# Generative Adversarial Networks (GANs)
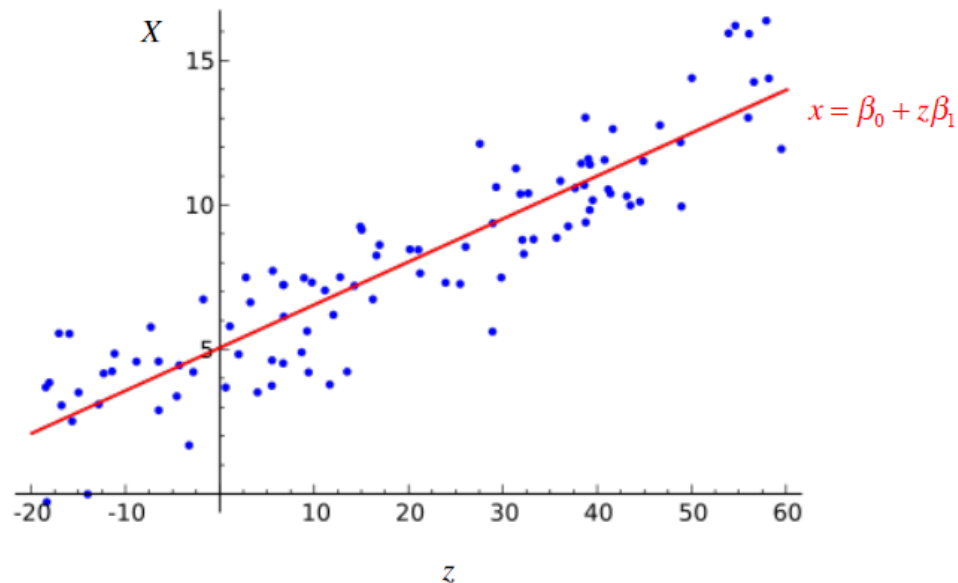
Prof. Seungchul Lee

POSTECH

# Source

- 1시간만에 GAN(Generative Adversarial Network) 완전 정복하기
  - by 최윤제 (고려대 석사생)
  - YouTube: https://www.youtube.com/watch?v=odpjk7_tGY0
  - Slides: https://www.slideshare.net/NaverEngineering/1-gangenerative-adversarial-network

- CSC321 Lecture 19: GAN
  - By Prof. Roger Grosse at Univ. of Toronto
  - http://www.cs.toronto.edu/~rgrosse/courses/csc321_2018/

- CS231n: CNN for Visual Recognition
  - Lecture 13: Generative Models
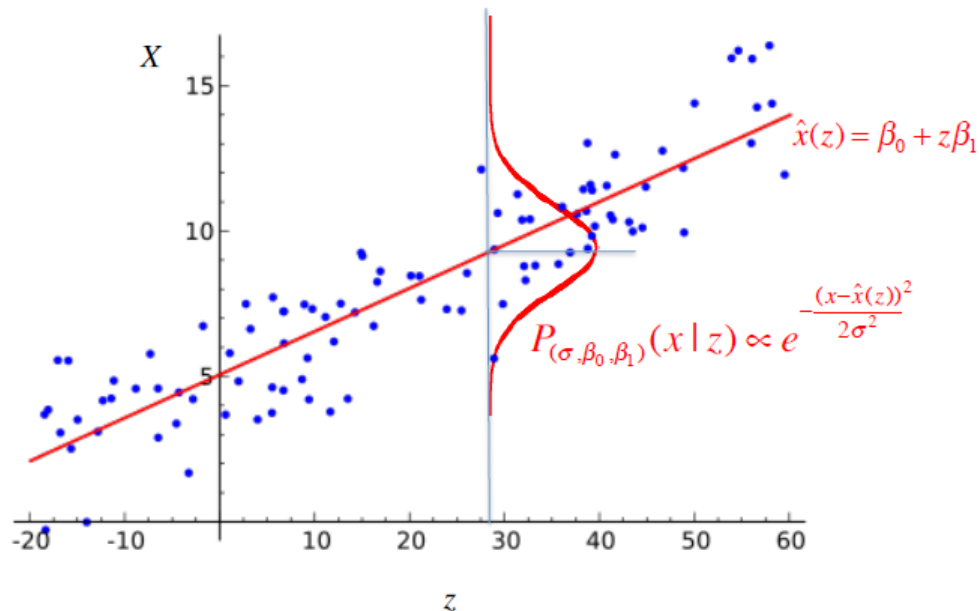  - By Prof. Fei-Fei Li at Stanford University
  - http://cs231n.stanford.edu/

# Recap: Linear Regression

- Most people think of linear regression as points and a straight line:
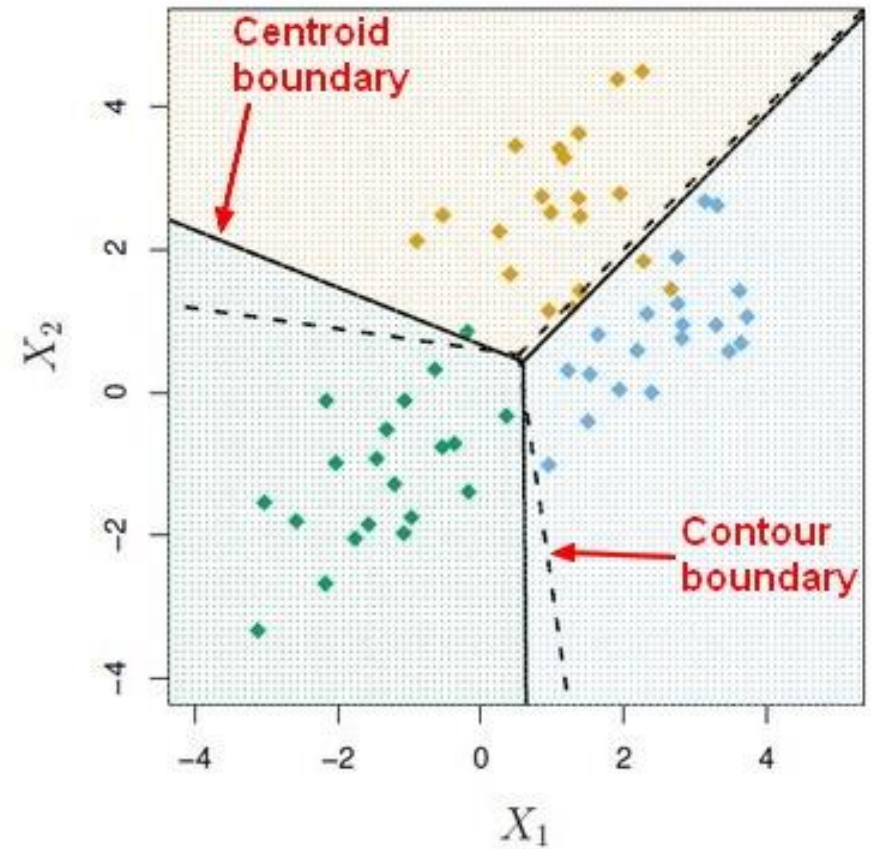


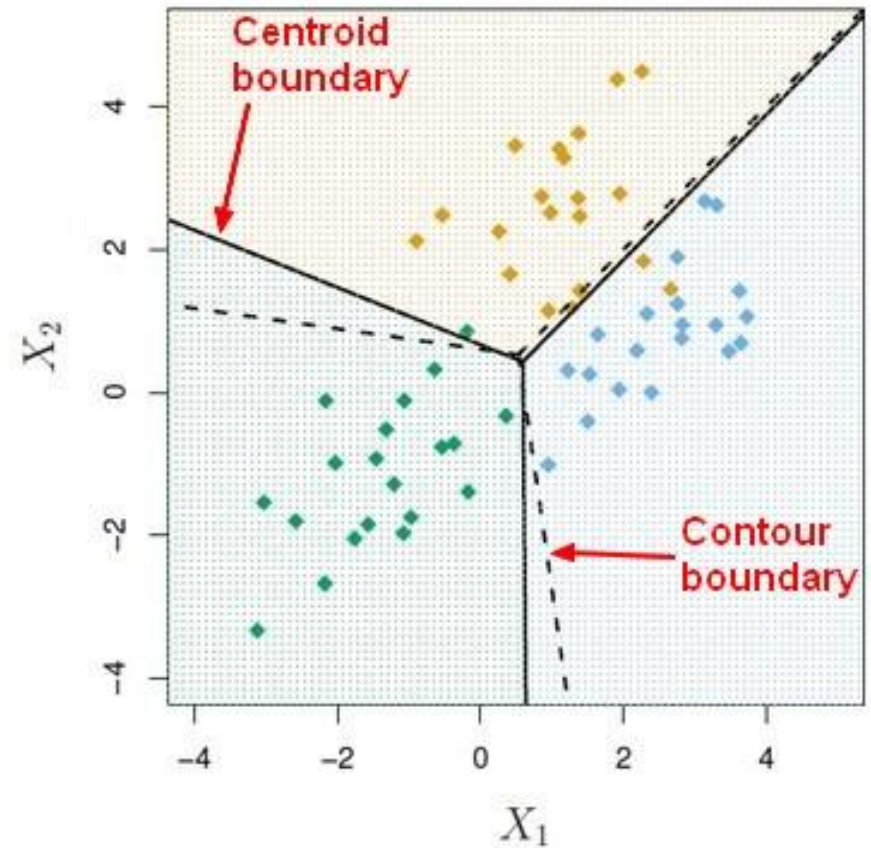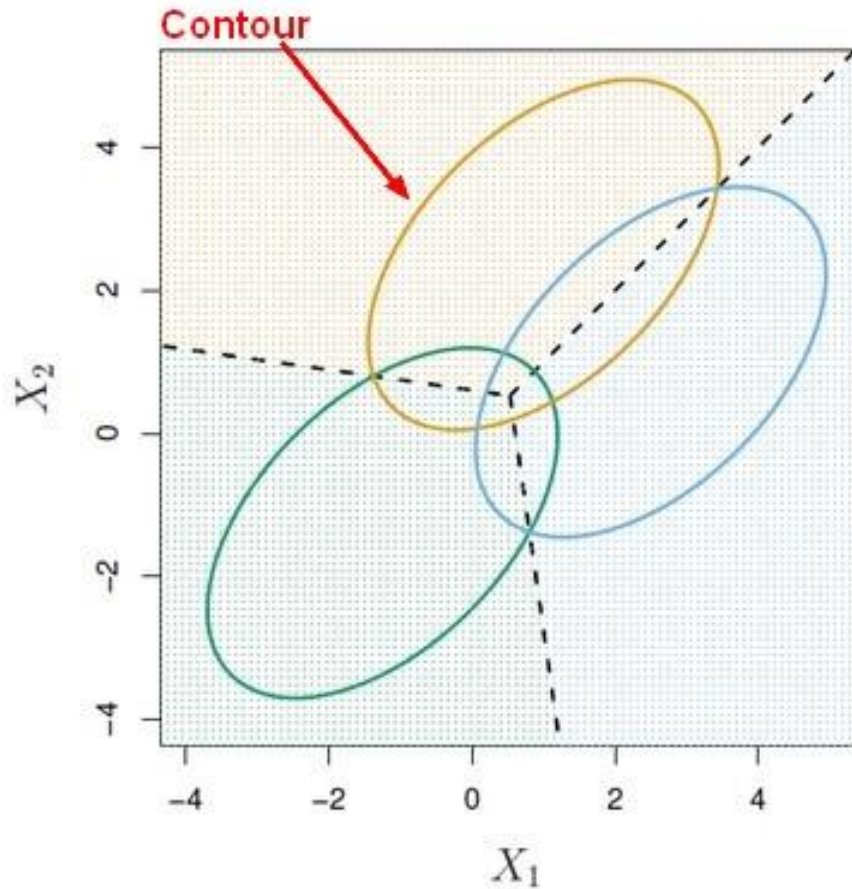$$x = \beta_0 + z\beta_1$$

# Recap: Linear Regression

- Statisticians additionally have $P_\theta(X|Z)$
- Benefits of having an error model:
  - How likely is a data point
  - Confidence bounds
  - Compare models



$$\hat{x}(z) = \beta_0 + z\beta_1$$

$$P_{(\sigma, \beta_0, \beta_1)}(x|z) \propto e^{-\frac{(x-\hat{x}(z))^2}{2\sigma^2}}$$
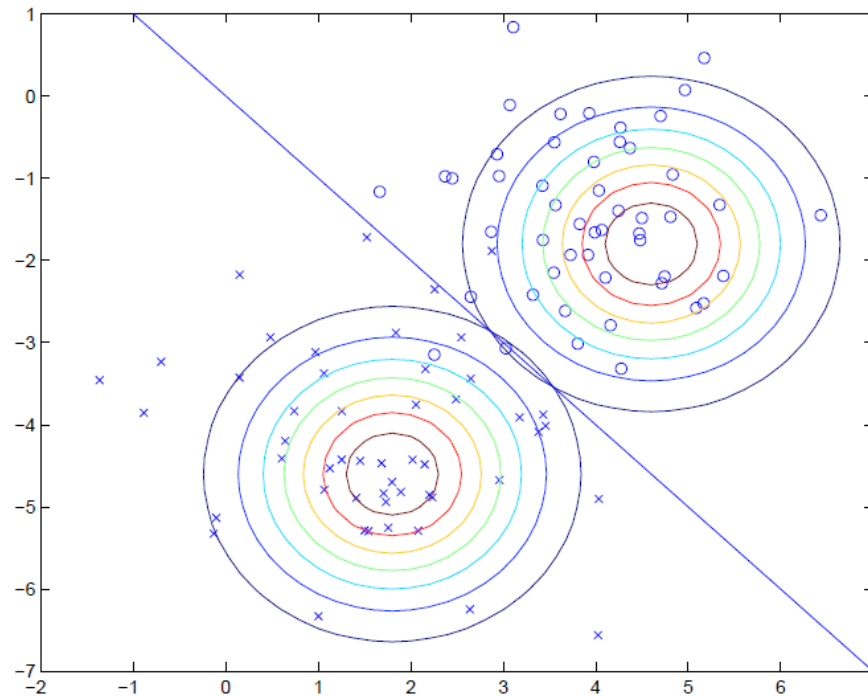
# Recap: Linear Classifier

# Recap: Linear Classifier
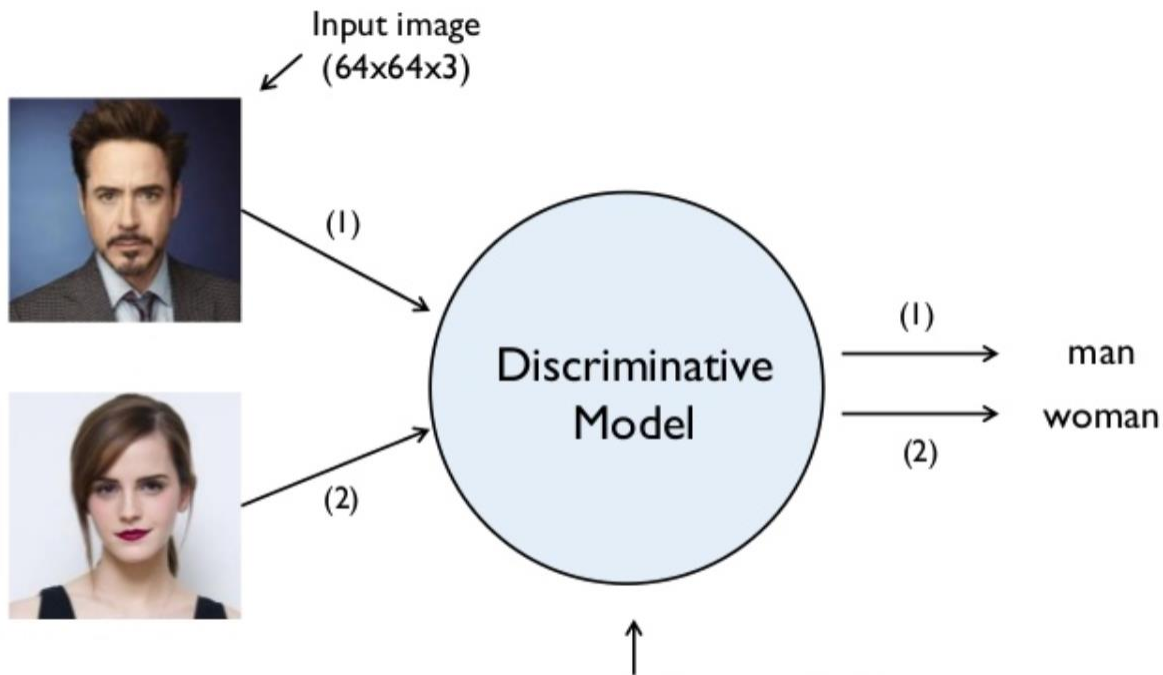
# Recap: Linear Classifier

- Think about how data is generated

$$
\begin{aligned}
y &\sim \text{Bernoulli}(\phi) \\
x|y=0 &\sim \mathcal{N}(\mu_0, \Sigma) \\
x|y=1 &\sim \mathcal{N}(\mu_1, \Sigma)
\end{aligned}
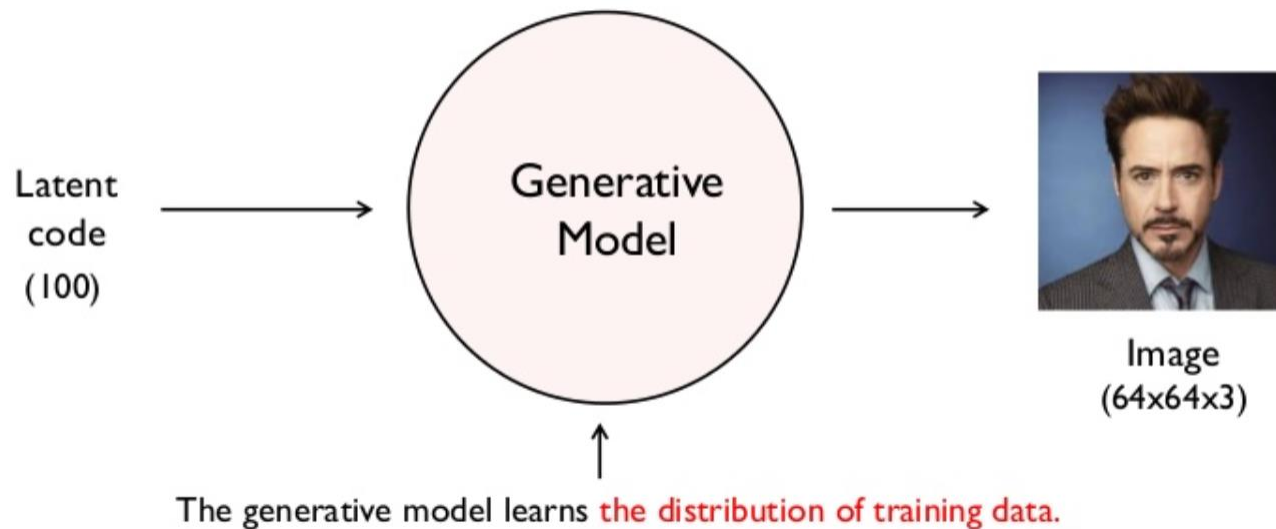$$

# Supervised Learning

- Discriminative model



Input image (64x64x3)

Discriminative Model

(1) → man
(2) → woman

The discriminative model learns how to classify input to its class.

# Unsupervised Learning

- Generative model



Latent code (100) → Generative Model → Image (64x64x3)
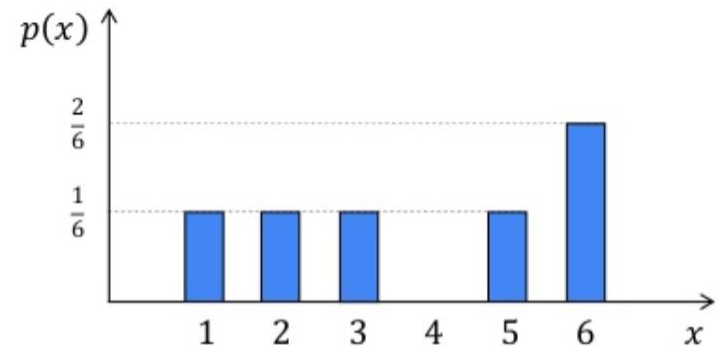
The generative model learns the distribution of training data.

# Probability Distribution

Probability Basics (Review)



Random variable

| $X$ | 1 | 2 | 3 | 4 | 5 | 6 |
|------|------|------|------|------|------|------|
| $P(X)$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{0}{6}$ | $\frac{1}{6}$ | $\frac{2}{6}$ |

Probability mass function

# Probability Distribution

What if $x$ is actual images in the training data?

At this point, $x$ can be represented as a (for example) 64x64x3 dimensional vector.
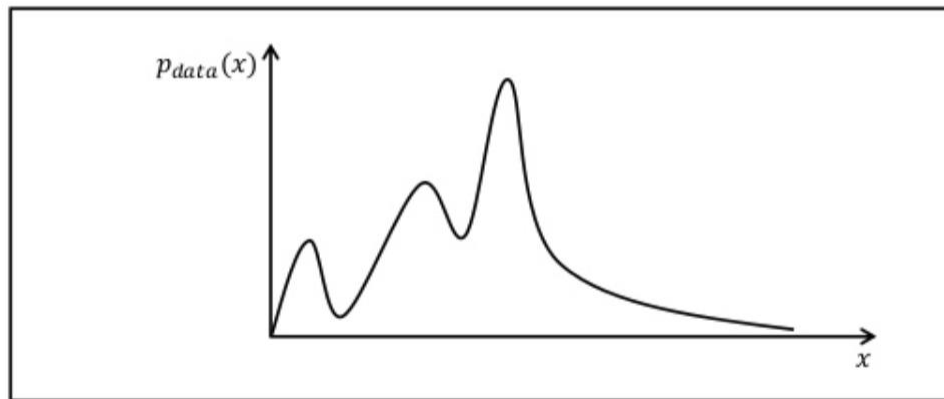
# Probability Distribution

Probability density function

There is a $p_{data}(x)$ that represents the distribution of actual images.
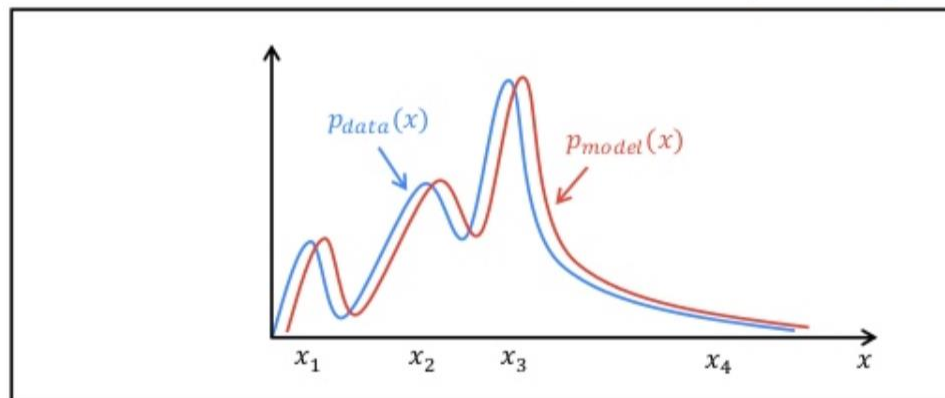
# Density Estimation

- Probability density estimation problem

The goal of the generative model is to find a $p_{model}(x)$ that
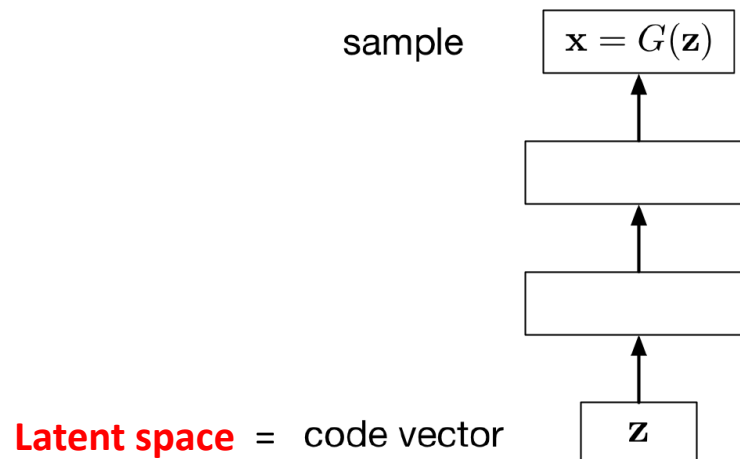approximates $p_{data}(x)$ well.

→ Distribution of images generated by the model

↳ Distribution of actual images

$p_{data}(x)$

$p_{model}(x)$

$x_1$ $x_2$ $x_3$ $x_4$ $x$

- If $P_{model}(x)$ can be estimated as close to $P_{data}(x)$, then data can be generated by sampling from $P_{model}(x)$
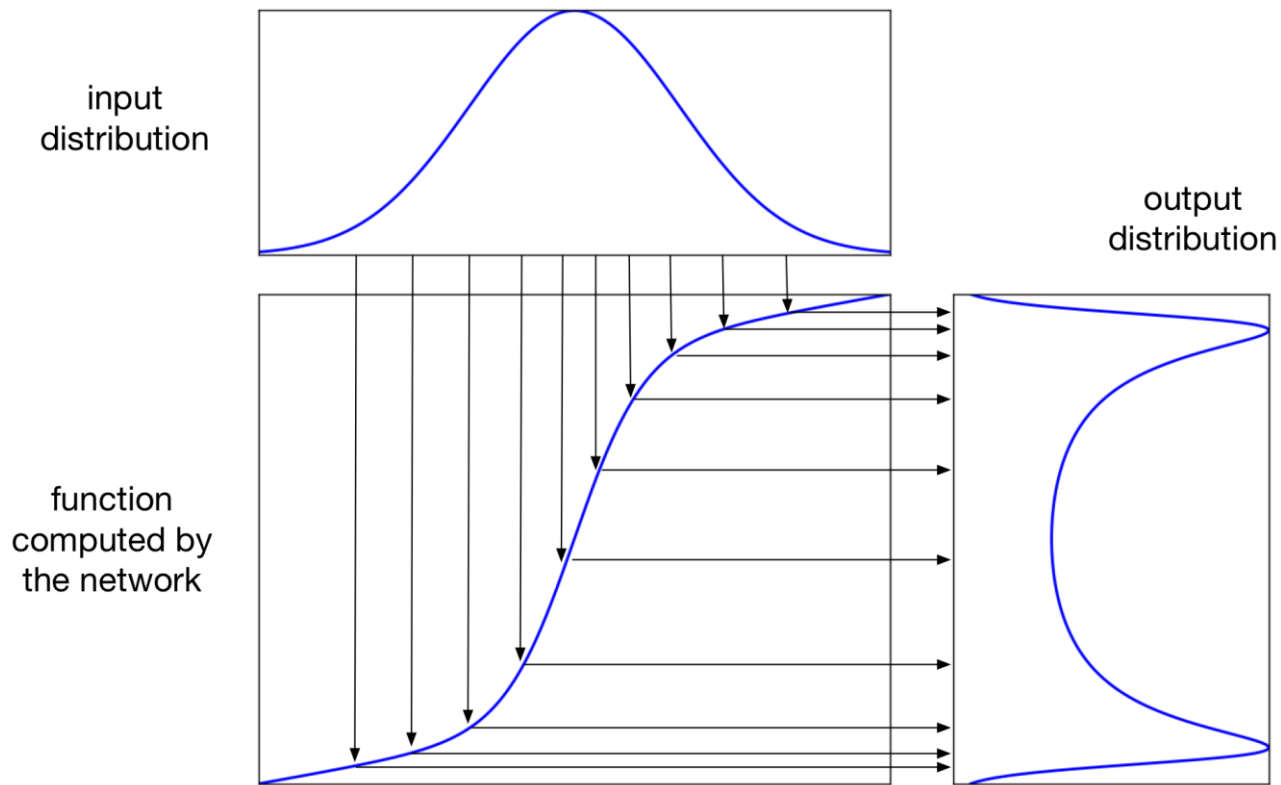
# Generative Models from Lower Dim.

- Learn transformation via a neural network

- Start by sampling the code vector $z$ from a fixed, simple distribution (e.g. uniform distribution or Gaussian distribution)
- Then this code vector is passed as input to a deterministic generator network $G$, which produces an output sample $x = G(z)$
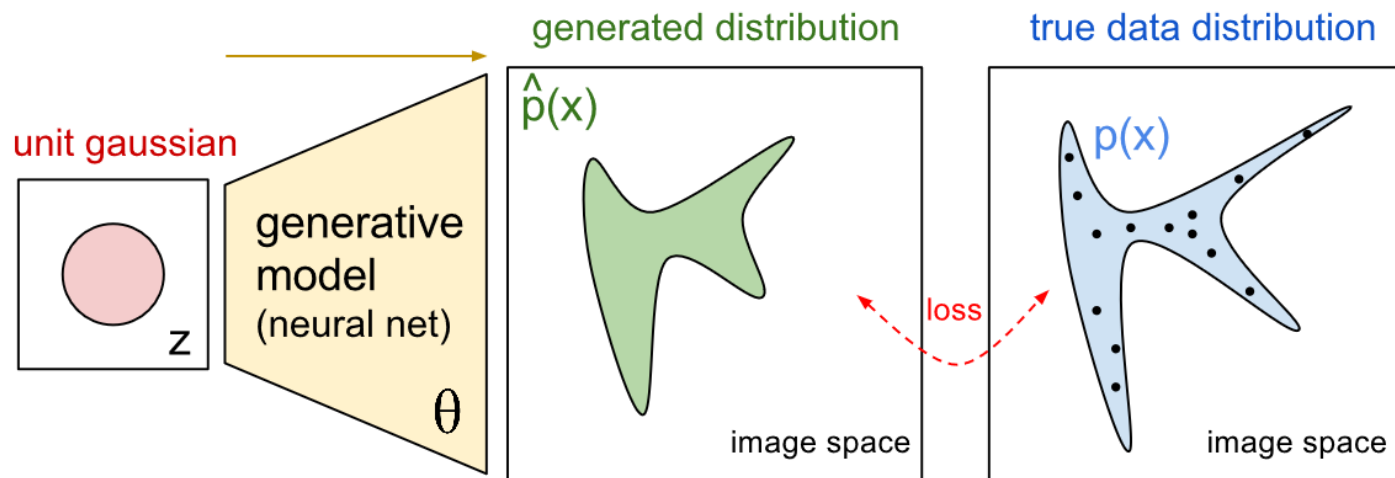
sample     $\mathbf{x} = G(\mathbf{z})$

**Latent space** = code vector    $\mathbf{z}$

# Deterministic Transformation (by Network)

- 1-dimensional example:

input
distribution

output
distribution

function
computed by
the network

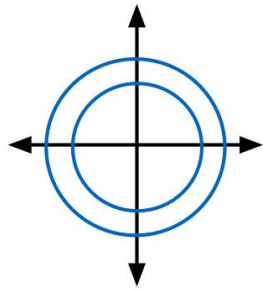# Deterministic Transformation (by Network)
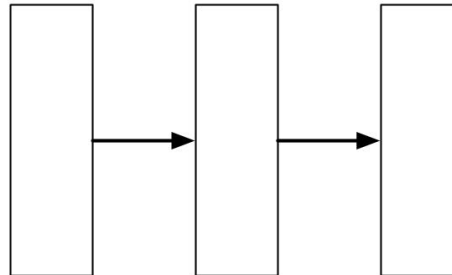
- High dimensional example:
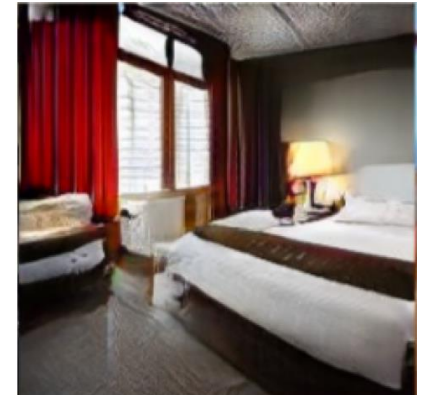
# Prob. Density Function by Deep Learning

- Generative model of image



Each dimension of the code vector is sampled independently from a simple distribution, e.g. Gaussian or uniform.
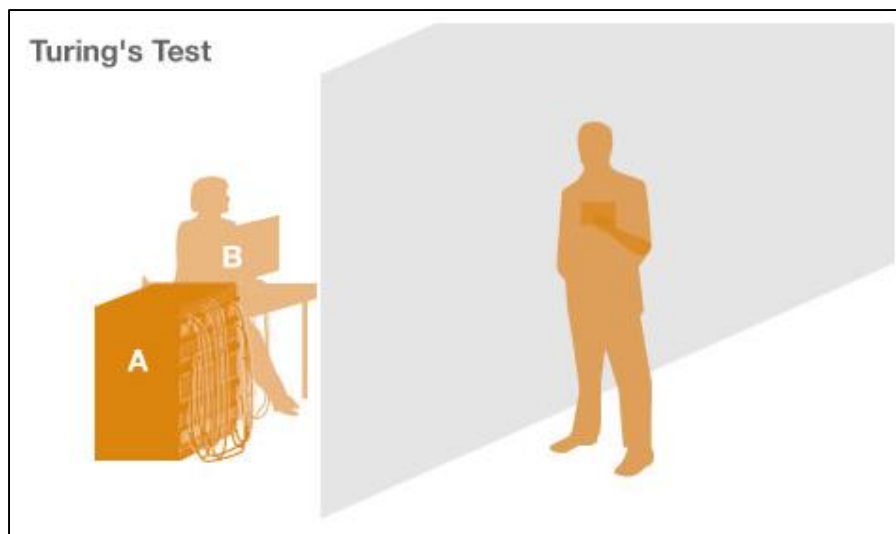
This is fed to a (deterministic) generator network.
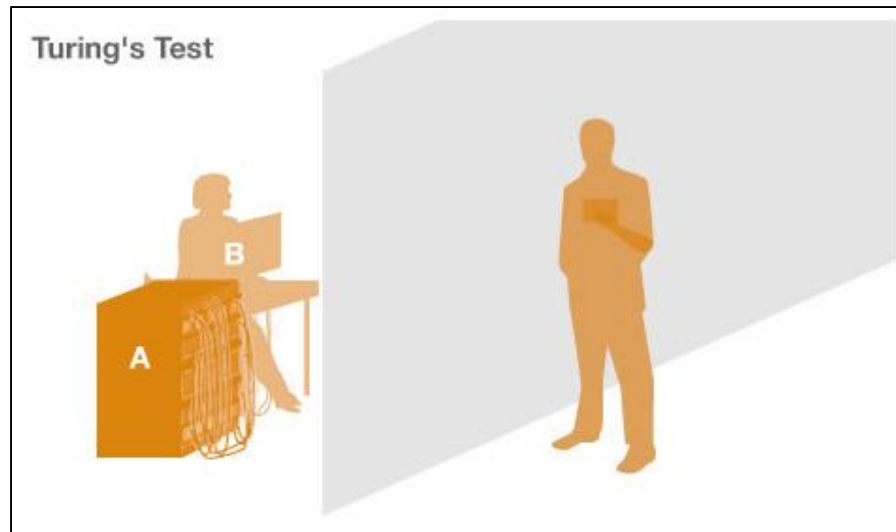
The network outputs an image.

# Generative Adversarial Networks (GANs)

- In generative modeling, we'd like to train a network that models a distribution, such as a distribution over images.

- GANs do not work with any explicit density function !
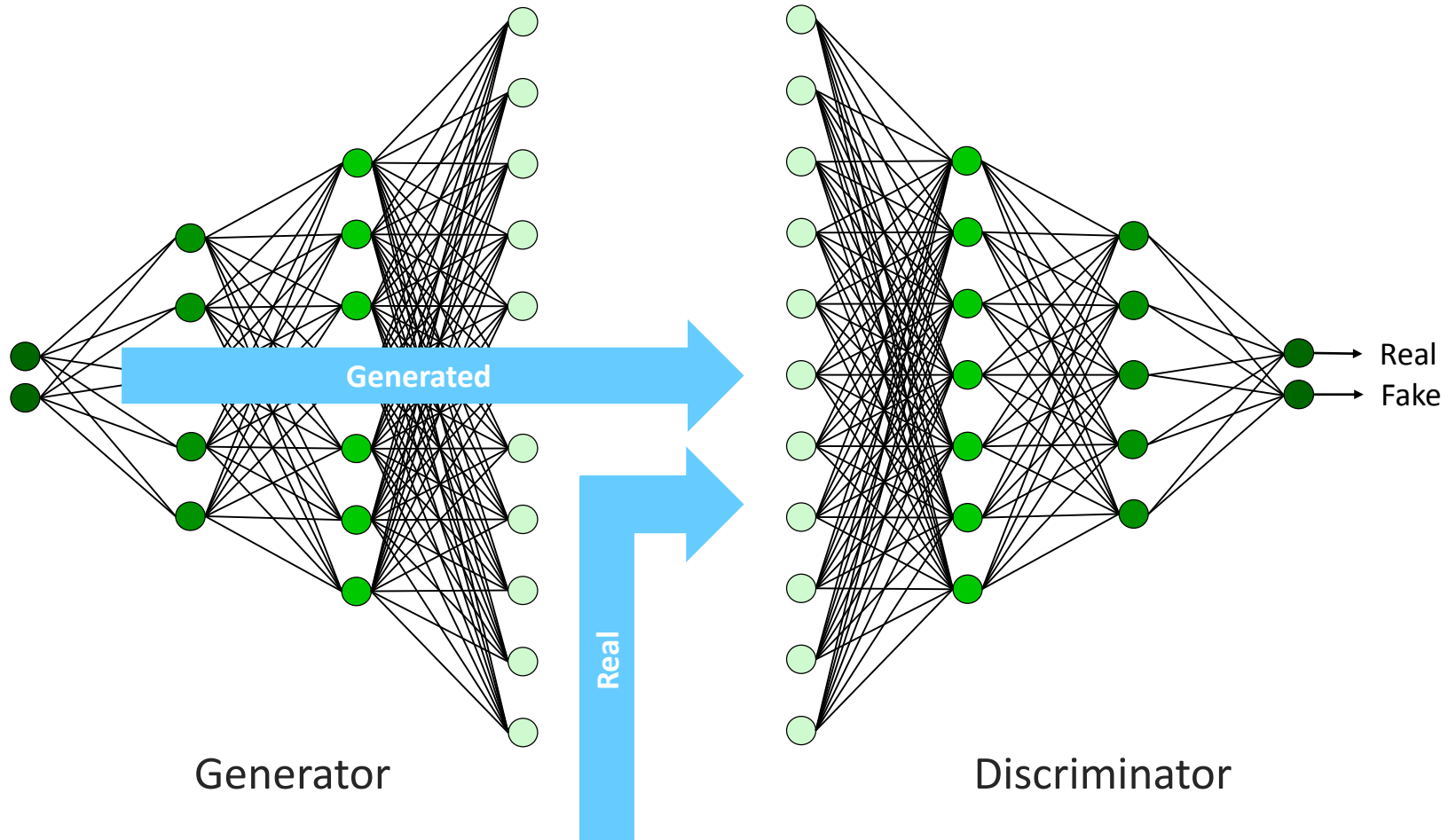  - Instead, take game-theoretic approach

# Turing Test

- One way to judge the quality of the model is to sample from it.
- GANs are based on a very different idea:
  - Model to produce samples which are indistinguishable from the real data, as judged by a discriminator network whose job is to tell real from fake



Turing's Test

# Generative Adversarial Networks (GAN)

**Analogous to Turing Test**



Generator

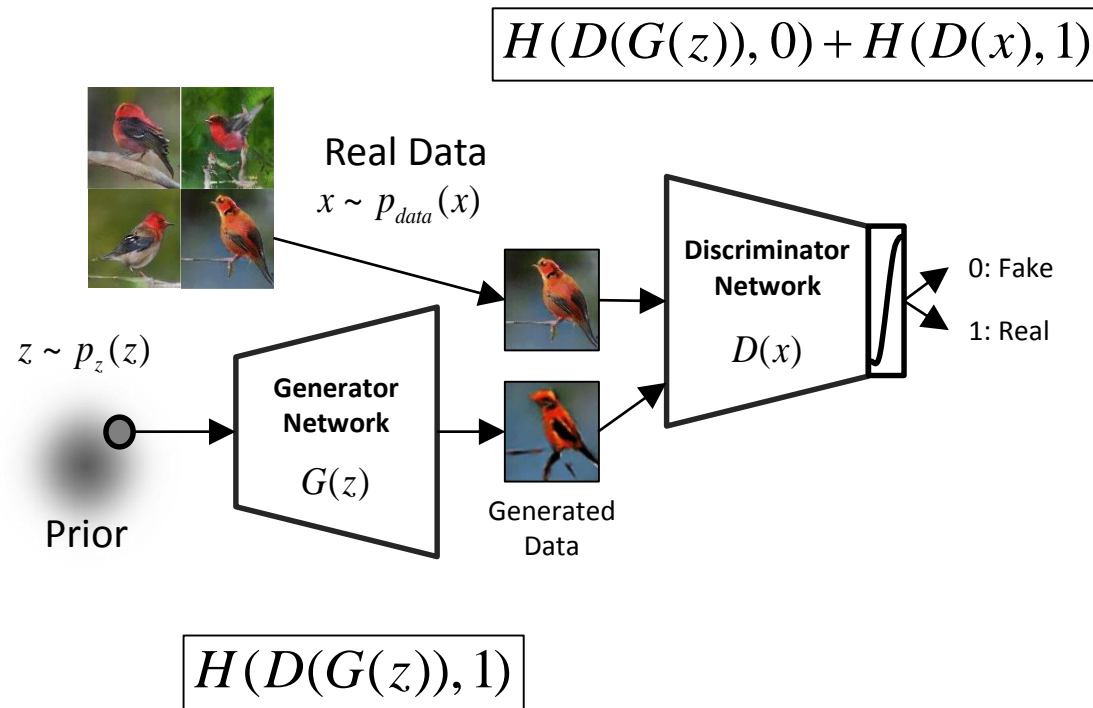Discriminator

Generated

Real

Real
Fake

# Generative Adversarial Networks (GAN)

- The idea behind Generative Adversarial Networks (GANs): train two different networks
  - Generator network: try to produce realistic-looking samples
  - Discriminator network: try to distinguish between real and fake data


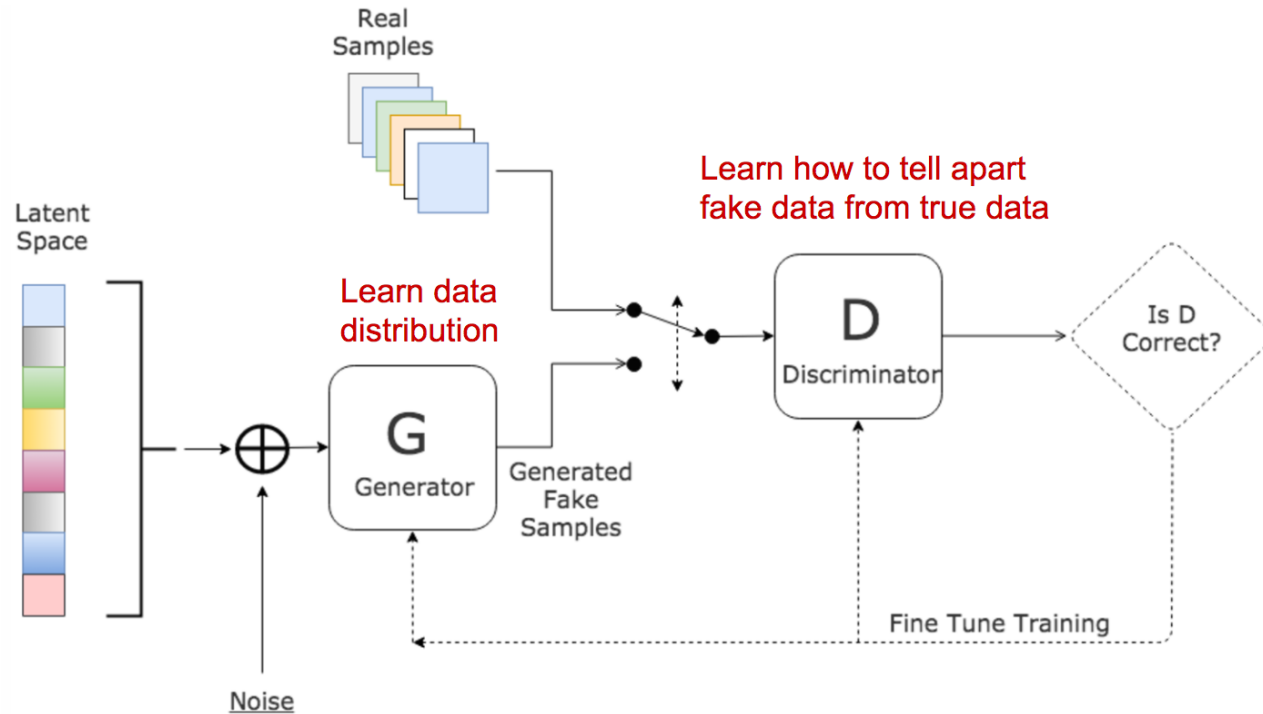- The generator network tries to fool the discriminator network

# Generative Adversarial Networks (GAN)

- How to generate data?
  - Train through competition
  - Generator vs. Discriminator

$$H(D(G(z)), 0) + H(D(x), 1)$$

Real Data
$x \sim p_{data}(x)$

$z \sim p_z(z)$

**Generator Network**

$G(z)$

Prior

Generated Data

**Discriminator Network**

$D(x)$

0: Fake

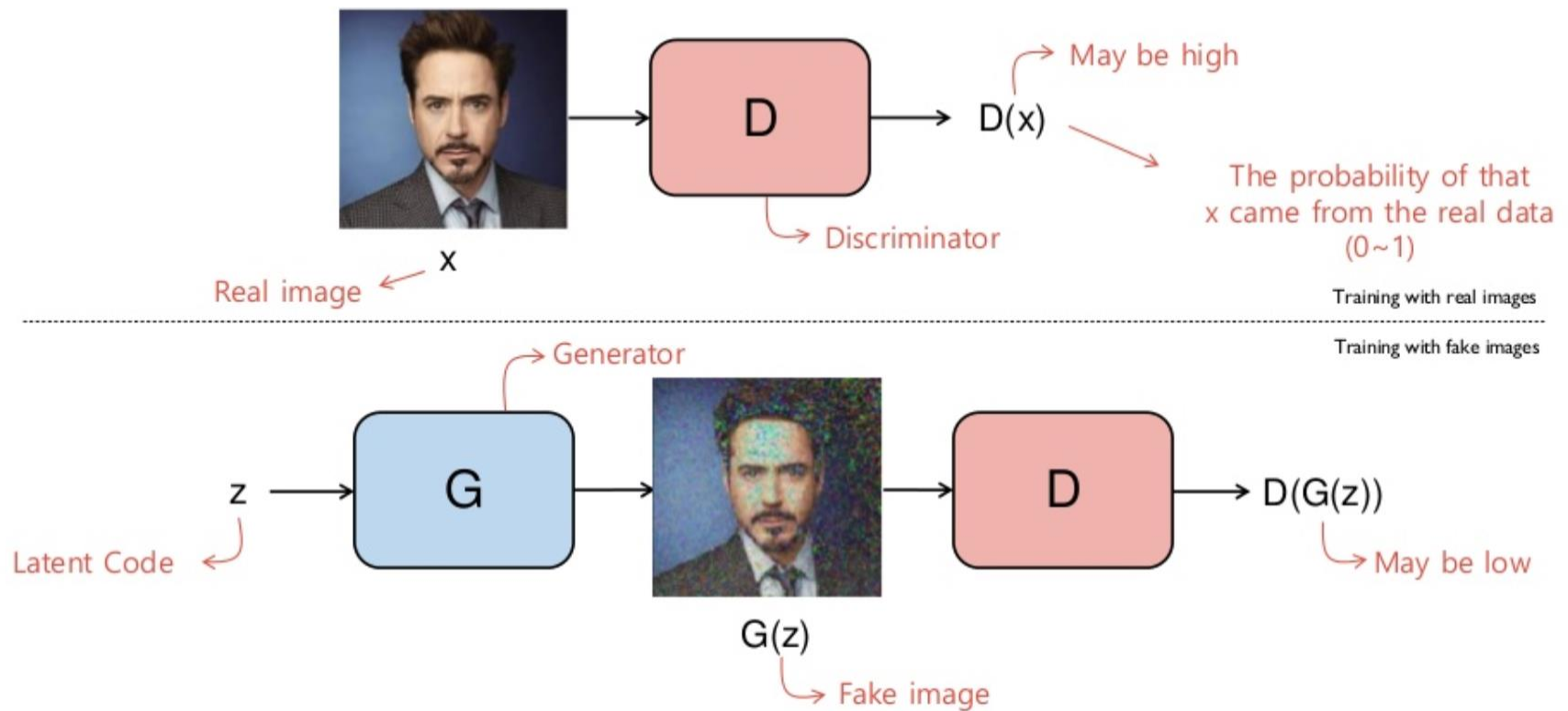1: Real

$$H(D(G(z)), 1)$$

# Generative Adversarial Networks (GAN)

- ## How to generate data?
  - Train through competition
  - Generator vs. Discriminator

# Intuition for GAN
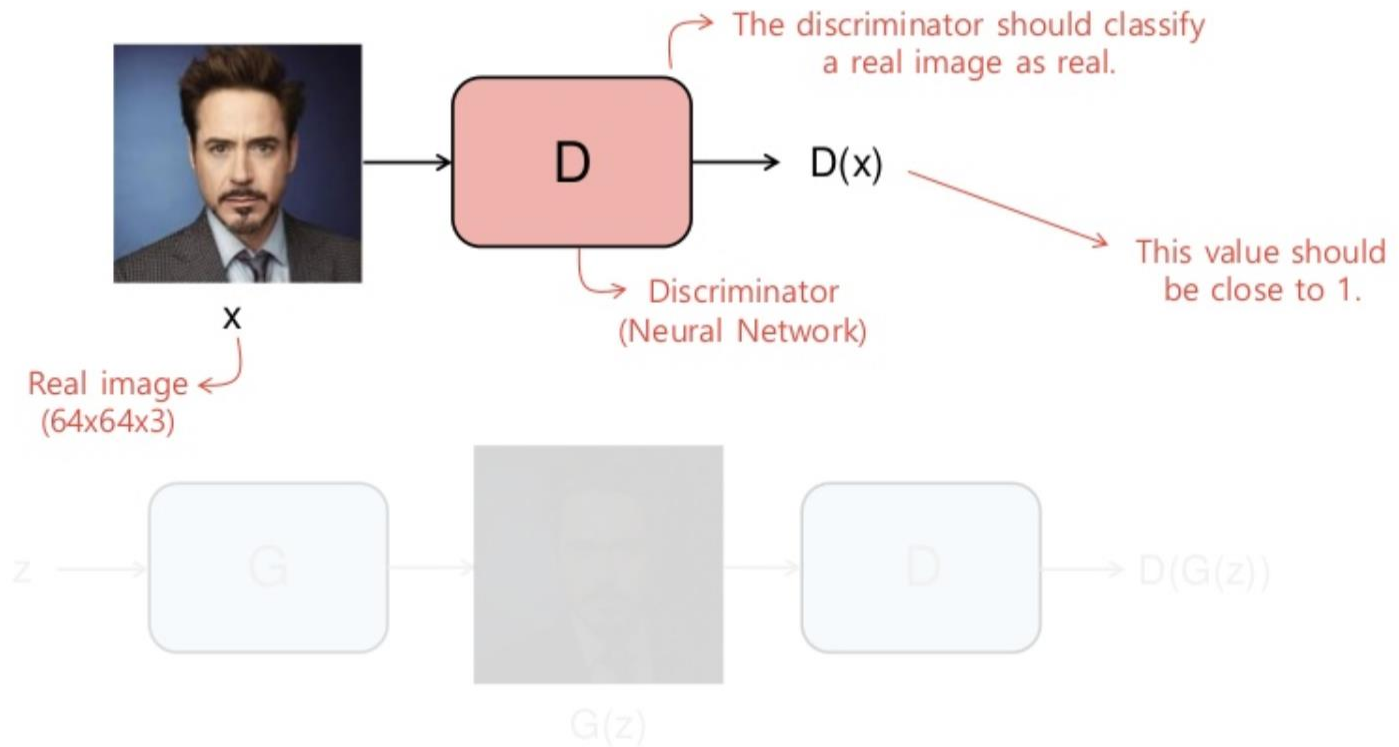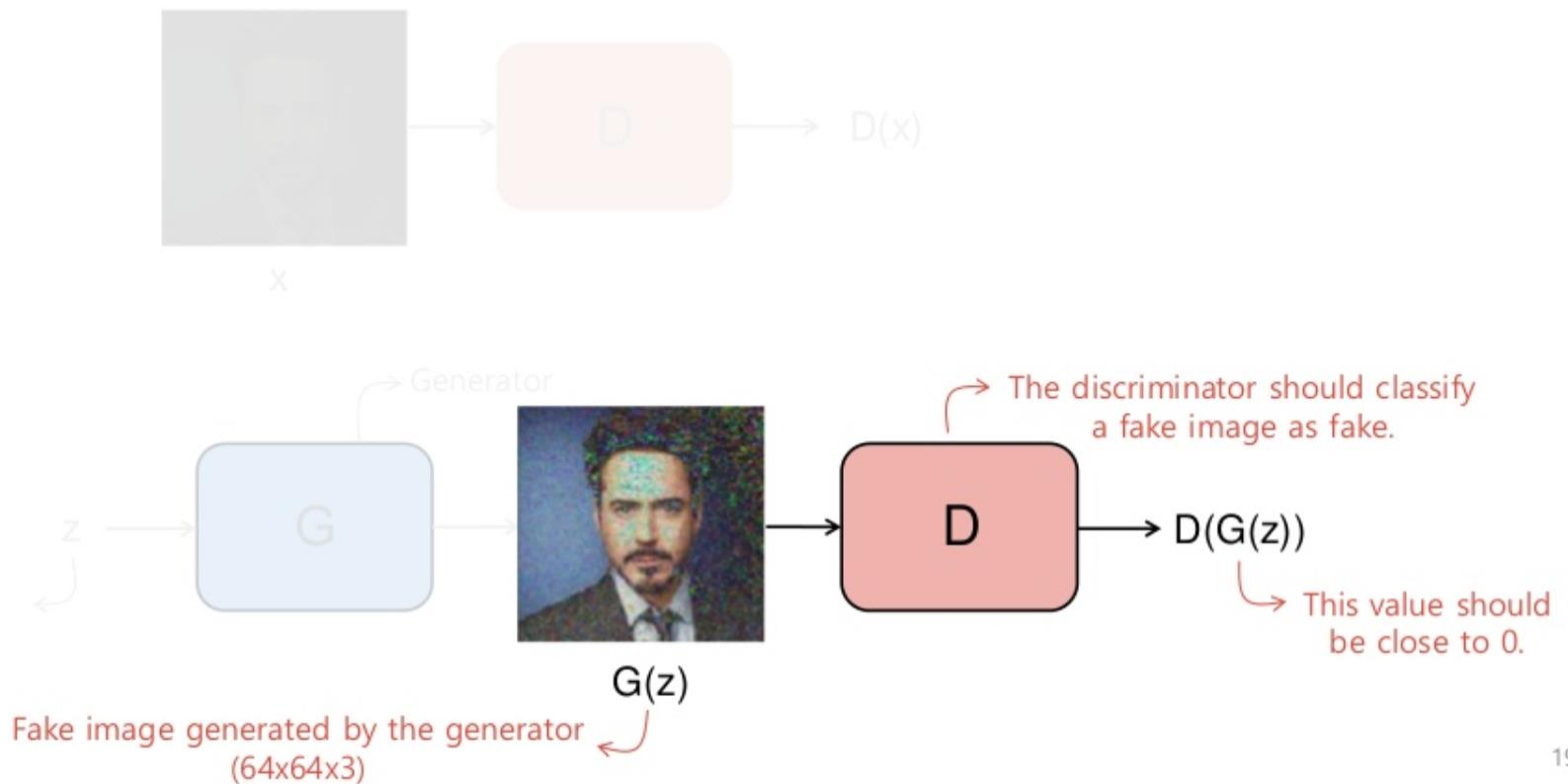
# Discriminator Perspective



The discriminator should classify a real image as real.

D(x)

This value should be close to 1.

Discriminator (Neural Network)

x

Real image (64x64x3)

z → G → G(z) → D → D(G(z))

18

# Discriminator Perspective



G(z)

Fake image generated by the generator
(64x64x3)

The discriminator should classify
a fake image as fake.

D(G(z))

This value should
be close to 0.

19

# Generator Perspective



Generator
(Neural Network)

z → G → D → D(G(z))

Latent Code
(100)

The generator should create an image
that is indistinguishable from real to
deceive the discriminator

G(z)

Generated image
(64x64x3)

This value should
be close to 1.

20

# Objective Function of GAN

$$\text{loss} = -y \log h(x) - (1 - y) \log(1 - h(x))$$



Objective function

Training with real images

Training with fake images

# Objective Function of GAN



$$\min_{G} \max_{D} V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

$G$ is independent of this part

$G$ should minimize $V(D, G)$

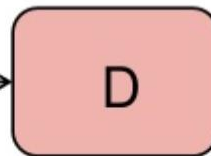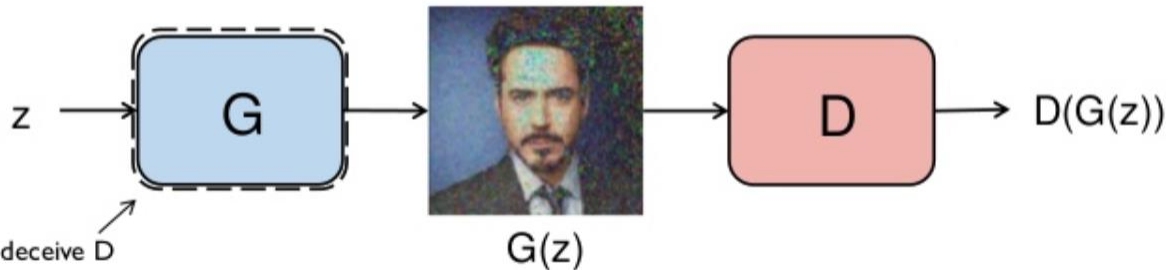Minimum when $D(G(z)) = 1$

Objective function

D → D(x)

x

Training with real images

Training with fake images

z → G → D → D(G(z))

Train G to deceive D

G(z)

22

# Non-Saturating Game

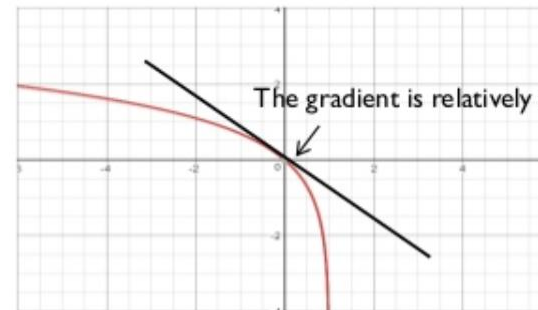$$\min_G E_{z \sim p_z(z)}[\log(1 - D(G(z))]$$

Objective function of G

At the beginning of training, the discriminator can clearly classify the generated image as fake because the quality of the image is very low.

This means that D(G(z)) is almost zero at early stages of training.

The gradient is relatively small at D(G(z))=0

$$y = \log(1 - x)$$

Images created by the generator at the beginning of training

32

# Non-Saturating Game

$$\min_{G} E_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

Modification (heuristically motivated)

$$\max_{G} E_{z \sim p_z(z)}[\log D(G(z))]$$

```
1    # tensorflow
2    tf.losses.sigmoid_cross_entropy()
3
4    # pytorch
5    nn.BCELoss()
```

- Practical Usage

Use binary cross entropy loss function with fake label (1)

$$\min_{G} E_{z \sim p_z(z)}[-y \log D(G(z)) - (1 - y) \log(1 - D(G(z)))]$$

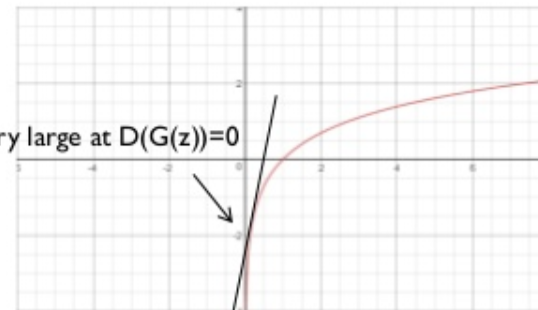$$\downarrow \quad y = 1$$

$$\min_{G} E_{z \sim p_z(z)}[-\log D(G(z))]$$

The gradient is very large at D(G(z))=0



$$y = \log(x)$$

33

# Solving a Minmax Problem

Step 1: Fix $G$ and perform a gradient step to

$$\max_D E_{x \sim p_{\text{data}}(x)} \left[ \log D(x) \right] + E_{x \sim p_z(z)} \left[ \log(1 - D(G(z))) \right]$$

Step 2: Fix $D$ and perform a gradient step to

$$\max_G E_{x \sim p_z(z)} \left[ \log D(G(z)) \right]$$

OR

Step 1: Fix $G$ and perform a gradient step to
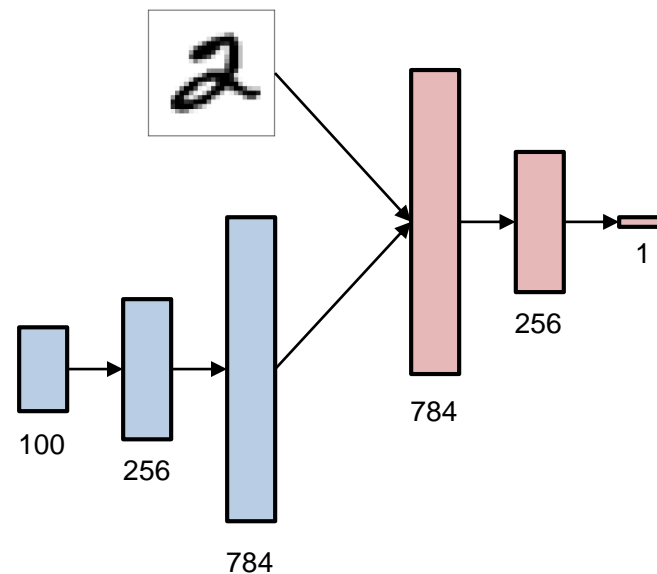
$$\min_D E_{x \sim p_{\text{data}}(x)} \left[ -\log D(x) \right] + E_{x \sim p_z(z)} \left[ -\log(1 - D(G(z))) \right]$$

Step 2: Fix $D$ and perform a gradient step to
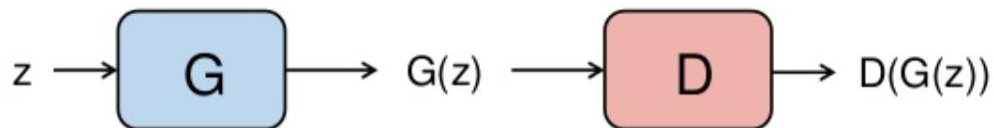
$$\min_G E_{x \sim p_z(z)} \left[ -\log D(G(z)) \right]$$

**POSTECH**

# TensorFlow Implementation



Training with real images

Training with fake images

# TensorFlow Implementation

Assume x is MNIST (784 dimension)

Define the discriminator

input size: 784
hidden size: 128
output size: 1

Output probability(1 dimension)

Discriminator

```
n_D_input = 28*28
n_D_hidden = 256
n_D_output = 1

n_G_input = 100
n_G_hidden = 256
n_G_output = 28*28
```

# TensorFlow Implementation



Define the generator

input size: 100
hidden size: 128
output size: 784

```
n_D_input = 28*28
n_D_hidden = 256
n_D_output = 1

n_G_input = 100
n_G_hidden = 256
n_G_output = 28*28
```

x → D → D(x)

Generator

z → G → G(z) → D → D(G(z))

Latent code (100 dimension)    Generated image (784 dimension)

100    256    784    784    256    1

# TensorFlow Implementation

Step 1: Fix $G$ and perform a gradient step to

$$\min_{D} E_{x \sim p_{\text{data}}(x)} \left[ -\log D(x) \right] + E_{x \sim p_z(z)} \left[ -\log(1 - D(G(z))) \right]$$

Step 2: Fix $D$ and perform a gradient step to
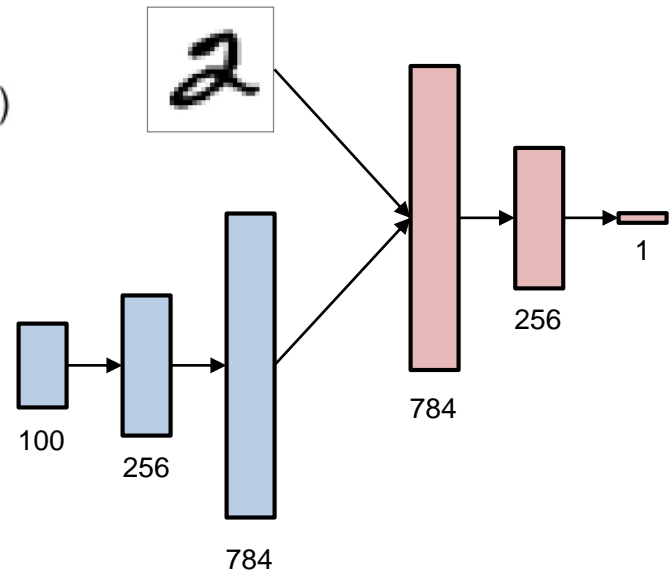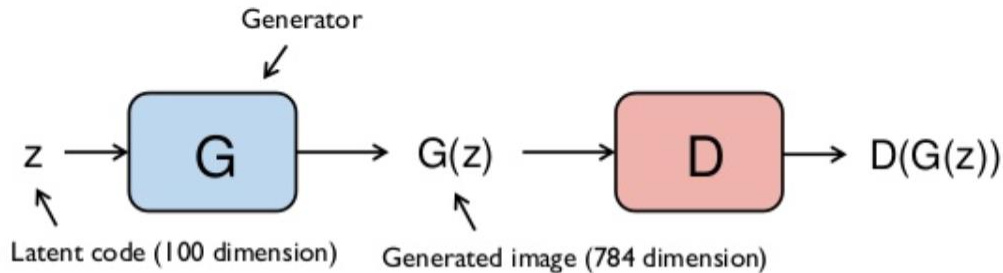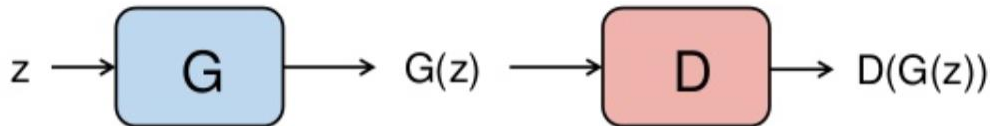
$$\min_{G} E_{x \sim p_z(z)} \left[ -\log D(G(z)) \right]$$

$x \longrightarrow$ D $\longrightarrow D(x)$

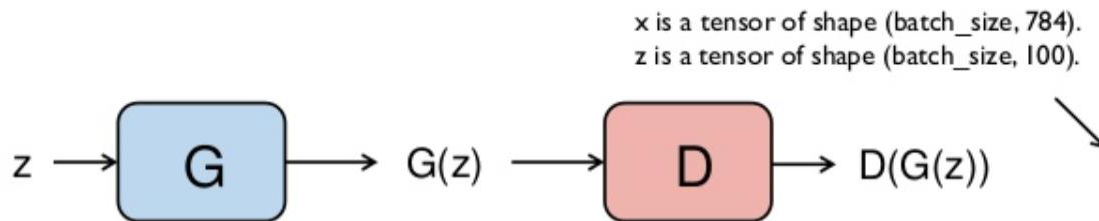$z \longrightarrow$ G $\longrightarrow G(z) \longrightarrow$ D $\longrightarrow D(G(z))$

```
cost_D = tf.reduce_mean(tf.log(D_real) + tf.log(1 - D_fake))
cost_G = tf.reduce_mean(tf.log(D_fake))
```

```
LR = 0.0002
D_train = tf.train.AdamOptimizer(LR).minimize(-cost_D, var_list = D_var_list)
G_train = tf.train.AdamOptimizer(LR).minimize(-cost_G, var_list = G_var_list)
```
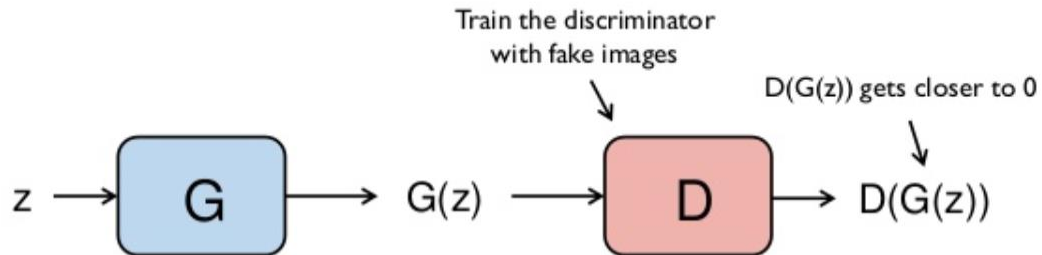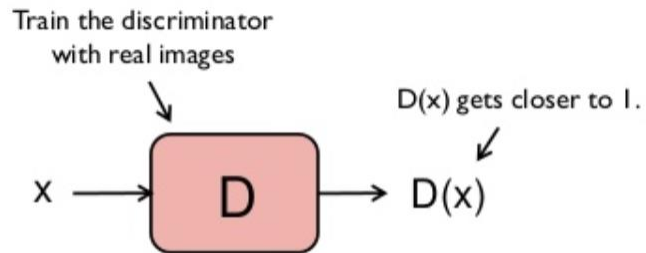
*POSTECH*

# TensorFlow Implementation



x is a tensor of shape (batch_size, 784).
z is a tensor of shape (batch_size, 100).

```
z = tf.placeholder(tf.float32, [None, n_G_input])
x = tf.placeholder(tf.float32, [None, n_D_input])
```

# TensorFlow Implementation

```
_, D_loss_val = sess.run([D_train, cost_D], feed_dict={x: train_x, z: noise})
```



Train the discriminator
with real images

$D(x)$ gets closer to 1.

$x \rightarrow$ D $\rightarrow D(x)$

Train the discriminator
with fake images

$D(G(z))$ gets closer to 0

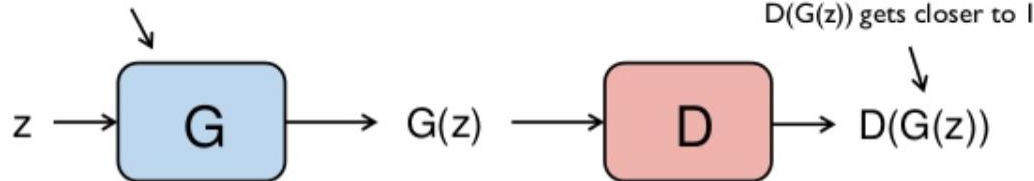$z \rightarrow$ G $\rightarrow G(z) \rightarrow$ D $\rightarrow D(G(z))$

# TensorFlow Implementation

```
_, D_loss_val = sess.run([D_train, cost_D], feed_dict={x: train_x, z: noise})
_, G_loss_val = sess.run([G_train, cost_G], feed_dict={z: noise})
```

# After Training

- After training, use generator network to generate new data



```
noise = make_noise(n_batch, n_G_input)
G_img = sess.run(G_output, feed_dict={z: noise})
plt.imshow(G_img[0,:].reshape(28,28),'gray')
plt.show()
```

**POSTECH**

# GAN Samples



2009



2015



2018





CelebA-HQ
1024 × 1024

Latent space interpolations