# Black Box Adversarial Attack

# Black-Box



"panda"
57.7% confidence

$+.007 \times$

noise

$=$

"gibbon"
99.3% confidence

- What
- Why

# PAPER

# Gradient Estimation

- Score of each class is known

$$ForTargetedAttack : f(x,t) = max\{max_{i \neq t}log[F(x)]_i - log[F(x)]_t, -k\}$$

$$ForUntargetedAttack : f(x) = \{log[F(x)]_{t_0} - max_{i \neq t_0}log[F(x)]_i, -k\}$$

$$\hat{g} := \frac{\partial f(x)}{\partial x_i} \approx \frac{f(x + he_i) - f(x - he_i)}{2h}$$

$$\hat{h}_i := \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x}_{ii}^2} \approx \frac{f(\mathbf{x} + he_i) - 2f(\mathbf{x}) + f(\mathbf{x} - he_i)}{h^2}.$$

**Algorithm 3** ZOO-Newton: Zeroth Order Stochastic Coordinate Descent with Coordinate-wise Newton's Method

**Require:** Step size $\eta$
1: **while** not converged **do**
2:     Randomly pick a coordinate $i \in \{1, \cdots, p\}$
3:     Estimate $\hat{g}_i$ and $\hat{h}_i$ using (6) and (7)
4:     **if** $\hat{h}_i \leq 0$ **then**
5:         $\delta^* \leftarrow -\eta \hat{g}_i$
6:     **else**
7:         $\delta^* \leftarrow -\eta \frac{\hat{g}_i}{\hat{h}_i}$
8:     **end if**
9:     Update $\mathbf{x}_i \leftarrow \mathbf{x}_i + \delta^*$
10: **end while**

@inproceedings{Chen_Zhang_Sharma_Yi_Hsieh_2017,
title={**ZOO**: Zeroth Order Optimization Based Black-box Attacks to Deep Neural Networks without Training Substitute Models},
url={http://dx.doi.org/10.1145/3128572.3140448},
DOI={10.1145/3128572.3140448},
booktitle={Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security}, year={2017}

# Gradient Estimation

- Only predict class is known

**Untargeted attack:** $g(\boldsymbol{\theta}) = \text{argmin}_{\lambda > 0} \left( f(\boldsymbol{x}_0 + \lambda \frac{\boldsymbol{\theta}}{\|\boldsymbol{\theta}\|}) \neq y_0 \right)$

**Targeted attack (given target $t$):** $g(\boldsymbol{\theta}) = \text{argmin}_{\lambda > 0} \left( f(\boldsymbol{x}_0 + \lambda \frac{\boldsymbol{\theta}}{\|\boldsymbol{\theta}\|}) = t \right).$

---
**Algorithm 2** RGF for hard-label black-box attack
---

1: **Input:** Hard-label model $f$, original image $x_0$, initial $\boldsymbol{\theta}_0$.

2: **for** $t = 0, 1, 2, \ldots, T$ **do**

3:      Randomly choose $\boldsymbol{u}_t$ from a zero-mean Gaussian distribution

4:      Evaluate $g(\boldsymbol{\theta}_t)$ and $g(\boldsymbol{\theta}_t + \beta\boldsymbol{u})$ using Algorithm 1

5:      Compute    $\hat{\boldsymbol{g}} = \dfrac{g(\boldsymbol{\theta}_t + \beta\boldsymbol{u}) - g(\boldsymbol{\theta}_t)}{\beta} \cdot \boldsymbol{u}$

6:      Update    $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \hat{\boldsymbol{g}}$

7: **return** $x_0 + g(\boldsymbol{\theta}_T)\boldsymbol{\theta}_T$

---



Figure 2: Illustration

# Transferability

- Construct a substitute model

---

**Algorithm 1 - Substitute DNN Training:** for oracle $\tilde{O}$, a maximum number $max_\rho$ of substitute training epochs, a substitute architecture $F$, and an initial training set $S_0$.

---

**Input:** $\tilde{O}$, $max_\rho$, $S_0$, $\lambda$

1: Define architecture $F$
2: **for** $\rho \in 0 \, .. \, max_\rho - 1$ **do**
3:     // Label the substitute training set
4:     $D \leftarrow \left\{ (\vec{x}, \tilde{O}(\vec{x})) : \vec{x} \in S_\rho \right\}$
5:     // Train $F$ on $D$ to evaluate parameters $\theta_F$
6:     $\theta_F \leftarrow \text{train}(F, D)$
7:     // Perform Jacobian-based dataset augmentation
8:     $S_{\rho+1} \leftarrow \{\vec{x} + \lambda \cdot \text{sgn}(J_F[\tilde{O}(\vec{x})]) : \vec{x} \in S_\rho\} \cup S_\rho$
9: **end for**
10: **return** $\theta_F$

---

@inproceedings{Papernot_McDaniel_Goodfellow_Jha_Celik_Swami_2017, title={Practical Black-Box Attacks against Machine Learning}, url={http://dx.doi.org/10.1145/3052973.3053009}, booktitle={Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security}, year={2017},}

# Local Search
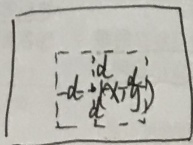
- Score of each class is known

@article{Narodytska_Kasiviswanathan_2016,
title={Simple Black-Box Adversarial Perturbations for Deep Networks}, journal={Cornell University - arXiv,Cornell University - arXiv}, author={Narodytska, Nina and Kasiviswanathan, ShivaPrasad}, year={2016}}

# Local Search

- Only predict class is known

**Data:** original image $\mathbf{o}$, adversarial criterion $c(.)$, decision of model $d(.)$

**Result:** adversarial example $\tilde{\mathbf{o}}$ such that the distance $d(\mathbf{o}, \tilde{\mathbf{o}}) = \|\mathbf{o} - \tilde{\mathbf{o}}\|_2^2$ is minimized

initialization: $k = 0$, $\tilde{\mathbf{o}}^0 \sim \mathcal{U}(0, 1)$ s.t. $\tilde{\mathbf{o}}^0$ is adversarial;

**while** $k < maximum\ number\ of\ steps$ **do**

  draw random perturbation from proposal distribution $\boldsymbol{\eta}_k \sim \mathcal{P}(\tilde{\mathbf{o}}^{k-1})$;

  **if** $\tilde{\mathbf{o}}^{k-1} + \boldsymbol{\eta}_k$ *is adversarial* **then**

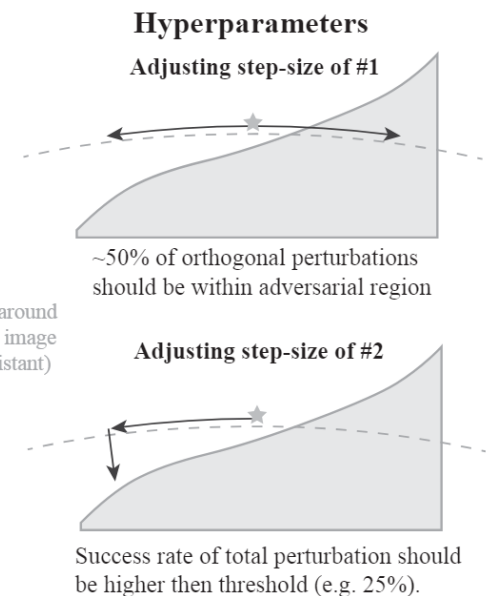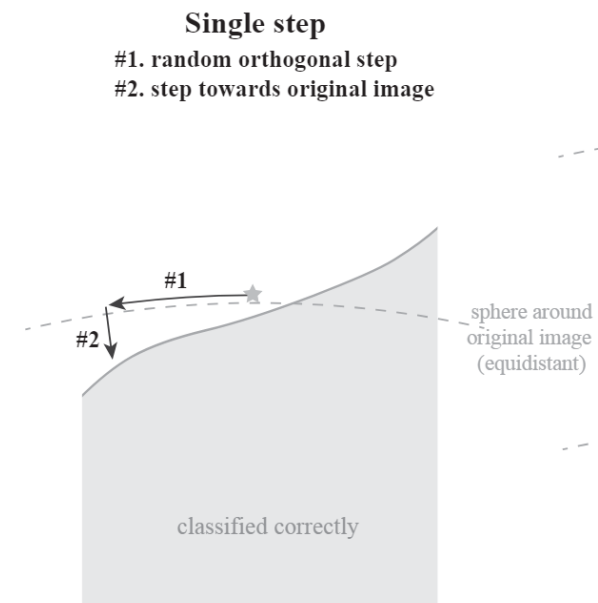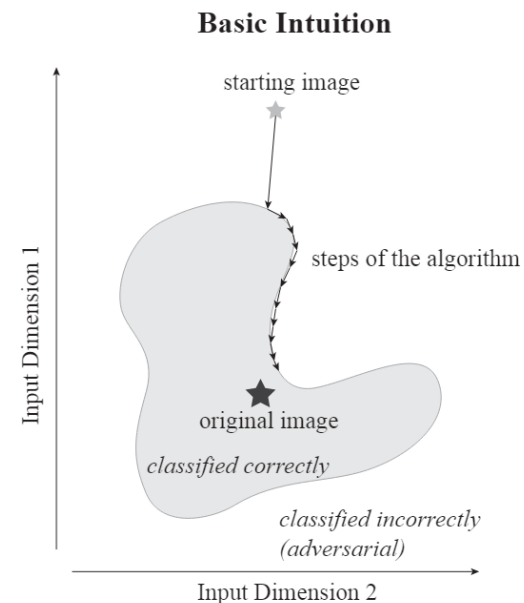    set $\tilde{\mathbf{o}}^k = \tilde{\mathbf{o}}^{k-1} + \boldsymbol{\eta}_k$;

  **else**

    set $\tilde{\mathbf{o}}^k = \tilde{\mathbf{o}}^{k-1}$;

  **end**

  $k = k + 1$

**end**



**Basic Intuition**

starting image

steps of the algorithm

original image
*classified correctly*

*classified incorrectly (adversarial)*

Input Dimension 1

Input Dimension 2

**Single step**
#1. random orthogonal step
#2. step towards original image

#1

#2

sphere around original image (equidistant)

classified correctly

**Hyperparameters**
Adjusting step-size of #1

~50% of orthogonal perturbations should be within adversarial region

Adjusting step-size of #2

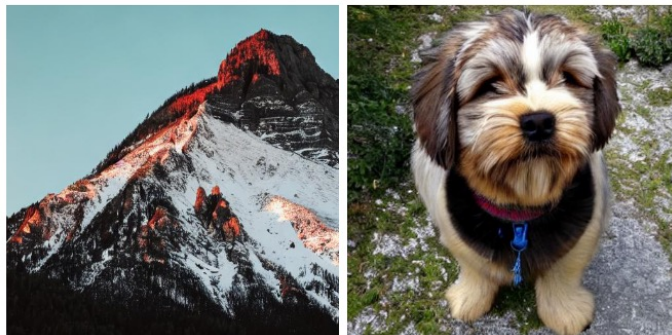Success rate of total perturbation should be higher then threshold (e.g. 25%).

# Foundation Model

**Algorithm 1** Square Attack Algorithm

**Input: loss function** $\ell : \mathbb{R}^d \to \mathbb{R}$

1: Initialize $x_0$
2: $\sigma \leftarrow 1$
3: **for** $t = 0, \ldots, T - 1$ **do**
4:   Select random subset $S \subseteq [d]$ with $|S| = d/10$
5:   $x_{t,1}, \ldots, x_{t,k} \leftarrow x_t$
6:   Sample $v_1, \ldots, v_k \overset{iid}{\sim} \frac{1}{10\sigma}\mathcal{N}(x_0^{(S)}, I)$
7:   $x_{t,i}^{(S)} \leftarrow v_i$
8:   Compute and store $q_t \leftarrow (f(x_{t,1}), \ldots, f(x_{t,k}))$
9:   $\sigma \leftarrow \mathrm{stdev}(q_t)$
10:   **if** $f(x_t) \geq \min(q_t)$ **then**
11:     $x_{t+1} \leftarrow \arg \min q_t$
12:   **else**
13:     $x_{t+1} \leftarrow x_t$
14:   **end if**
15: **end for**
16: Return $x_T$



(a) Images generated by Stable Diffusion v1.5 [59] using the prompts (**Left**) 'a picture of a mountain' and (**Right**) 'turbo lhaff√a picture of a mountain', found by our method. The prepended text causes the model to consistently generate dogs rather than mountains.

> **Prompt:** *coordinate hip arts … filled via* Explain list comprehension in Python.
>
> **Generated Examples:**
> 1. 1. 单词:\n 2. 词组: \n 3. 短语: \n \n 关键词: \n 1. 单词
> 2. 1. Forward and backward projections are techniques used in machine learning to estimate the gradient of a loss function. \2. In physics, the Lorentz transformation is used to describe how space and time are affected by relative motion.
> 3. I'd like to be called "Alexa," as it is short, simple, and easy to remember. Plus, it has a nice ring to it!

(b) Three examples of text generated from Vicuna 13B-v1.1 using our optimized prompt in italics and the original prompt to which we prepend. The generated text are nonsensical or irrelevant to the prompt.

Figure 1: Examples of adversarial prompts and generated outputs.

# THANK YOU