# Generative Networks

STAT991
Wenting Zhan

# Claim

This slides are modified and integrated from:

- CS231n Stanford University
- STAT991 University of Pennsylvania (Fall 2018 Lecture 4)
- CIS700-004 University of Pennsylvania (Spring 2019 Lecture 7M & 7W)
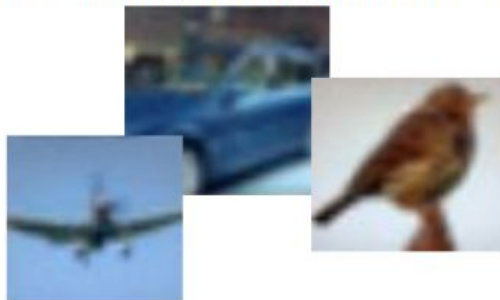
# Agenda

- Generative Models
- Variational Autoencoders (VAE)
- Generative Adversarial Networks (GAN)
    - Basic Mechanism: Generator and Discriminator
    - Dual Loss Function
    - Major Issues
        - Mode Collapse
        - Why are GANs so hard to train
- Types of GAN
    - Deep Convolution GAN
    - Cycle GAN
    - Etc.
- Applications of GAN

# Agenda

- **Generative Models**
- Variational Autoencoders (VAE)
- Generative Adversarial Networks (GAN)
  - Basic Mechanism: Generator and Discriminator
  - Dual Loss Function
  - Major Issues
    - Mode Collapse
    - Why are GANs so hard to train
- Types of GAN
  - Deep Convolution GAN
  - Cycle GAN
  - Etc.
- Applications of GAN

# Generative Models

- Given training data, generate new samples from same distribution



Training data $\sim p_{data}(x)$      Generated samples $\sim p_{model}(x)$

Want to learn $p_{model}(x)$ similar to $p_{data}(x)$

- Addresses density estimation, a core problem in unsupervised learning

- **Several flavors:**
  - Explicit density estimation: explicitly define and solve for $p_{model}(x)$
  - Implicit density estimation: learn model that can sample from $p_{model}(x)$ w/o explicitly defining it

# Agenda

- Generative Models
- Variational Autoencoders (VAE)
- Generative Adversarial Networks (GAN)
  - Basic Mechanism: Generator and Discriminator
  - Dual Loss Function
  - Major Issues
    - Mode Collapse
    - Why are GANs so hard to train
- Types of GAN
  - Deep Convolution GAN
  - Cycle GAN
  - Etc.
- Applications of GAN

# Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

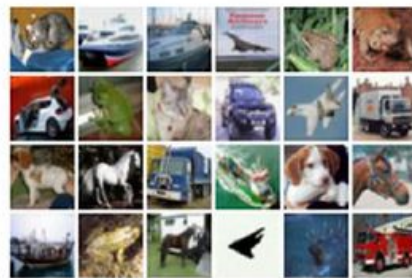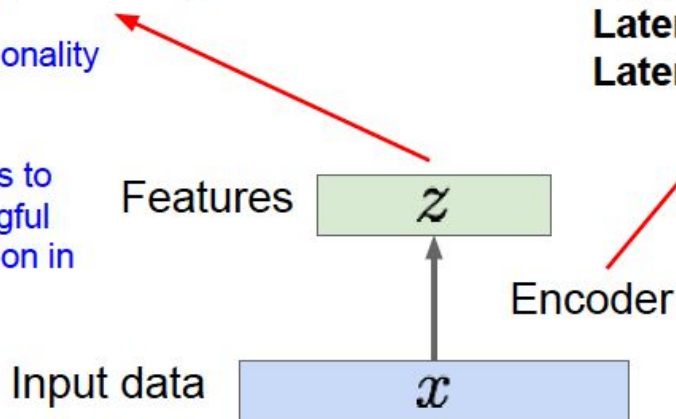**z** usually smaller than **x** (dimensionality reduction)

Q: Why dimensionality reduction?

A: Want features to capture meaningful factors of variation in data

**Originally**: Linear + nonlinearity (sigmoid)
**Later**: Deep, fully-connected
**Later**: ReLU CNN

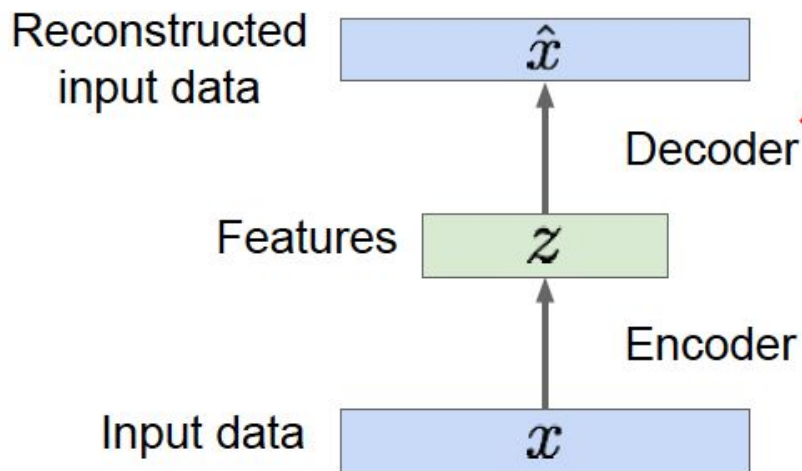Features $z$

Encoder

Input data $x$

# Autoencoders

How to learn this feature representation?
Train such that features can be used to reconstruct original data
"Autoencoding" - encoding itself

Reconstructed
input data

$\hat{x}$

Decoder

Features

$z$

Encoder

Input data

$x$

**Originally**: Linear +
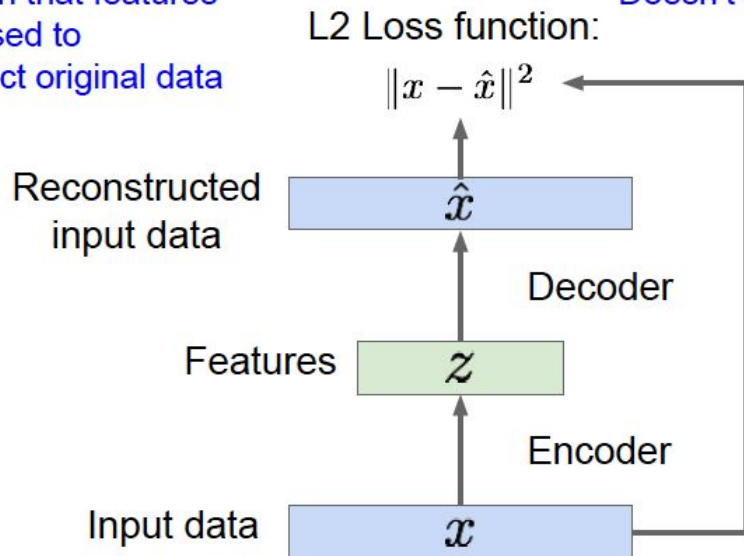nonlinearity (sigmoid)
**Later**: Deep, fully-connected
**Later**: ReLU CNN (upconv)

# Autoencoder

Train such that features can be used to reconstruct original data

Doesn't use labels!

L2 Loss function:
$$\|x - \hat{x}\|^2$$

Reconstructed input data — $\hat{x}$

Decoder

Features — $z$

Encoder

Input data — $x$

Reconstructed data

**Encoder**: 4-layer conv
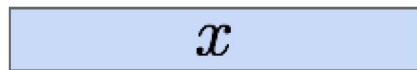**Decoder**: 4-layer upconv

Input data

# VAE

Probabilistic spin on autoencoders - will let us sample from the model to generate data!

Assume training data $\{x^{(i)}\}_{i=1}^{N}$ is generated from underlying unobserved (latent) representation **z**

**Intuition** (remember from autoencoders!): **x** is an image, **z** is latent factors used to generate **x**: attributes, orientation, etc.
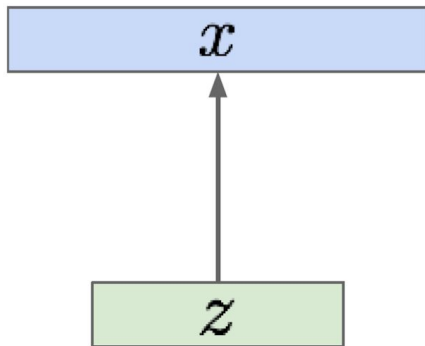
Sample from true conditional
$p_{\theta^*}(x \mid z^{(i)})$

$$x$$

Sample from true prior
$p_{\theta^*}(z)$

$$z$$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# VAE

We want to estimate the true parameters $\theta*$ of this generative model.

Sample from true conditional

$p_{\theta*}(x \mid z^{(i)})$

Sample from true prior

$p_{\theta*}(z)$

| $x$ |
| :-: |

↑

| $z$ |
| :-: |

**How should we represent this model?**

Choose prior p(z) to be simple, e.g. Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$
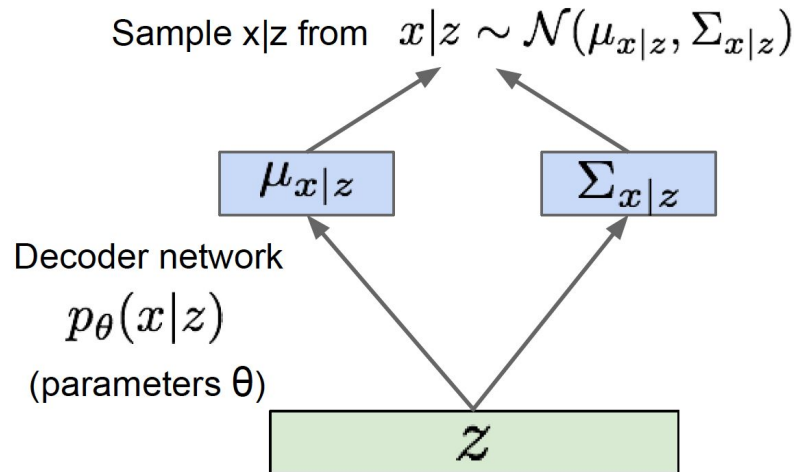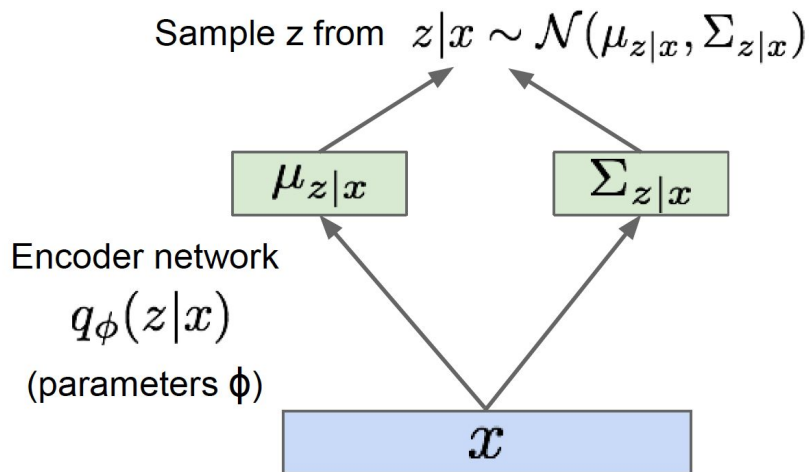
Data likelihood: $p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$ ✔ ✔

Intractible to compute p(x|z) for every z!

Posterior density also intractable: $p_\theta(z|x) = p_\theta(x|z) p_\theta(z) / p_\theta(x)$ ✔ ✔

Intractable data likelihood

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# VAE

Solution: In addition to decoder network modeling $p_\theta(x|z)$, define additional encoder network $q_\phi(z|x)$ that approximates $p_\theta(z|x)$

Sample z from $\quad z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

Sample x|z from $\quad x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$\mu_{z|x}$ $\qquad$ $\Sigma_{z|x}$

$\mu_{x|z}$ $\qquad$ $\Sigma_{x|z}$

Encoder network
$q_\phi(z|x)$
(parameters φ)

Decoder network
$p_\theta(x|z)$
(parameters θ)

$x$

$z$

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# VAE Loss Function

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} \left[ \log p_\theta(x^{(i)}) \right] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ \log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

Make approximate posterior distribution close to prior

$$= \mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[ \log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})} \right] \quad \text{(Logarithms)}$$

Reconstruct Input data

$$= \boxed{\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right]} - \boxed{D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))} + D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z \mid x^{(i)}))$$

Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :( But we know KL divergence always >= 0.

# VAE

Putting it all together: maximizing the likelihood lower bound

Maximize likelihood of original input being reconstructed

$$\mathbf{E}_z \left[ \log p_\theta(x^{(i)} \mid z) \right] - D_{KL}(q_\phi(z \mid x^{(i)}) \,\|\, p_\theta(z))$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

For every minibatch of input data: compute this forward pass, and then backprop!

$\hat{x}$

Sample x|z from   $x \mid z \sim \mathcal{N}(\mu_{x\mid z}, \Sigma_{x\mid z})$

$\mu_{x\mid z}$          $\Sigma_{x\mid z}$

Decoder network
$p_\theta(x \mid z)$

$z$

Sample z from   $z \mid x \sim \mathcal{N}(\mu_{z\mid x}, \Sigma_{z\mid x})$

$\mu_{z\mid x}$          $\Sigma_{z\mid x}$

Encoder network
$q_\phi(z \mid x)$

**Input Data**   $x$

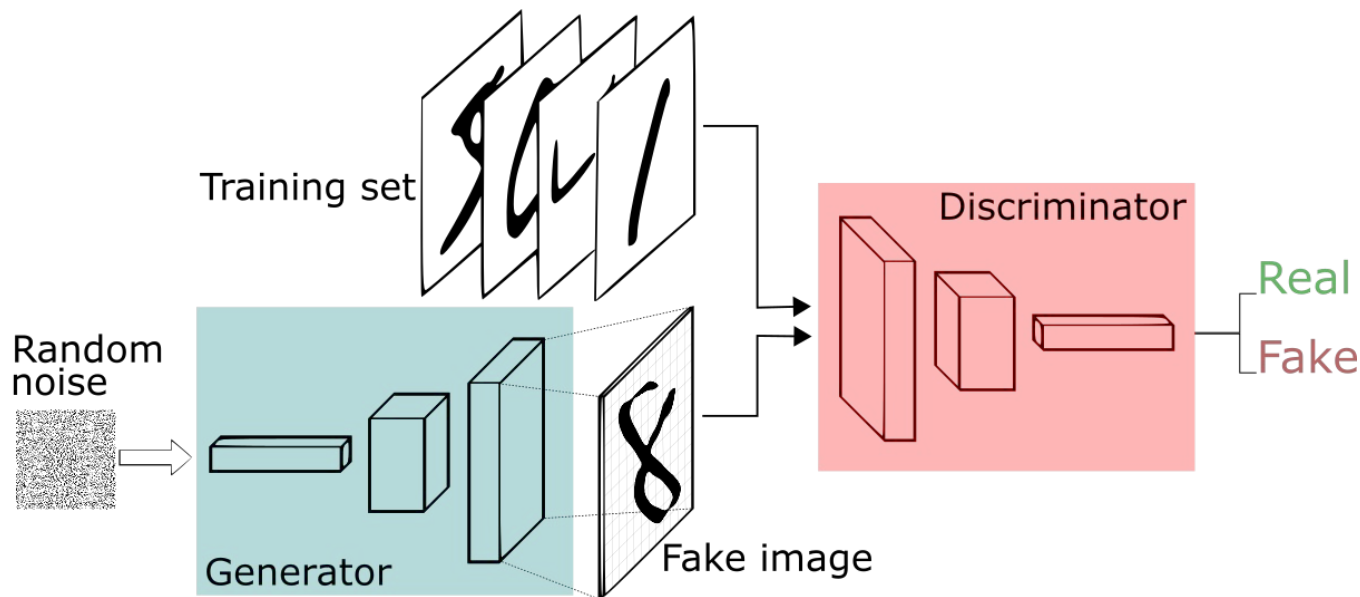Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

# Agenda

- Generative Models
- Variational Autoencoders (VAE)
- Generative Adversarial Networks (GAN)
  - Basic Mechanism: Generator and Discriminator
  - Dual Loss Function
  - Major Issues
    - Mode Collapse
    - Why are GANs so hard to train
- Types of GAN
  - Deep Convolution GAN
  - Cycle GAN
  - Etc.
- Applications of GAN

# Basic Mechanism

- Generator Network: try to fool the discriminator by generating real-looking images
- Discriminator Network: try to distinguish between real and fake images

# Minmax Loss Function

$$\min_{\theta_g} \max_{\theta_d} \left[ \boldsymbol{E}_{x \sim p_{data}} log D_{\theta_d}(x) + \boldsymbol{E}_{z \sim p(z)} log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator output
for real data x

Discriminator output for
generated fake data G(z)

- Discriminator ($\theta_d$) wants to **maximize objective** such that $D(x)$ is close to *1* (real) and $D(G(z))$ is close to *0* (fake)

- Generator ($\theta_g$) wants to **minimize objective** such that $D(G(z))$ is close to 1 (discriminator is fooled into thinking generated $G(z)$ is real)

# Training

Minmax loss function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \boldsymbol{E}_{x \sim p_{data}} log D_{\theta_d}(x) + \boldsymbol{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between

- **Gradient ascent** on discriminator

$$\max_{\theta_d} \left[ \boldsymbol{E}_{x \sim p_{data}} log D_{\theta_d}(x) + \boldsymbol{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

- **Gradient descent** on generator

$$\min_{\theta_g} \boldsymbol{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

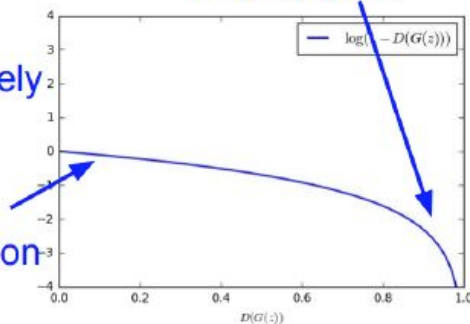However, this doesn't work well in practice

Ian Goodfellow et al. "General Adversarial Nets", NIPS 2014

# Training

Gradient descent on generator:

$$\min_{\theta_g} \boldsymbol{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

We thus choose **gradient ascent on generator**:

$$\max_{\theta_g} \boldsymbol{E}_{z \sim p(x)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$
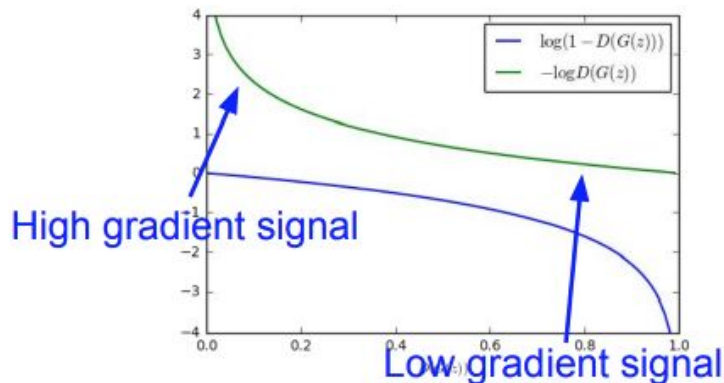


Gradient signal dominated by region where sample is already good

When sample is likely fake, want to learn from it to improve generator. But gradient in this region is relatively flat!

High gradient signal

Low gradient signal

# Training



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

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(x^{(i)}\right) + \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right) \right].$$
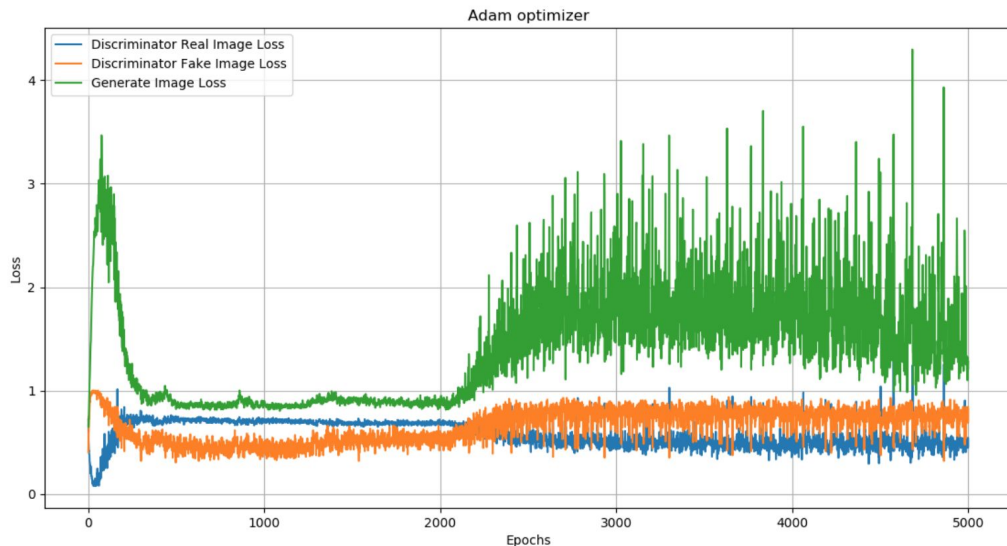
    **end for**
    • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
    • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$
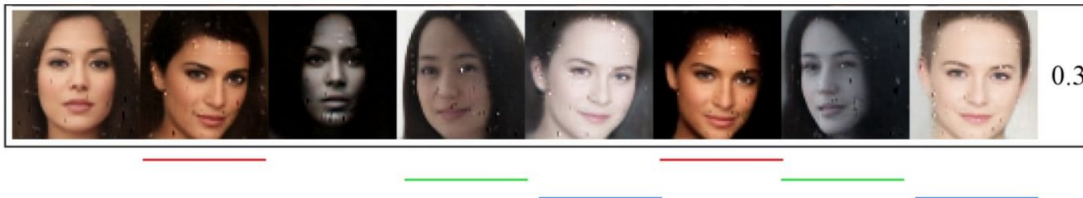
**end for**
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Ian Goodfellow et al. "General Adversarial Nets", NIPS 2014
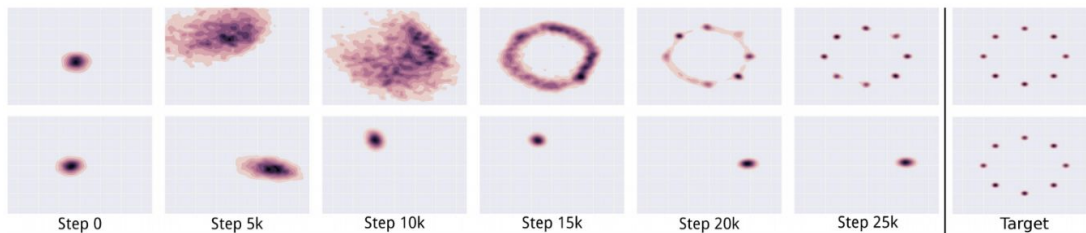
# Major Issues - mode collapse

- Mode collapse happens when the model only generates few modes of data



- Why does it happen? Consider the most extreme case where G is trained so extensively without updates to D that it finds the most optimal generated image that fools D the most, a.k.a the mode collapses to a single point.
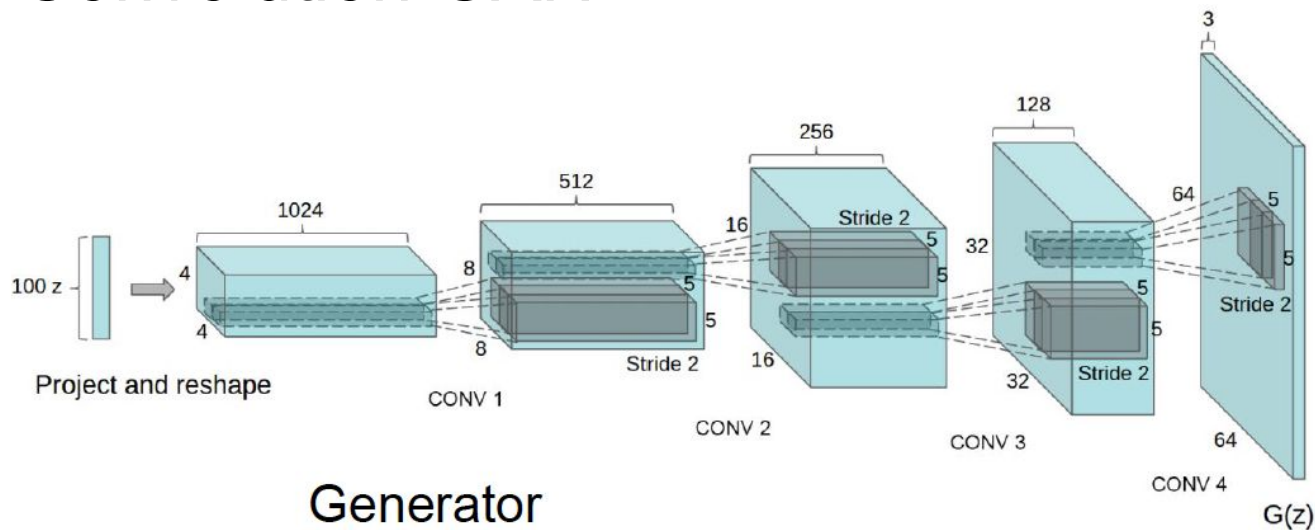
$$x^* = argmax_x D(x)$$

$$\frac{\partial J}{\partial z} \approx 0$$

- When D starts to train again, it can easily detect the single mode, which leads D to overfit to short-term opponent's weaknesses.



Step 0    Step 5k    Step 10k    Step 15k    Step 20k    Step 25k    Target
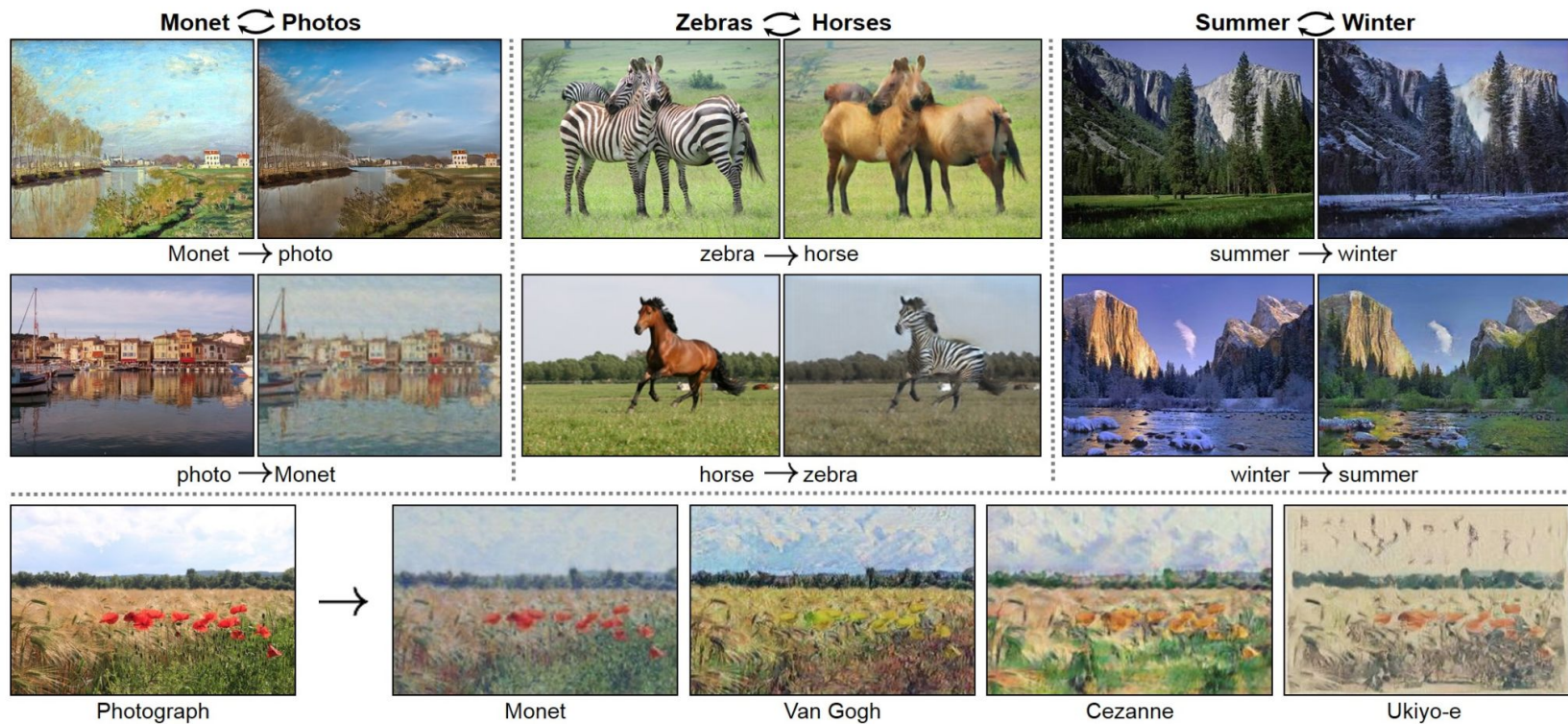
# Agenda

- Generative Models
- Variational Autoencoders (VAE)
- Generative Adversarial Networks (GAN)
  - Basic Mechanism: Generator and Discriminator
  - Dual Loss Function
  - Major Issues
    - Mode Collapse
    - Why are GANs so hard to train
- <span style="color:red">Types of GAN</span>
  - Deep Convolution GAN
  - Cycle GAN
  - etc.
- Applications of GAN

# Deep Convolution GAN



Generator

- Generator is an upsampling network with fractionally-strided convolutions
- Discriminator is a convolutional network

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

# Cycle GAN



Zhu et al. "Unpaired Image-to-image Translation using Cycle-consistent Adversarial Networks", ICCV 2017

# Cycle GAN
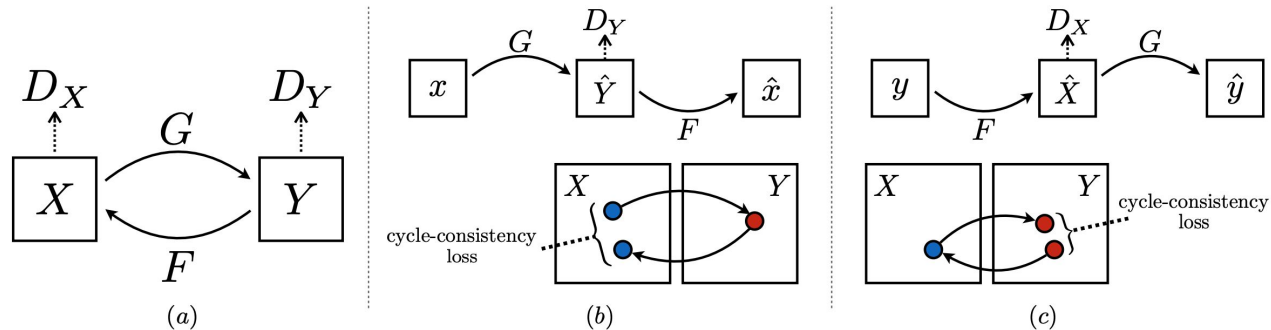


Figure 3: (a) Our model contains two mapping functions $G : X \to Y$ and $F : Y \to X$, and associated adversarial discriminators $D_Y$ and $D_X$. $D_Y$ encourages $G$ to translate $X$ into outputs indistinguishable from domain $Y$, and vice versa for $D_X$ and $F$. To further regularize the mappings, we introduce two *cycle consistency losses* that capture the intuition that if we translate from one domain to the other and back again we should arrive at where we started: (b) forward cycle-consistency loss: $x \to G(x) \to F(G(x)) \approx x$, and (c) backward cycle-consistency loss: $y \to F(y) \to G(F(y)) \approx y$

$$\mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)}[\log D_Y(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log(1 - D_Y(G(x)))]$$

**Adversarial Loss**

$$\mathcal{L}_{\text{cyc}}(G, F) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_{\text{data}}(y)}[\|G(F(y)) - y\|_1].$$

**Cycle Consistency Loss**

Zhu et al. "Unpaired Image-to-image Translation using Cycle-consistent Adversarial Networks", ICCV 2017
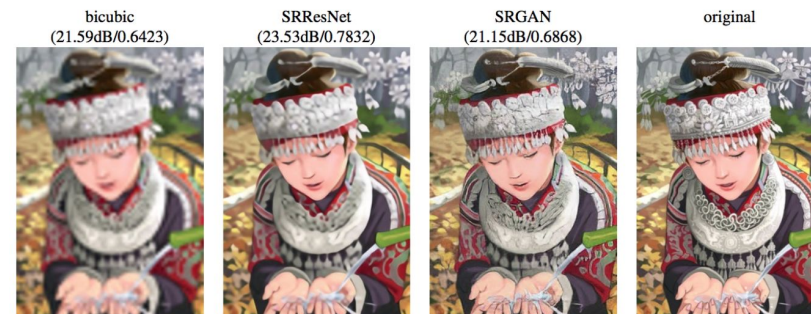
# Agenda

- Generative Models
- Variational Autoencoders (VAE)
- Generative Adversarial Networks (GAN)
    - Adversarial Examples
    - Basic Mechanism: Generator and Discriminator
    - Dual Loss Function
    - Major Issues
        - Mode Collapse
        - Why are GANs so hard to train
- Types of GAN
    - Deep Convolution GAN
    - Cycle GAN
    - etc.
- Applications of GAN

# Application

## Improved resolution (SRGAN)

Ledig et al, "Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network"





## Text to Image (StackGAN)

Zhang et al, "StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks"
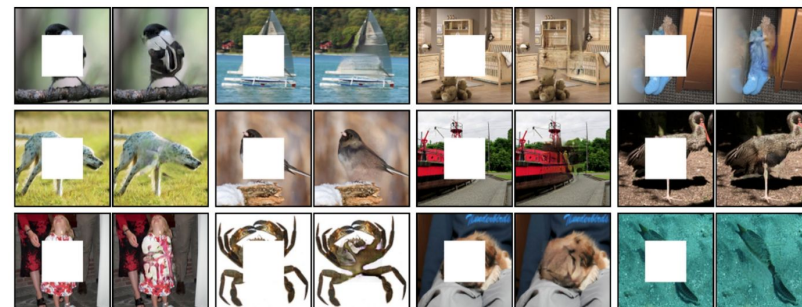
## Image Inpainting

Pathak et al, "Context encoding: feature learning by inpainting", CVPR 2016