

## Adversarial

$(x, y_{true})$  - object  $f_{\theta}(x)$  - neural network

$L(\theta, x, y)$  - loss value

$f_{\theta}(x) = y_{true}$   $\tilde{x} \approx x: f_{\theta}(\tilde{x}) \neq y_{true}$

$\|\tilde{x} - x\|_{\infty} < \frac{\epsilon}{255}$   $\epsilon$  - attack budget

$$\|\tilde{x} - x\|_{\infty} = \max_i |\tilde{x}_i - x_i|$$

- White Box - gradient are available
- Black Box - gradient not available

- Targeted attack:  $f_{\theta}(x) = y_{target}$
- Untargeted attack:  $f_{\theta}(x) \neq y_{true}$

FGSM (Fast Gradient Sign Method)

- untargeted

$$\tilde{x} = x + \epsilon \cdot \text{sign}(\nabla_x L(\theta, x, y_{true}))$$

- targeted

$$\tilde{x} = x - \epsilon \text{sign}(\nabla_x L(\theta, x, y_{target}))$$

## I-FGSM (iterative FGSM)

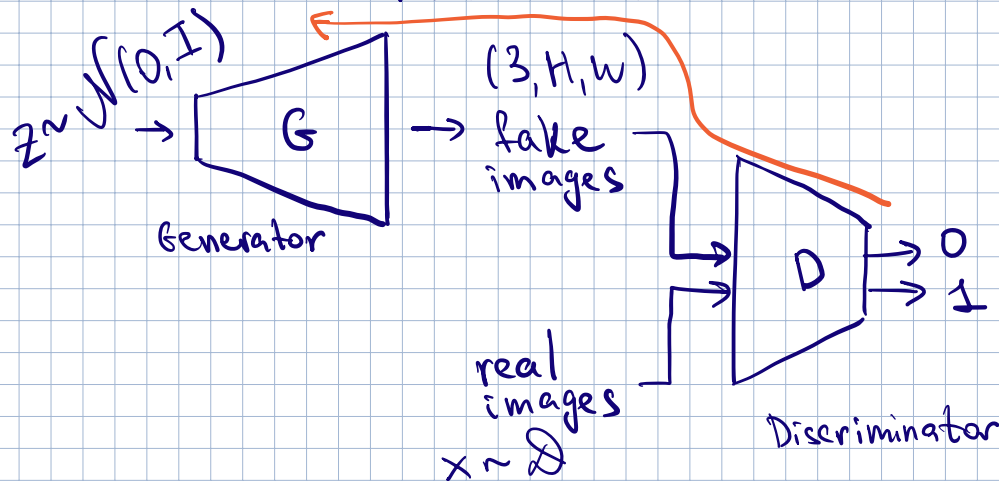
$$\tilde{x}_0 = x$$

for  $i = 1, \dots, N$

$$\tilde{x}_i = \tilde{x}_{i-1} + \frac{\epsilon}{N} \text{sign} \left( D_x h(\theta, \tilde{x}_{i-1}, y_{\text{true}}) \right)$$

$$\tilde{x} = \tilde{x}_N$$

## Generative Adversarial Networks (GANs)



$$D(x) \in [0, 1]$$

$$\min_G \max_D \left\{ \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim \mathcal{N}(0, I)} [\log (1 - D(G(z)))] \right\}$$

## StyleGAN

### cGAN (conditional GAN)

$$G_\theta(z, c)$$

$D_\phi^1(x)$  - head for GAN loss

$D_\phi^2(x)$  - head for classification loss

Img 2 Img models

Pix 2 Pix

$$X_S \sim p_S(x), X_T \sim p_T(x)$$

$$G(X_S) \approx X_T$$

$$\text{MSE loss: } \mathcal{L}_{\text{MSE}} = \mathbb{E}_{X_S, X_T} \|G(X_S) - X_T\|_2^2 \rightarrow \min_G$$

$$\text{GAN loss: } \mathcal{L}_{\text{GAN}}(G, D)$$

$$\text{MSE + GAN loss: } \alpha \mathcal{L}_{\text{MSE}} + (1-\alpha) \mathcal{L}_{\text{GAN}}$$

$$\mathcal{L}_{\text{GAN}} = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{x \sim p_{\text{gen}}} [\log (1-D(x))] \rightarrow \max_D$$

$$D^*(x) = ?$$

$$\mathcal{L}_{\text{GAN}} = \int_{x \in \mathcal{X}} [p_{\text{data}}(x) \log D(x) + p_{\text{gen}}(x) \log (1-D(x))] dx$$

$$g(c) = p_{\text{data}}(x) \cdot \log c + p_{\text{gen}}(x) \log (1-c)$$

$$g'(c) = \frac{p_{\text{data}}(x)}{c} - \frac{p_{\text{gen}}(x)}{1-c} = 0$$

$$(1-c) p_{\text{data}}(x) - c p_{\text{gen}}(x) = 0$$

$$D^*(x) = c = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\text{gen}}(x)}$$

$$\times \quad \begin{array}{l} p_{\text{gen}}(x) \gg 1 \\ p_{\text{data}}(x) \approx 0 \end{array} \quad D^*(x) \approx 0$$

$$L_{GAN} = \int \left[ P_{data}(x) \log \frac{P_{data}(x)}{P_{data}(x) + P_{gen}(x)} + P_{gen}(x) \log \frac{P_{gen}(x)}{P_{data}(x) + P_{gen}(x)} + P_{data}(x) \log 2 - P_{data}(x) \log 2 + P_{gen}(x) \log 2 - P_{gen}(x) \log 2 \right] dx$$

$$\int \left[ P_{data}(x) \log \frac{2 P_{data}(x)}{P_{data}(x) + P_{gen}(x)} + P_{gen}(x) \log \frac{2 P_{gen}(x)}{P_{data}(x) + P_{gen}(x)} \right] dx - \log 2 \int (P_{data}(x) + P_{gen}(x)) dx =$$

$= 2$

$$= KL(P_{data} || \frac{P_{data} + P_{gen}}{2}) + KL(P_{gen} || \frac{P_{data} + P_{gen}}{2}) - 2 \log 2 = 2 JS(P_{data}, P_{gen}) - 2 \log 2$$

Mode collapsing

