


Adversarial attack methods

- White box attacks

: The network is "transparent" to the attacker. Both the architecture and weights are known

- Black box attacks

: The attacker has only access to the input and output of the network

- Gray box attacks

: The attacker knows the network architecture but the weights

Goal: $\max_{\delta} \text{loss}(\theta, x + \delta, y)$, where θ : model parameters
 x : input
 y : correct label

Examples: δ that is small w.r.t

- ℓ_p norm
- Rotation and/or translation
- VGG (Visual Geometry Group)

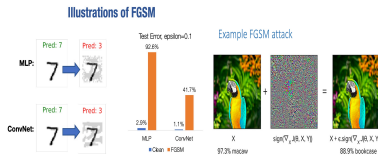
Evasion Attack 1: Fast Gradient Sign Method (FGSM): white box attack

• Classifier (e.g. ResNet-50): $\hat{y} = f(\theta, x)$

• Find adversarial image x' that maximizes the loss: $\mathcal{L}(\hat{x}, y) = \mathcal{L}(f(\theta, \hat{x}), y)$

• Bounded perturbation: $\|x' - x\|_{\infty} \leq \epsilon$, ϵ : the attack strength

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



Evasion Attack 2: Iterative Fast Gradient Sign Method (IFGSM): Whitebox attacks

• Similar to FGSM

• Generates enhanced attacks

$$x^{(m)} = x^{(m-1)} + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x^{(m-1)}, y))$$

With $x^{(0)} = x$, and $x' = x^{(M)}$, where M is the number of iterations

• Both FGSM, and IFGSM are fix-perturbation attacks

Evasion Attack 3: Least Likely (LL) and Iterative Least Likely (ILL) Attacks : White box

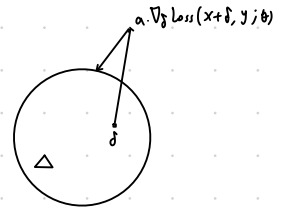
• Similar to FGSM

$x' = x - \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x, y_{\text{LL}}))$, where y_{LL} is the least likely class predicted by the network on clean image x

• Strong attack as it emphasizes least likely class

Evasion Attack 4: Projected Gradient Descent (PGD) : White box

- Recall that we are optimizing : $\max_{\delta \in \Delta} \text{Loss}(x+\delta, y; \theta)$
- We can employ a projected gradient descent method, take gradient step and project back into feasible set Δ : $\delta := P_{\Delta}[\delta + \nabla_{\delta} \text{Loss}(x+\delta, y; \theta)]$
- Projected gradient descent applied to L2 ball, repeat : $\delta := \text{clip}_{\epsilon}[\delta + \alpha \nabla_{\delta} J(\delta)]$
- Slower than FGM (requires multiple iterations), but typically able to find better optima



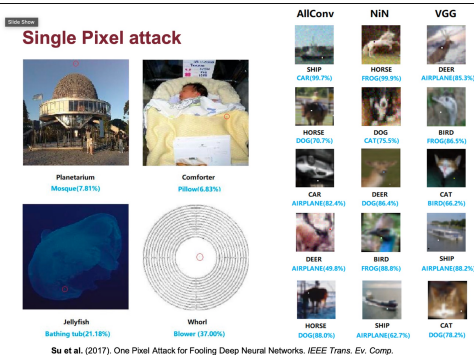
Evasion Attack 5: Carlini and Wagner attack (CW) : White box

- Zero confidence attack
- It tries to find the adversarial image that will be classified as t by solving the problem : $\min \|\delta\|_2^2$
subject to
 $f(x+\delta) = t, x, \delta \in [0, 1]^n$
- Finding the exact solution is difficult
- CW-L2 Attack
- Relaxed version : $\min_{\delta} \|\delta\|_2^2 + c \cdot g(x+\delta)$
s.t.
 $x+\delta \in [0, 1]^n, c \geq 0$
- Letting $Z(x)$ be the neural net activations before the output layer (logits)
 $g(x) = \max_{i \neq t} (Z(x)_i) - Z(x)_t, 0$
- Resists many defense methods

Evasion Attack 6: Universal Adversarial Perturbations

- A single perturbation to fool a network on any image with a high probability (e.g. 0.84)
- Perturbations generalise well across different models, posing a threat to Deep Learning in practice

Evasion Attack 7: Single Pixel Attack



Adversarially Robust Neural Networks - Adversarial Training

Standard generalization: $\mathbb{E}_{(x,y) \sim D} [\text{loss}(\theta, x, y)]$

Adversarially robust generalization: $\mathbb{E}_{(x,y) \sim D} [\max_{\delta \in \Delta} \text{loss}(\theta, x+\delta, y)]$

The outer minimization problem

Inner Maximization:

$$\max_{\delta \in \Delta} \text{loss}(x+\delta, y; \theta)$$



Outer minimization

$$\min_{\theta} \sum_{x,y \in S} \max_{\delta \in \Delta} \text{loss}(x+\delta, y; \theta)$$

① Local search (lower bound on objective)

② Combinatorial Optimization

③ Convex relaxation (upper bound on objective)

① Adversarial training

② Robust training

Darsch's Theorem:

A fundamental result in optimization:

$$\nabla_{\theta} \max_{\delta \in \Delta} \text{loss}(x+\delta, y; \theta) = \nabla_{\theta} \text{loss}(x+\delta^*, y; \theta)$$

$$\text{where } \delta^* = \arg \max_{\delta \in \Delta} \text{loss}(x+\delta, y; \theta)$$

Seems obvious, but it is a very subtle result; means we can optimize through the max by just finding its maximizing value

Adversarial training (Goodfellow et al., 2014)

Repeat

1. Select minibatch B

2. For each $(x,y) \in B$, compute adversarial example $\delta^*(x)$

3. Update θ

$$\theta := \theta - \frac{\alpha}{|B|} \sum_{x,y \in B} \nabla_{\theta} \text{loss}(x+\delta^*(x), y; \theta)$$

Another core idea on Adversarial Robustness: **Denoising**

Defence Mechanisms - Defence GANs

Main idea: Train a GAN that generates unperturbed images. Instead of classifying a given input image, use the closest image generated by the GAN

Pros: Effective against white-box and black-box attacks

No accuracy drop (theoretically)

Cons: Complex method

Difficult to train GAN.

Other types of attack: Adversarial Noise

Adversarial Rotation

Adversarial Photomontage