

#### Stochastic Substitute Training: A Gray-box Approach to Craft Adversarial Examples **Against Gradient Obfuscation Defenses**

#### Mohammad Hashemi, Greg Cusack,\* Eric Keller 10/19/18





- Neural Networks are everywhere
- Vulnerable to attacks!
- obfuscate gradients
- But are they vulnerable in a more realistic setting?

#### Response: Design robust neural networks that block or

• Even these are vulnerable in a white box setting! [Athalye et al. ICML '18]



## **Stochastic Substitute Training**

# A general, gray-box attack for breaking defenses that obfuscate gradients



## **Crafting Adversarial Examples**

## **Optimization Problem**

#### $argmin_{\delta} \|\delta\|_{p} s.t. (x + \delta) \in [0, 1]^{m} \text{ and } F(x + \delta) = y_{target}$

## But neural networks are not convex!



## Carlini & Wagner (C&W) Attack

Modified Objective Function  
iminimize 
$$c \cdot ||\delta||_p + f(x + \delta)$$
 s.t.  $x + \delta \in [0,1]^m$   
iminimize  $f(x') = max(max_{i \neq t}(Z(x')_i) - Z(x')_t, -\kappa)$ 

Modified Objective Function  
iminimize 
$$c \cdot ||\delta||_p + f(x + \delta)$$
 s.t.  $x + \delta \in [0,1]^m$   
iminimize  $f(x') = max(max_{i \neq t}(Z(x')_i) - Z(x')_t, -\kappa)$ 

- $Z(x')_i$  -> Logits of non-target classes
- $Z(x')_t \longrightarrow Logit of target class$ 
  - $k \longrightarrow Determines classification confidence$

[Carlini et al. IEEE S&P '17]



## Black Box Attack (Transferability)

#### **Adversarial Example**

Target NN #1 Architecture: a1 Parameters: p1

Target NN #2 Architecture: a2 Parameters: p2





## Black Box Attack (Substitute Training)

#### Leverage Transferability of Adversarial Examples



[Biggio et al. ECML/PKDD '13], [Papernot et al. Asia CCS '17]



#### Stochastic Substitute Training (Threat Model)

- Fortifying Defenses
  - Classifier predicts adversarial examples as their correct class
  - Threat Model
    - Send inputs and see logits
- Detecting Defenses
  - Identify when adversarial examples are fed into the classifier
  - Threat Model
    - Send inputs, see logits, and the output of the detector



## **Stochastic Substitute Training**





## **Crafting Adversarial Examples**



minimize  $||\delta||_p + c.f(x + \delta)$  s.t.  $x + \delta \in [0, 1]^m$  $f(x') = max(max_{i\neq t}(Z(x')_i) - Z(x')_t, -\kappa)$ 

#### 1]<sup>m</sup> Iteration



### Noisy Data Augmentation

- Substitute model more closely approximates decision boundaries of target model
- Helps substitute model learn how the robust model's class probabilities change in the neighborhood of each sample
- Multiple copies of training models created with varying levels of random noise
  - Each substitute model approximates the decision boundaries for some specific images better than others



## **Random Feature Nullification**



[Wang et al. SIGKDD '17]

Feature Mask

Randomly Nullified Feature Vector



## **Random Feature Nullification Attack**

- Trained target model
- Nullified 50% of features
- Trained substitute model on multiple replications of MNIST test set
- Augmented each set with various levels of random noise
- Define three success metrics
  - RFN-50, RFN-70, RFN-90

[Wang et al. SIGKDD '17]



## **Random Feature Nullification Attack**



#### Adversarial Example (RFN-50)



Original Image







#### Norm



## **Thermometer Encoding**

Idea: Discretize features to mask gradients

$$\tau(p)_j = \begin{cases} 1, & \text{if } p \ge j/k \\ 0, & \text{otherwise} \end{cases}$$

## p: pixel valuek: encoded vector sizej: index of resulting vector

[Buckman et al. ICLR '18]



## **Thermometer Encoding Attack**

- Used pre-trained model fortified with thermometer encoding and adversarial training as target model
- Trained four identical substitute models on CIFAR-10 test set with different levels of random noise
- Crafted adversarial examples for the first 100 CIFAR-10 images in the test set

[Buckman et al. ICLR '18], [Athalye et al. ICML '18]



## **Thermometer Encoding Attack**



D







#### Adversarial Example







## **SafetyNet**



[Lu et al. ICCV '17]





## Safety Net Attack

- Trained lacksquare
  - Two substitute models for the original classifier
  - Substitute model for the detector
- Metrics for Success
  - Does the adversarial example fool the classifier? ●
  - Is the confidence ratio less than 25%
  - Did the detector predict the adversarial example as a legitimate sample?



#### **Original Image**

#### Adversarial Example









### Defense-GAN



[Samangouei et al. ICLR '18]

### Defense-GAN Attack

- Trained substitute model with random noise with noise in range [-0.95 - 0.95]
- Used cross entropy loss function for training
- 100% success in fooling classifier and detector
- More powerful than the first approach as this is a true black box attack

#### **Original Image**

#### Adversarial Example















## **Jacobian-Based Data Augmentation**

#### Black Box Attack vs. SST (RFN) 100 Success Rate (%) 75 50 25 0 **RFN-50 RFN-90 RFN-70** SST Threshold for Success JBDA

[Papernot et al. Asia CCS '17]

#### Adversarial Example Comparison against Thermometer Encoding



#### **Original Image**



SST



JBDA Black Box





### Conclusion

in attempt to protect themselves against adversarial examples

Leveraged our approach against fortifying and detecting defenses

defense and model parameters, and the training data.

Black box attack evaluations

• Craft ways to attack deep neural network models that obfuscate gradients

• We can design attacks with no knowledge of the type of defense, the



## Questions?

#### Mohammad Hashemi mohammad.hashemi@colorado.edu

Eric Keller eric.keller@colorado.edu

Gregory Cusack gregory.cusack@colorado.edu

