

# PixelDefend: Leveraging Generative Models to Understand and Defend against Adversarial Examples

Yang Song<sup>1</sup>, Taesup Kim<sup>2</sup>, Sebastian Nowozin<sup>3</sup>, Stefano Ermon<sup>1</sup>, Nate Kushman<sup>3</sup>

<sup>1</sup>Stanford University, <sup>2</sup>Université de Montréal, <sup>3</sup>Microsoft Research

## ABSTRACT

- We show that generative models can be used for detecting adversarially perturbed images and observe that most adversarial examples lie in low probability regions.
- We introduce a novel family of methods for defending against adversarial attacks based on the idea of purification.
- We show that a defensive technique from this family, PixelDefend, can achieve state-of-the-art results on a large number of attacking techniques, improving the accuracy against the strongest adversary on the CIFAR-10 dataset from 32% to 70%.

## ADVERSARIAL EXAMPLES

State-of-the-art classifiers can be fooled by adding quasi-imperceptible noise.



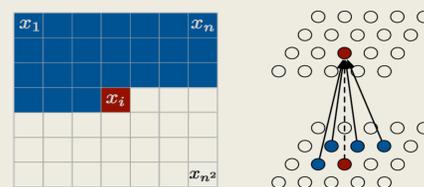
**Figure 1:** Various attacks of an image from CIFAR-10. The text above shows the attacking methods while the text below shows the predicted labels (of a ResNet).

## NEURAL DENSITY MODELS

PixelCNN a convolutional neural network that factorizes  $p(X)$  using the product rule

$$p(X) = \prod_{i=1}^n p(x_i | x_{1:(i-1)}),$$

where the pixel dependencies are in raster scan order.



**Figure 2:** PixelCNN.



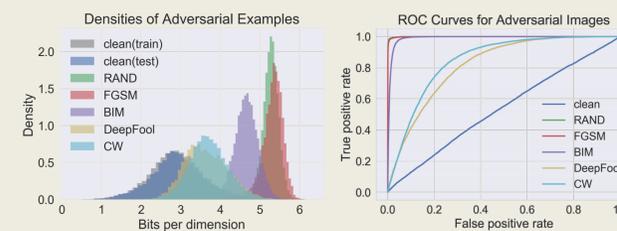
**Figure 3:** Sampled images for Fashion-MNIST and CIFAR-10. Above red line are real images. Below red line are PixelCNN samples.

## DETECTING ADVERSARIAL EXAMPLES

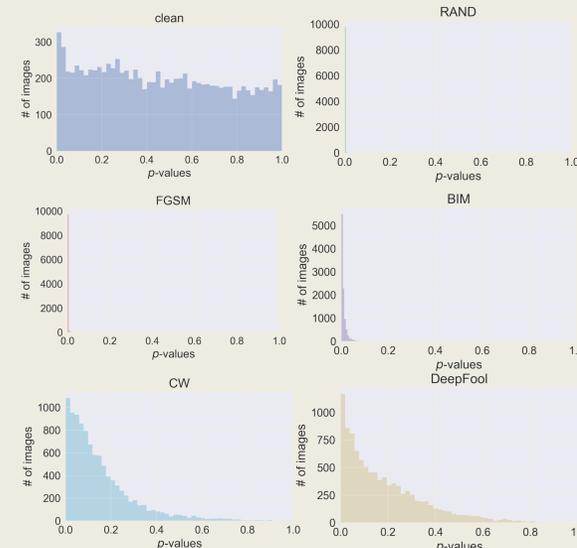
**Observation:** The PixelCNN density of an adversarial example is usually significantly lower than that of a clean example. Therefore,  $p(X)$  can be used as a test statistic to detect adversarial examples.

**Statistical test:** Given an input  $X' \sim q(X)$  and training images  $X_1, X_2, \dots, X_N \sim p_t(X)$ . The null hypothesis is  $H_0: p_t(X) = q(X)$  while the alternative is  $H_1: p_t(X) \neq q(X)$ . The p-value is computed as

$$p\text{-value} = \frac{1}{N+1} \left( \sum_{i=1}^N \mathbb{I}[p(X_i) \leq p(X')] + 1 \right)$$



**Figure 4:** (Left) Likelihoods of different adversarial examples. (Right) ROC curves for detecting various attacks.



**Figure 5:** Distributions of p-values for different attacks.

## PIXELDEFEND

**Intuition:** The harm of adversarial examples might be reduced if they can be modified to have higher likelihood.

### Algorithm 1 PixelDefend

**Input:** Image  $X$ , Defense parameter  $\epsilon_{\text{defend}}$ , Pre-trained PixelCNN model  $p_{\text{CNN}}$   
**Output:** Purified Image  $X^*$

- $X^* \leftarrow X$
- for** each row  $i$  **do**
- for** each column  $j$  **do**
- for** each channel  $k$  **do**
- $x \leftarrow X[i, j, k]$
- Set feasible range  $R \leftarrow [\max(x - \epsilon_{\text{defend}}, 0), \min(x + \epsilon_{\text{defend}}, 255)]$
- Compute the 256-way softmax  $p_{\text{CNN}}(X^*)$ .
- Update  $X^*[i, j, k] \leftarrow \arg \max_{z \in R} p_{\text{CNN}}[i, j, k, z]$
- end for**
- end for**
- end for**

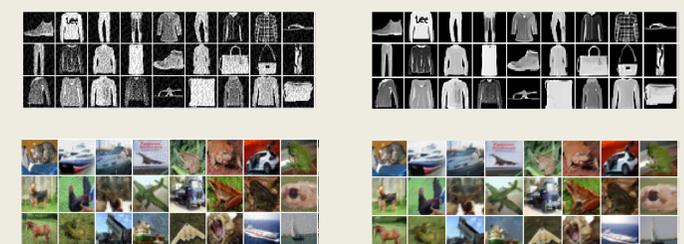
## EXPERIMENTS

Table 1: Fashion MNIST ( $\epsilon_{\text{attack}} = 8/25$ ,  $\epsilon_{\text{defend}} = 32$ )

| NETWORK | TRAINING TECHNIQUE                      | CLEAN | RAND  | FGSM  | BIM   | DEEP FOOL | CW    | STRONGEST ATTACK |
|---------|---|-------|-------|-------|-------|-----------|-------|------------------|
| ResNet  | Normal                                  | 93/93 | 89/71 | 38/24 | 00/00 | 06/06     | 20/01 | 00/00            |
| VGG     | Normal                                  | 92/92 | 91/87 | 73/58 | 36/08 | 49/14     | 43/23 | 36/08            |
| ResNet  | Adversarial FGSM                        | 93/93 | 92/89 | 85/85 | 51/00 | 63/07     | 67/21 | 51/00            |
|         | Adversarial BIM                         | 92/91 | 92/91 | 84/79 | 76/63 | 82/72     | 81/70 | 76/63            |
|         | Label Smoothing                         | 93/93 | 91/76 | 73/45 | 16/00 | 29/06     | 33/14 | 16/00            |
|         | Feature Squeezing                       | 84/84 | 84/70 | 70/28 | 56/25 | 83/83     | 83/83 | 56/25            |
| ResNet  | Adversarial FGSM + Feature Squeezing    | 88/88 | 87/82 | 80/77 | 70/46 | 86/82     | 84/85 | 70/46            |
|         | Normal + PixelDefend                    | 88/88 | 88/89 | 85/74 | 83/76 | 87/87     | 87/87 | 83/74            |
| VGG     | Normal + PixelDefend                    | 89/89 | 89/89 | 87/82 | 85/83 | 88/88     | 88/88 | 85/82            |
| ResNet  | Adversarial FGSM + PixelDefend          | 90/89 | 91/90 | 88/82 | 85/76 | 90/88     | 89/88 | 85/76            |
|         | Adversarial FGSM + Adaptive PixelDefend | 91/91 | 91/91 | 88/88 | 85/84 | 89/90     | 89/84 | 85/84            |

Table 2: CIFAR-10 ( $\epsilon_{\text{attack}} = 2/8/16$ ,  $\epsilon_{\text{defend}} = 16$ )

| NETWORK | TRAINING TECHNIQUE                      | CLEAN    | RAND     | FGSM     | BIM      | DEEP FOOL | CW       | STRONGEST ATTACK |
|---------|---|----------|----------|----------|----------|-----------|----------|------------------|
| ResNet  | Normal                                  | 92/92/92 | 92/87/76 | 33/15/11 | 10/00/00 | 12/06/06  | 07/00/00 | 07/00/00         |
| VGG     | Normal                                  | 89/89/89 | 89/88/80 | 60/46/30 | 44/02/00 | 57/25/11  | 37/00/00 | 37/00/00         |
| ResNet  | Adversarial FGSM                        | 91/91/91 | 90/88/84 | 88/91/91 | 24/07/00 | 45/00/00  | 20/00/07 | 20/00/00         |
|         | Adversarial BIM                         | 87/87/87 | 87/87/86 | 80/52/34 | 74/32/06 | 79/48/25  | 70/42/08 | 74/32/06         |
|         | Label Smoothing                         | 92/92/92 | 91/88/77 | 73/54/28 | 59/08/01 | 56/20/10  | 30/02/02 | 30/02/01         |
|         | Feature Squeezing                       | 84/84/84 | 83/82/76 | 31/20/18 | 13/00/00 | 75/75/75  | 78/78/78 | 13/00/00         |
| ResNet  | Adversarial FGSM + Feature Squeezing    | 86/86/86 | 85/84/81 | 73/07/55 | 55/02/00 | 85/85/85  | 83/83/83 | 55/02/00         |
|         | Normal + PixelDefend                    | 85/85/88 | 82/83/84 | 73/46/24 | 71/46/25 | 80/80/80  | 78/78/78 | 71/46/24         |
| VGG     | Normal + PixelDefend                    | 82/82/82 | 82/82/84 | 80/62/52 | 80/61/48 | 81/76/76  | 81/79/79 | 80/61/48         |
| ResNet  | Adversarial FGSM + PixelDefend          | 88/88/86 | 86/86/87 | 81/68/67 | 81/69/56 | 85/85/85  | 84/84/84 | 81/69/56         |
|         | Adversarial FGSM + Adaptive PixelDefend | 90/90/90 | 86/87/87 | 81/70/67 | 81/70/56 | 82/81/82  | 81/80/81 | 81/70/56         |



**Figure 6:** Adversarial images (left) and purified images after PixelDefend (right).