

Image recognition: Defense adversarial attacks using Generative Adversarial Network (GAN)

Speaker: Guofei Pang

Division of Applied Mathematics
Brown University

Presentation after reading the paper:

Ilyas, Andrew, et al. "The Robust Manifold Defense: Adversarial Training using Generative Models." arXiv preprint arXiv:1712.09196 (2017).

Outline

- Adversarial attacks
- Generative Adversarial Network (GAN)
- How to defense attacks using GAN
- Numerical results

Adversarial Attacks



x

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$

“nematode”

8.2% confidence

=



$x +$

$\epsilon \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$

“gibbon”

99.3 % confidence

Adversarial Attacks

$P(\text{man}) > 0.99$



$P(\text{woman}) > 0.99$



Adversarial Attacks

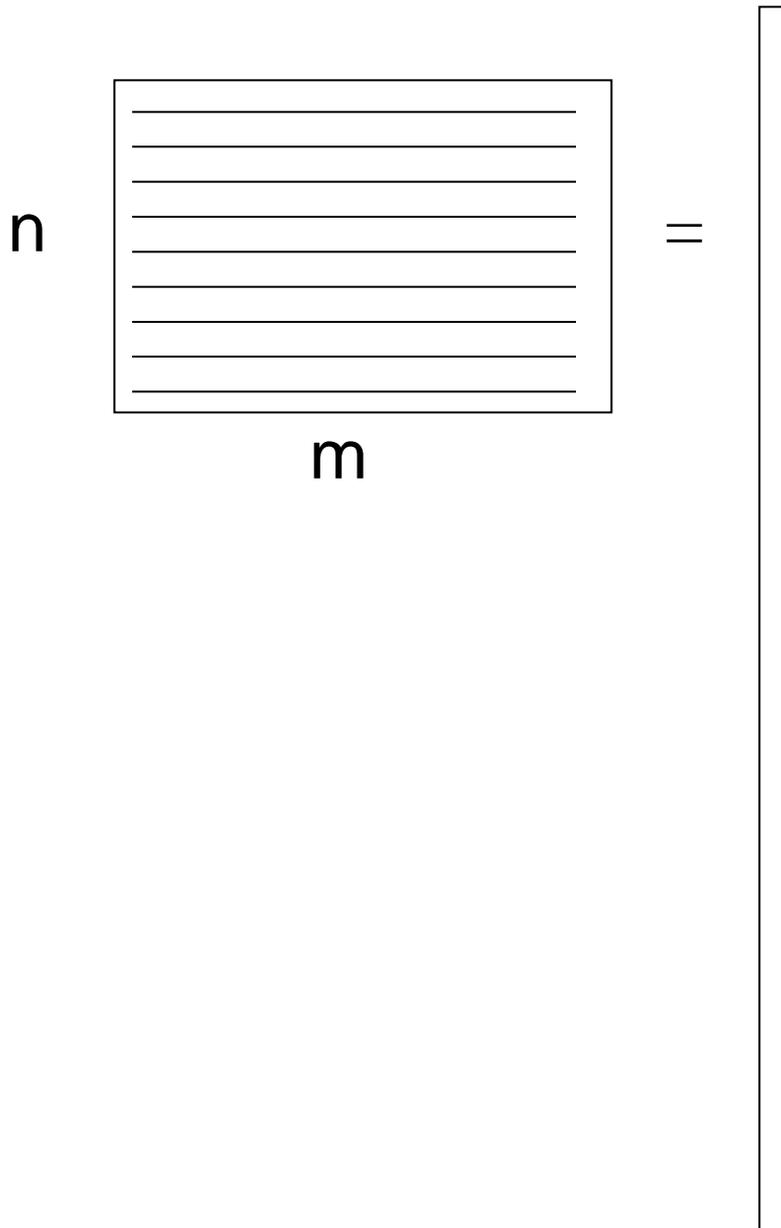


Image as a vector:

$$\mathbf{x} = \{\mathbf{x}_j\},$$
$$j = 1, 2, \dots, n * m$$

Adversarial Attacks



$$\|x_1 - x_2\|_2 < \epsilon_0 \implies |C(x_1) - C(x_2)| > f_0$$

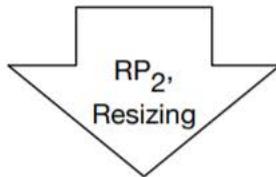
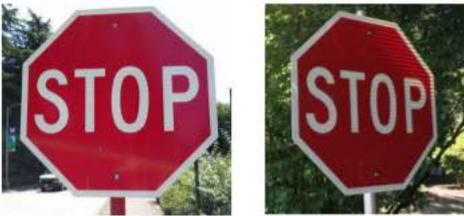
Adversarial examples for a classifier $C()$:

- A pair of input x_1 and x_2
- A person says they are of the same class
- But a classifier will they are completely different!

Adversarial Attacks

Robust Physical Perturbation

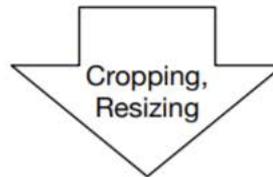
Sequence of physical road signs under different conditions



Different types of physical adversarial examples

Lab (Stationary) Test

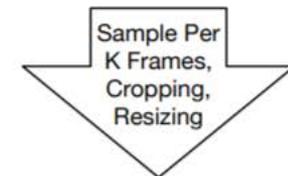
Physical road signs with adversarial perturbation under different conditions



Stop Sign → Speed Limit Sign

Field (Drive-By) Test

Video sequences taken under different driving speeds

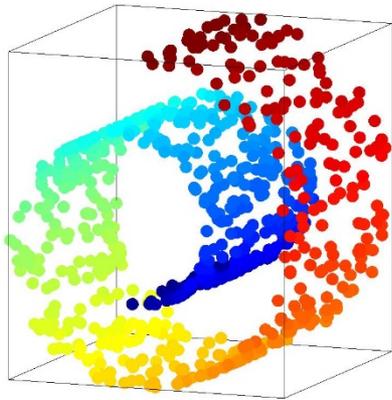


Stop Sign → Speed Limit Sign

Why does classifier become fool for these examples?

Adversarial Attacks

Why does classifier become fool for these examples?



An intuition from the authors:

- Natural image: Low-dimensional manifold
- Noisy image: High-dimensional manifold
- High dimensionality is **tough** for classifier.



+ .007 ×



=



Generative adversarial network (GAN)

- x and x' have similar PDF
- $G()$ has learned the underlying distribution of image dataset after training GAN
- The DNN $G()$ is a nonlinear mapping from low-dimensional space, z , to high-dimensional space, x'

Original image x

GAN



Synthetic image /Generative model

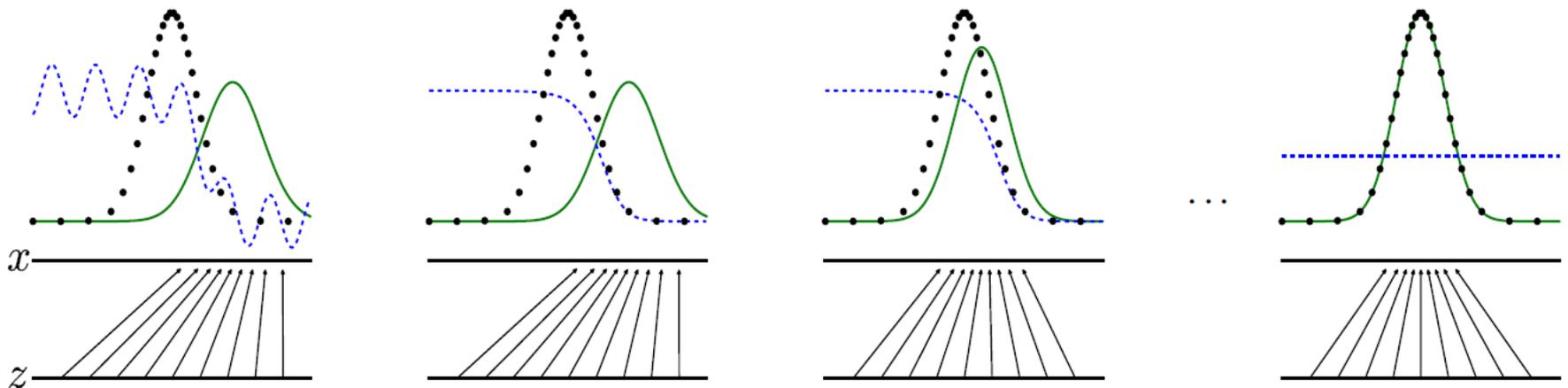
$$x' = G(z)$$

Generator $G(z)$

Noisy input z , say, $z \sim N(0, I)$

Generative adversarial network (GAN)

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$



- **Convergence state: $p_{\text{data}}(\mathbf{x}) = p_G(\mathbf{x})$**
- Green solid line: probability density function (PDF) of the generator $G()$
- Black dotted line: PDF of original image \mathbf{x} , i.e., $p_{\text{data}}(\mathbf{x})$
- Blue dash line: PDF of discriminator $D()$

Generative adversarial network (GAN)

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log(1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(\mathbf{z}^{(i)}))).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

How to defense attacks using GAN

- **G()** is pre-trained and has learned the underlying distribution of the training (image) dataset after training GAN

Synthetic image $\mathbf{x}' = \mathbf{G}(z^*)$
(**Preserve low-dimensional manifold**)

Original image \mathbf{x}
(Could include high-dimensional manifold when noise enters)

Invert and Classify

Classifier C()

$$z^* = \arg \min_z \|G(z) - x\|_2$$

How to defense attacks using GAN

- **G()** is pre-trained and has learned the underlying distribution of the training (image) dataset after training GAN

Enhanced Invert and Classify

Synthetic image $\mathbf{x}' = \mathbf{G}(z^*)$
(Preserve low-dimensional manifold)



Classifier $C()$
(retrain the classifier)

Upper bound of attack magnitude

Classification loss

$$\inf_{\theta} \mu \left(\sup_{z, z'} \|C_{\theta}(G(z)) - C_{\theta}(G(z'))\|_2^2 \right) + (1 - \mu) \left(\frac{1}{N} \sum_{i=1}^N f(y^{(i)}, C_{\theta}(x^{(i)})) \right)$$

s.t. $\|G(z) - G(z')\|_2^2 \leq \eta^2$.

$$z^* = \arg \min_z \|G(z) - x\|_2$$

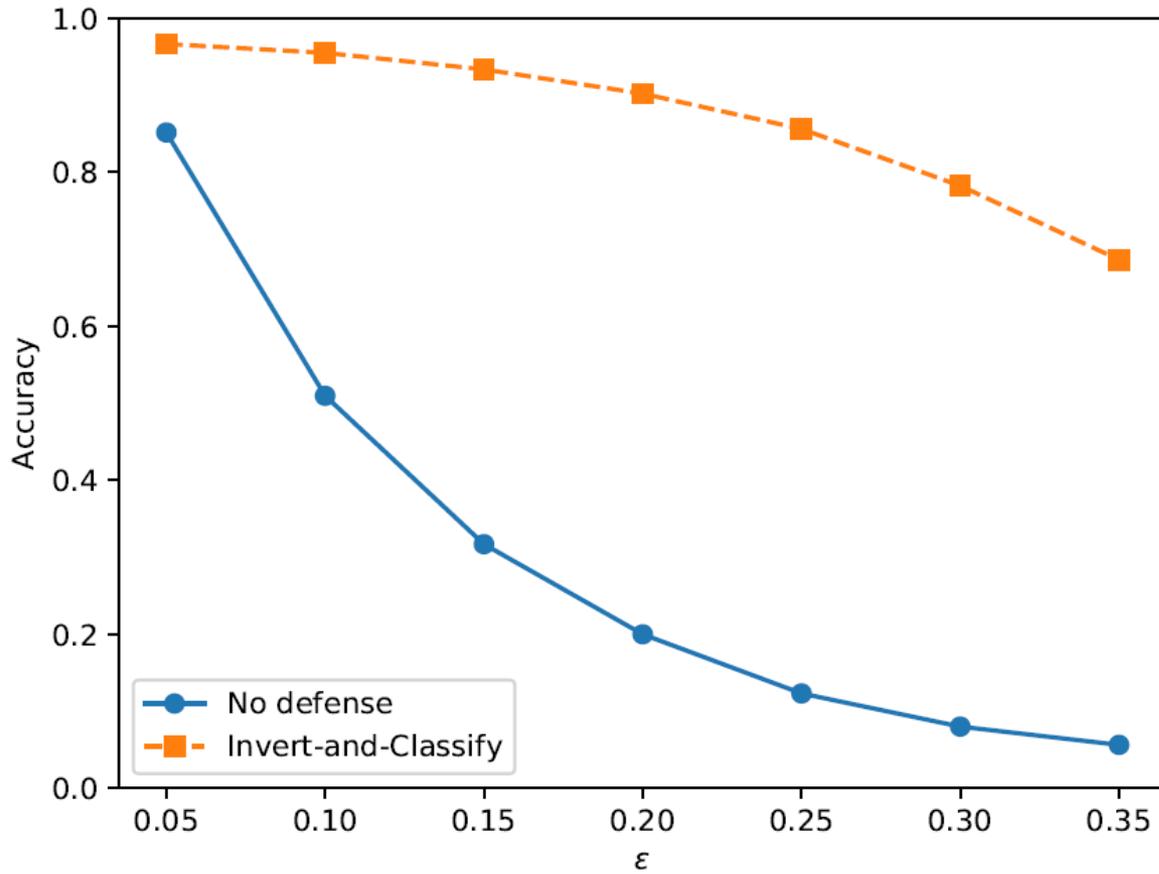
Numerical results

First-order classifier attacks for handwritten digit classification

$$\delta = +\epsilon \cdot \text{sign}(\nabla_X L(y, C_\theta(X))|_{X=x}),$$

Numerical results

First-order classifier attacks for handwritten digit classification



$$\delta = +\epsilon \cdot \text{sign}(\nabla_X L(y, C_\theta(X))|_{X=x}),$$

Numerical results

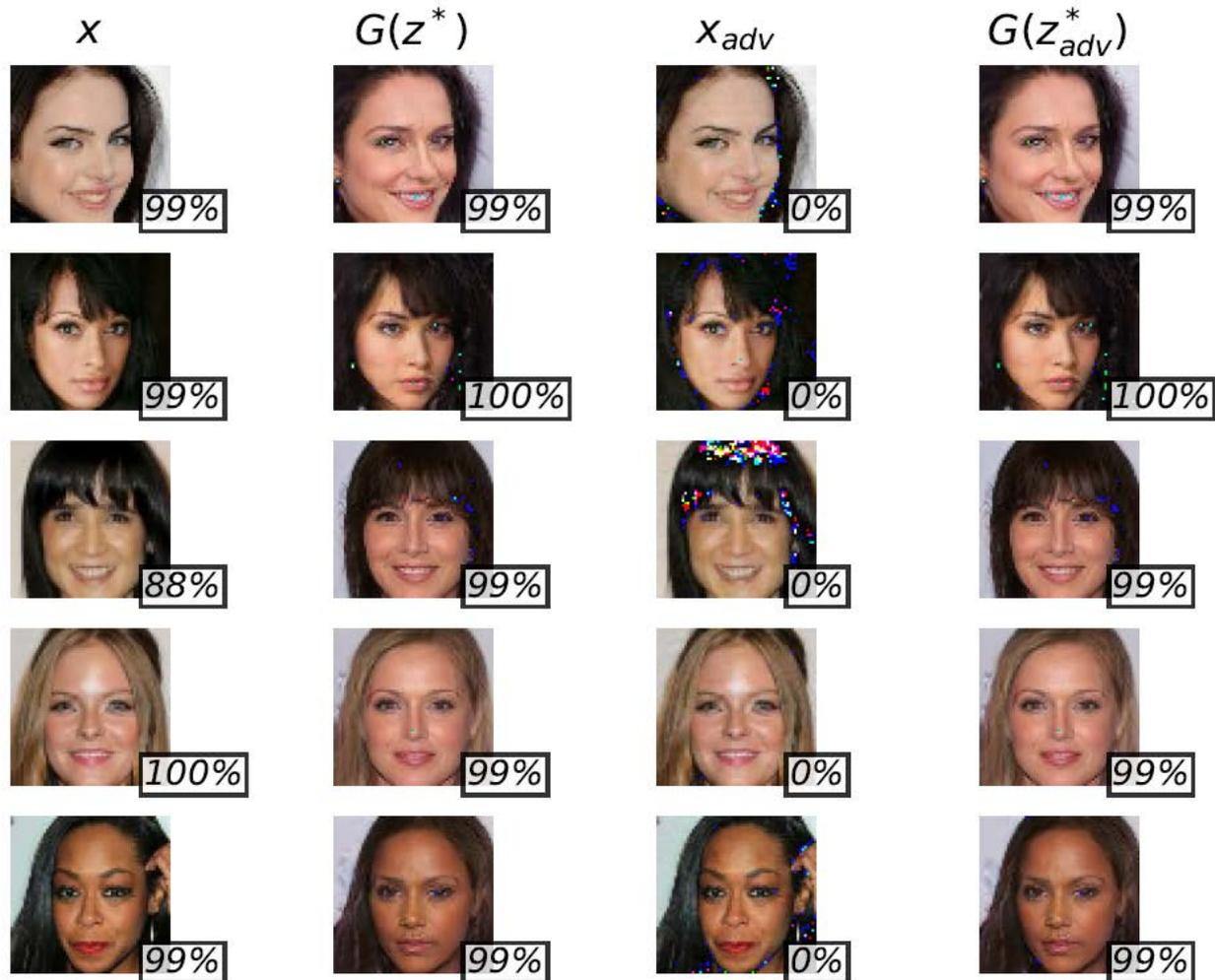
First-order classifier attacks for handwritten digit classification

ϵ	No defense	Invert and Classify
Clean Data	97%	84%
FGSM ($\epsilon = 0.05$)	1%	82%
FGSM ($\epsilon = 0.1$)	0%	80%
FGSM ($\epsilon = 0.2$)	0%	73%
Carlini-Wagner ℓ_2	0%	77%
Carlini-Wagner ℓ_0	0%	65%
Carlini-Wagner ℓ_∞	0%	66%

$$\delta = +\epsilon \cdot \text{sign}(\nabla_X L(y, C_\theta(X))|_{X=x}),$$

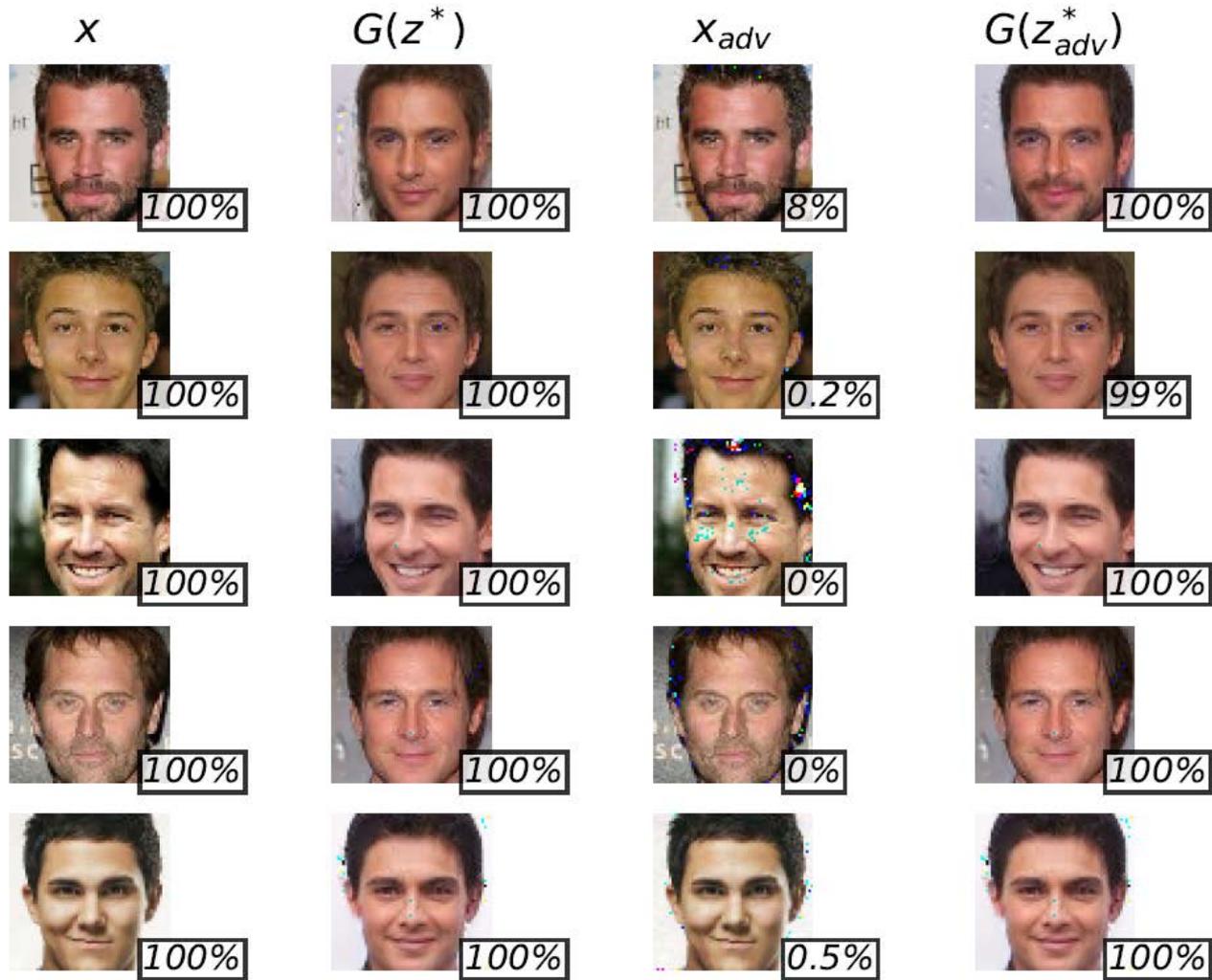
Numerical results

First-order classifier attacks for gender classification



Numerical results

First-order classifier attacks for gender classification



Numerical results

Substitute model attacks

Results from [Invert and Classify](#)

$P(\text{man}) > 0.99$



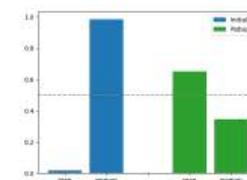
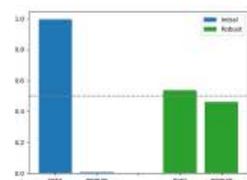
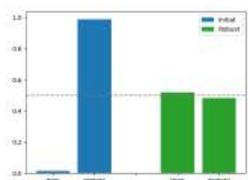
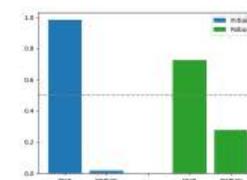
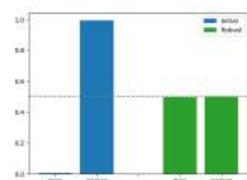
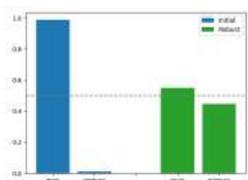
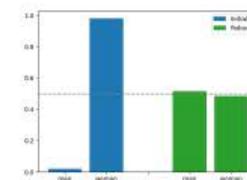
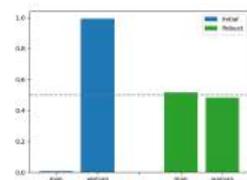
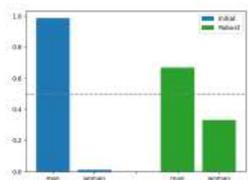
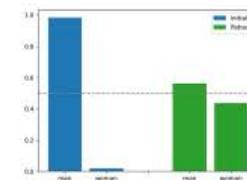
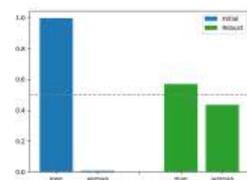
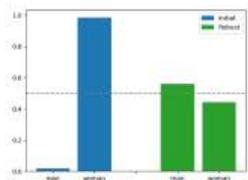
$P(\text{woman}) > 0.99$



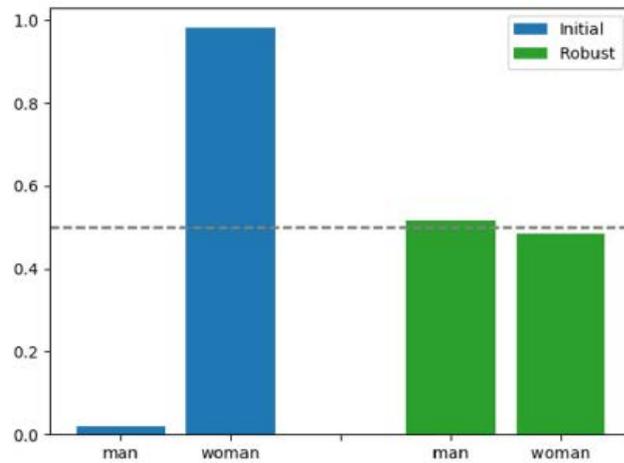
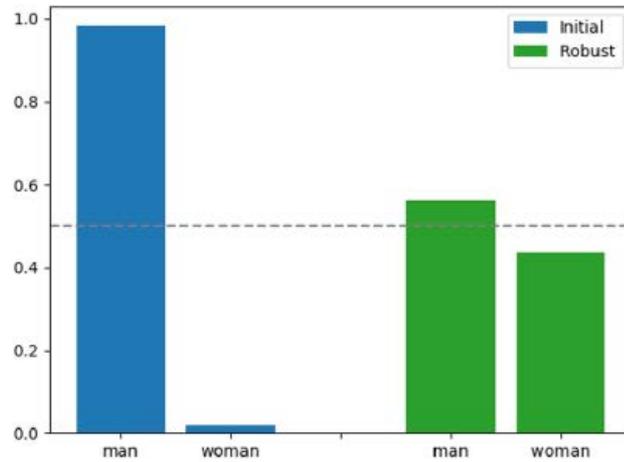
Numerical results

Comparison between
Invert and Classify and Enhanced Invert and Classify

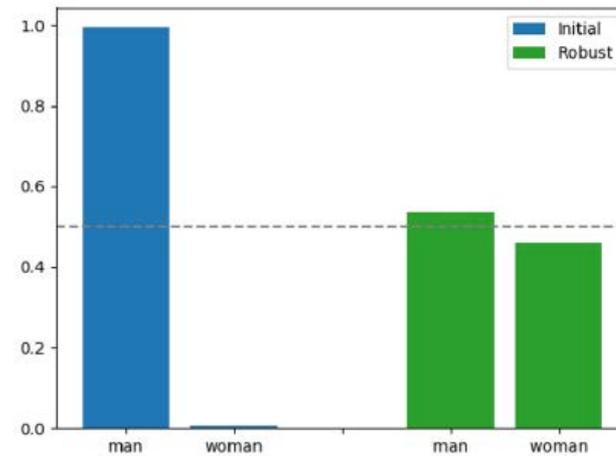
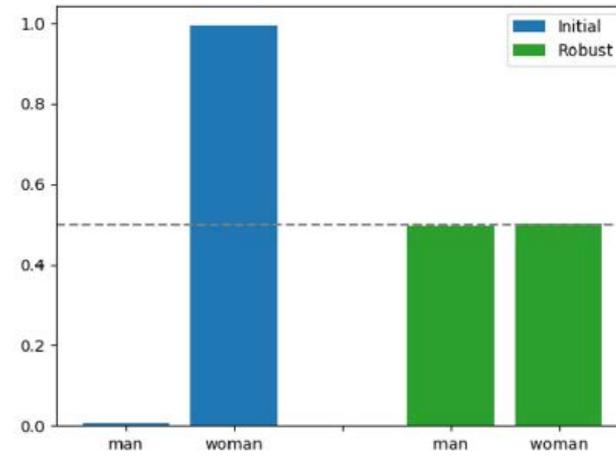
Numerical results



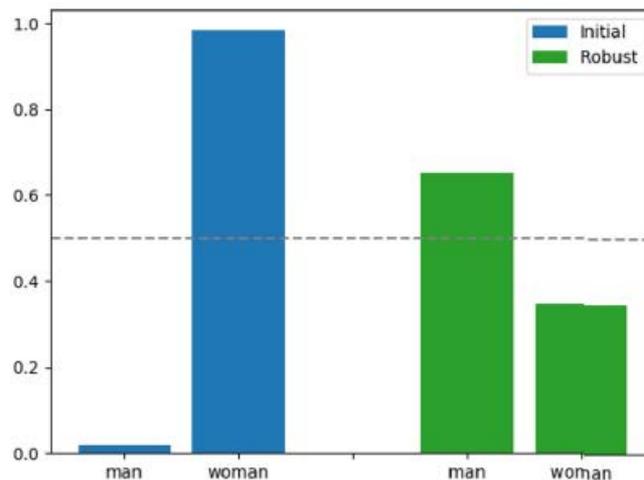
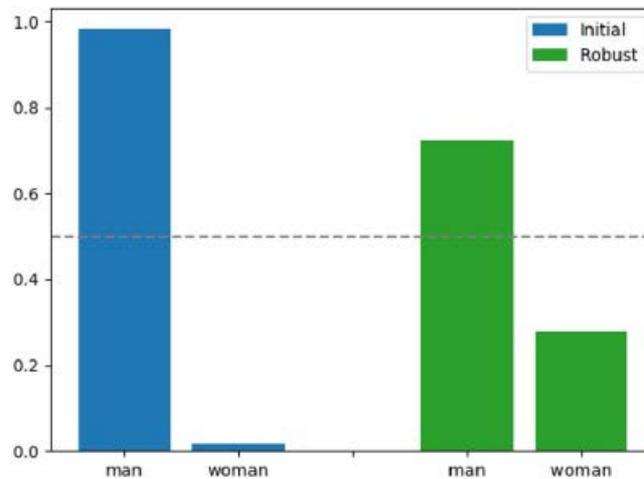
Numerical results



Numerical results



Numerical results



- GAN for regression problems?
- GAN versus other neural networks?
- One defense strategy for all types of attacks?