank You. Questions?

My blog on adversarial examples:

1. Tricking a Machine Learning into Thinking You're Milla Jovovich
(<https://medium.com/element-ai-research-lab/tricking-a-machine-into-thinking-youre-milla-jovovich-b19bf322d55c>)
2. Securing Machine Learning Models Against Adversarial Attacks
(<https://medium.com/element-ai-research-lab/securing-machine-learning-models-against-adversarial-attacks-b6cd5d2be8e2>)