# Practical Machine Learning
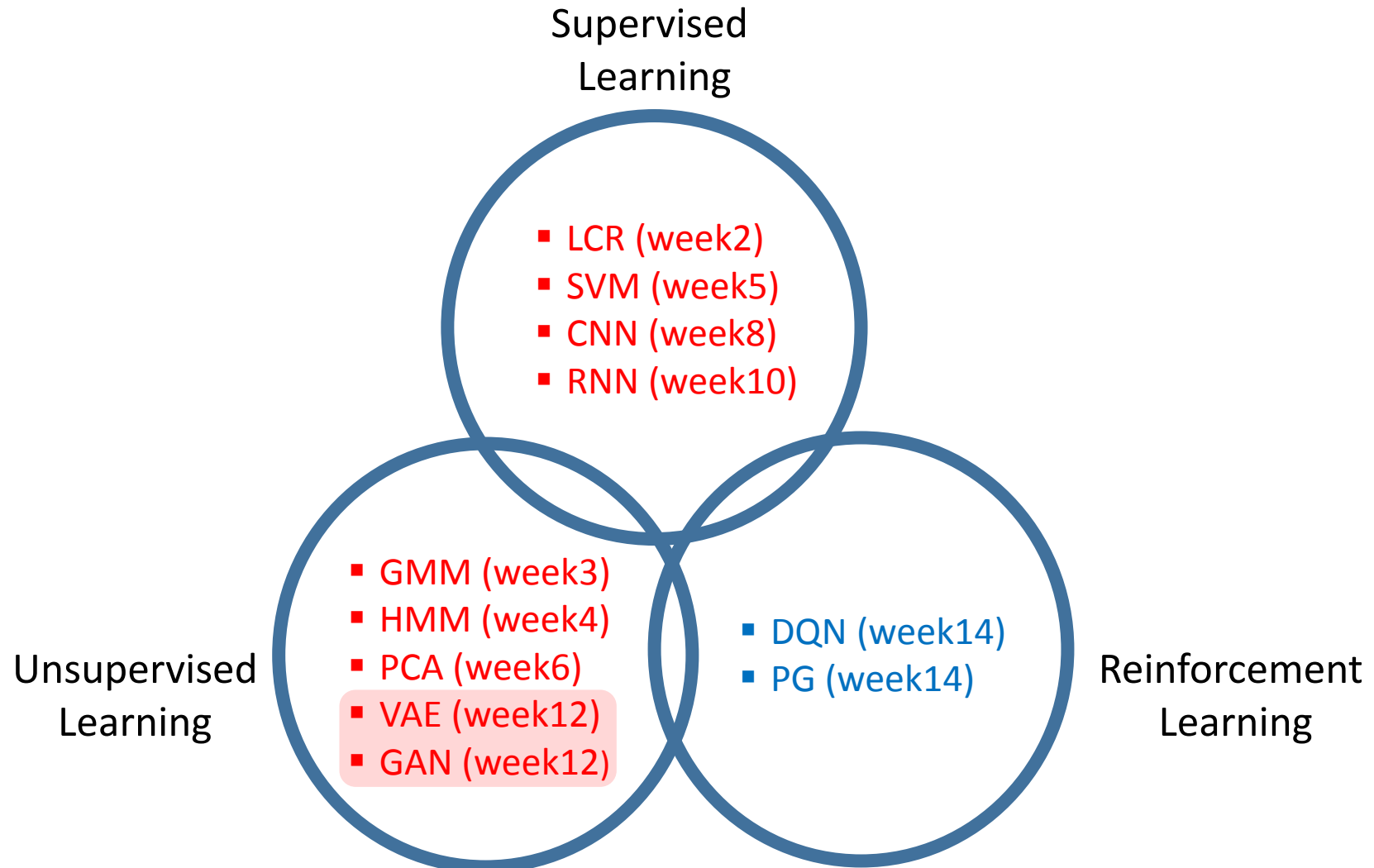
**Lecture 12**

**Generative Models: Variational Auto Encoder (VAE) and Generative Adversarial Networks (GAN)**

Dr. Suyong Eum

OSAKA UNIVERSITY

Supervised Learning

- LCR (week2)
- SVM (week5)
- CNN (week8)
- RNN (week10)

Unsupervised Learning

- GMM (week3)
- HMM (week4)
- PCA (week6)
- VAE (week12)
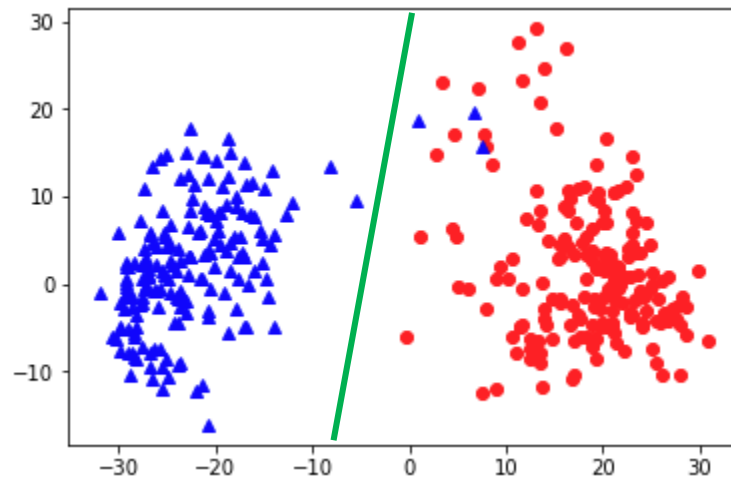- GAN (week12)

Reinforcement Learning
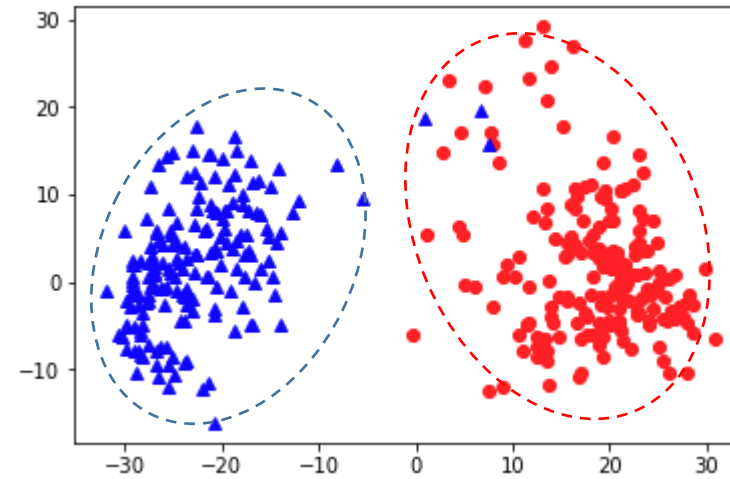
- DQN (week14)
- PG (week14)

# You are going to learn

❑ The basic concept of generative models

❑ Two generative models:
1) Variational Auto Encoder (VAE)
2) Generative Adversarial Networks (GAN)

❑ Some applications you may be interested

❑ There are two types of models:
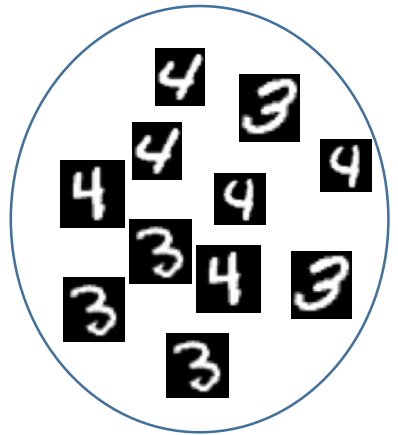1) Discriminative models
2) Generative models



Determinative model, e.g., SVM

Generative model, e.g, GMM

4

❑ Literally speaking, a sample can be generated from generative models.
- Of course, the model needs to be trained in advance to generate such a sample which you are interested.

e.g., GMM



Each data point
has 64 dimension

Project the data points in
2 dimension
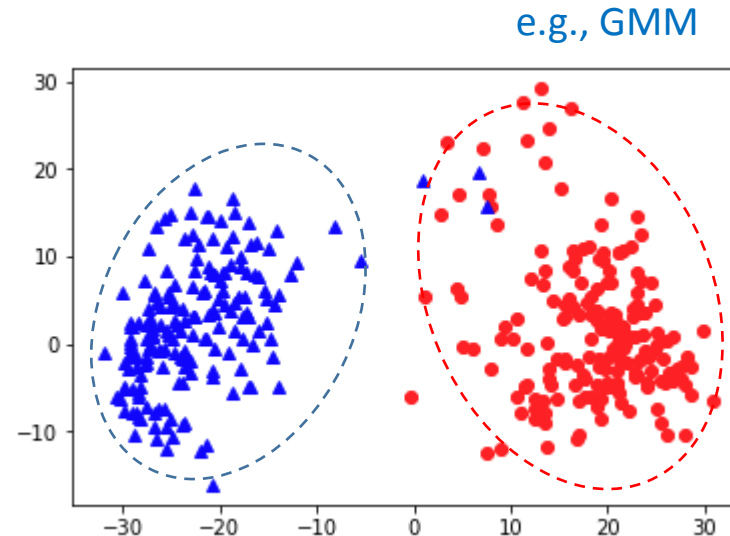
❑ Literally speaking, a sample can be generated from generative models.
- Of course, the model needs to be trained in advance to generate such a sample which you are interested.
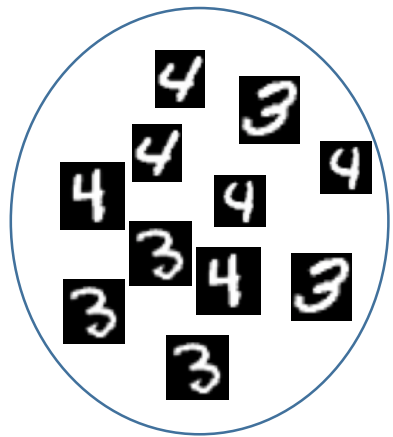
e.g., GMM



a sample can be generated
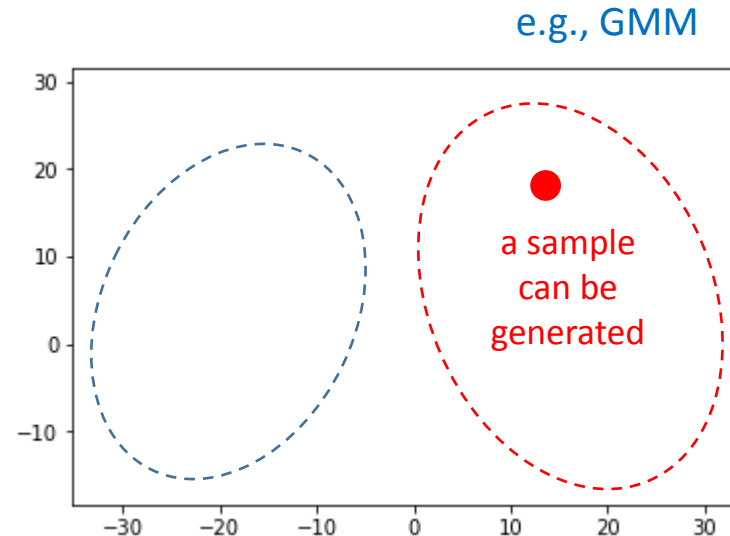
Each data point has 64 dimension

Project the data point in 2 dimension

8

8

64 dimensions

- We can generate a new image which corresponds to the sample.
- It is NOT one of training data points!

6

# Variational AutoEncoder (VAE)

❑ To build a method which does the following procedure systematically.



Each data point
has 64 dimension
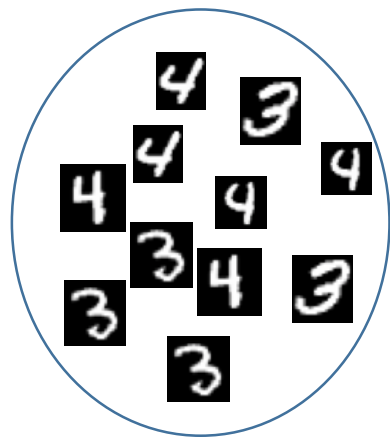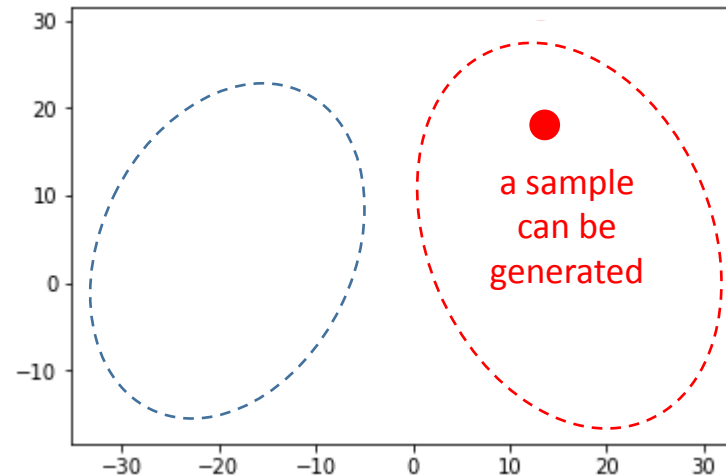
Project the data point in
2 dimension

a sample
can be
generated

8

8

64 dimensions

▪ We can generate a new image which corresponds to the sample.
▪ It is NOT one of training data points!

# Idea of VAE

- ❑ Assuming that there is a complex model parameterized with "θ"
- ❑ The model generates data $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$ given a latent variable "z": $p_\theta(x|z)$
- ❑ Also, the model maps data set into the latent space: $p_\theta(z|x)$

Data space (x)　　　　　Latent space (z)　　　　　Data space (x)

a sample

It is intractable

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

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

- ❑ This Intractability is well known, which can be handled with 1) Markov Chain Monte Carlos (MCMC) and 2) Variational Inference (VI).

- ❑ VAE uses the idea of Variational Inference and so the term "Variational" is in the name.

Data space (x)          Latent space (z)          Data space (x)



a sample

It is intractable      $p_\theta(z^{(i)} \mid \mathbf{x}^{(i)})$          $p_\theta(\mathbf{x}^{(i)} \mid z^{(i)})$

❑ Auto Encoder is a neural network which reproduces its input

Data space (x)　　　　　　　　Latent space (z)　　　　　　　　Data space (x)

$$X$$

Encoder　　　　　　　　Decoder

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

❑ Auto Encoder is a neural network which reproduces its input



Data space (x)

Latent space (z)

Data space (x)

X

Encoder

Decoder

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

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

approximation

$$q_\varphi(z^{(i)} \mid x^{(i)})$$

$z_1$    $z_2$

▪ Assuming $q_\varphi()$ as Gaussian distribution

12

❑ Auto Encoder is a neural network which reproduces its input

Data space (x)  Latent space (z)  Data space (x)

X

X

Encoder  Decoder

$$q_\varphi(z^{(i)} \mid \mathrm{x}^{(i)})$$  $$p_\theta(\mathrm{x}^{(i)} \mid z^{(i)})$$

2 dimensions

8

8  ![8]

**X**

pixel1
pixel2
⋮
pixel63
pixel64

$\mu_{z1}$
$\sigma_{z1}^2$
$\mu_{z2}$
$\sigma_{z2}^2$

**sampling**

$z1$
$z2$

$z$

$p_\theta(z) \sim N(0, 1)$

**sampling**

$\mu_{x1}$
$\sigma_{x1}^2$
⋮
$\mu_{x64}$
$\sigma_{x64}^2$

pixel1
pixel2
⋮
pixel63
pixel64

8

8  ![8]

**X**

$q_\varphi(z|x) \sim N(\mu, \sigma^2)$

$\{w, b\} \in \varphi$

$p_\theta(x|z) \sim N(\mu, \sigma^2)$ or Bernoulli

$\{w, b\} \in \theta$

# A loss function for AutoEncoder

❑  How can we train the network to obtain the parameter $\varphi$ and $\theta$ ?

-  To train the neural network, <span style="color:red">a loss function is necessary</span>. Then, the parameter "$\varphi$ and $\theta$" can be calculated through a backpropagation.

❑  Let's derive the loss function from the likelihood function $p_\theta(x)$

$$p_\theta(\mathrm{x}^{(1)}, \cdots, \mathrm{x}^{(N)}) = \prod_{i=1}^{N} \log p_\theta(\mathrm{x}^{(i)})$$  ➡️  <span style="color:#2E74B5">Likelihood function showing the probability that given batch data set occur with the parameter θ in the neural network.</span>

$$\log p_\theta(\mathrm{x}^{(1)}, \cdots, \mathrm{x}^{(N)}) = \sum_{i=1}^{N} \log p_\theta(\mathrm{x}^{(i)})$$  ➡️  <span style="color:#2E74B5">(log) likelihood</span>

$$\log p_\theta(\mathrm{x}^{(i)}) = L(\theta, \varphi; \mathrm{x}^{(i)}) + D_{KL}(q_\varphi(z \mid \mathrm{x}^{(i)}) \parallel p_\theta(z \mid \mathrm{x}^{(i)}))$$

❑  Since $D_{KL} \geq 0$, "L" is the lower bound of the likelihood function, which is called "ELBO" (Evidence Lower Bound)

❑  Kullback-Leibnitz divergence showing how difference between two posterior distributions: true posterior p(z|x) and its approximate posterior q(z|x)

❑  <span style="color:red">This term is intractable because of p(z|x).</span>

❑  However, we know $D_{KL} \geq 0$.

$$\log p_\theta(\mathrm{x}^{(1)}, \cdots, \mathrm{x}^{(N)}) = \sum_{i=1}^{N} \boxed{\log p_\theta(\mathrm{x}^{(i)})}$$

$$\log p(\mathrm{x}) = \sum_z q(z \mid \mathrm{x}) \log p(\mathrm{x})$$

$$p(x, z) = p(z, x)$$
$$p(x, z) = p(z \mid x) p(x)$$
$$p(x) = \frac{p(x, z)}{p(z \mid x)}$$

$$= \sum_z q(z \mid \mathrm{x}) \log\left( \frac{p(z, \mathrm{x})}{\mathrm{p}(z \mid \mathrm{x})} \right)$$

$$= \sum_z q(z \mid \mathrm{x}) \log\left( \frac{p(z, \mathrm{x})}{\mathrm{q}(z \mid \mathrm{x})} \frac{q(z \mid \mathrm{x})}{\mathrm{p}(z \mid \mathrm{x})} \right)$$

$$= \sum_z q(z \mid \mathrm{x}) \log\left( \frac{p(z, \mathrm{x})}{\mathrm{q}(z \mid \mathrm{x})} \right) + \sum_z q(z \mid \mathrm{x}) \log\left( \frac{q(z \mid \mathrm{x})}{\mathrm{p}(z \mid \mathrm{x})} \right)$$

$$\boxed{\log p_\theta(\mathrm{x}^{(i)})} = L(\theta, \varphi; \mathrm{x}^{(i)}) + D_{KL}(q_\varphi(z \mid \mathrm{x}^{(i)}) \| p_\theta(z \mid \mathrm{x}^{(i)}))$$

16

❑ By maximizing "L", we can maximize the likelihood function as well

$$\log p_\theta(\mathrm{x}^{(i)}) \geq \boxed{L(\theta, \varphi; \mathrm{x}^{(i)})}$$

Lower bound of the likelihood function

$$L(\theta, \varphi; \mathrm{x}^{(i)}) = \sum_z q_\varphi(z \mid \mathrm{x}) \log\left(\frac{p_\theta(z, \mathrm{x})}{q_\varphi(z \mid \mathrm{x})}\right)$$

$$= \sum_z q_\varphi(z \mid \mathrm{x}) \log\left(\frac{p_\theta(\mathrm{x} \mid z) p_\theta(z)}{q_\varphi(z \mid \mathrm{x})}\right)$$

$$= \sum_z q_\varphi(z \mid \mathrm{x}) \log\left(\frac{p_\theta(z)}{q_\varphi(z \mid \mathrm{x})}\right) + \sum_z q_\varphi(z \mid \mathrm{x}) \log\left(p_\theta(\mathrm{x} \mid z)\right)$$

$$= -D_{KL}\left(q(z \mid \mathrm{x}^{(i)}) \| p(z)\right) + E_{q(z \mid \mathrm{x}^{(i)})}\left(\log p(\mathrm{x}^{(i)} \mid z)\right)$$

# A loss function for AutoEncoder

- ❑ By maximizing "L", we can maximize the likelihood function as well,
- ❑ In other words, the most likely model ($p_\theta$ and $p_\varphi$), which generates the observed data, can be obtained by maximizing the function below.

$$\log p_\theta(\mathrm{x}^{(i)}) \geq \boxed{-D_{KL}(q_\varphi(z \mid \mathrm{x}^{(i)}) \parallel p_\theta(z))} + \boxed{E_{q_\varphi(z \mid \mathrm{x}^{(i)})}(\log p_\theta(\mathrm{x}^{(i)} \mid z))}$$

- ▪ It can be computed in a closed form

- ❑ $q_\varphi(z|x) \sim N(\mu, \sigma^2)$
- ❑ $p_\theta(z) \sim N(0, I)$
- ❑ J: dimension of z

- ❑ $p_\theta(x|z) \sim N(\mu, \sigma^2)$ or Bernoulli
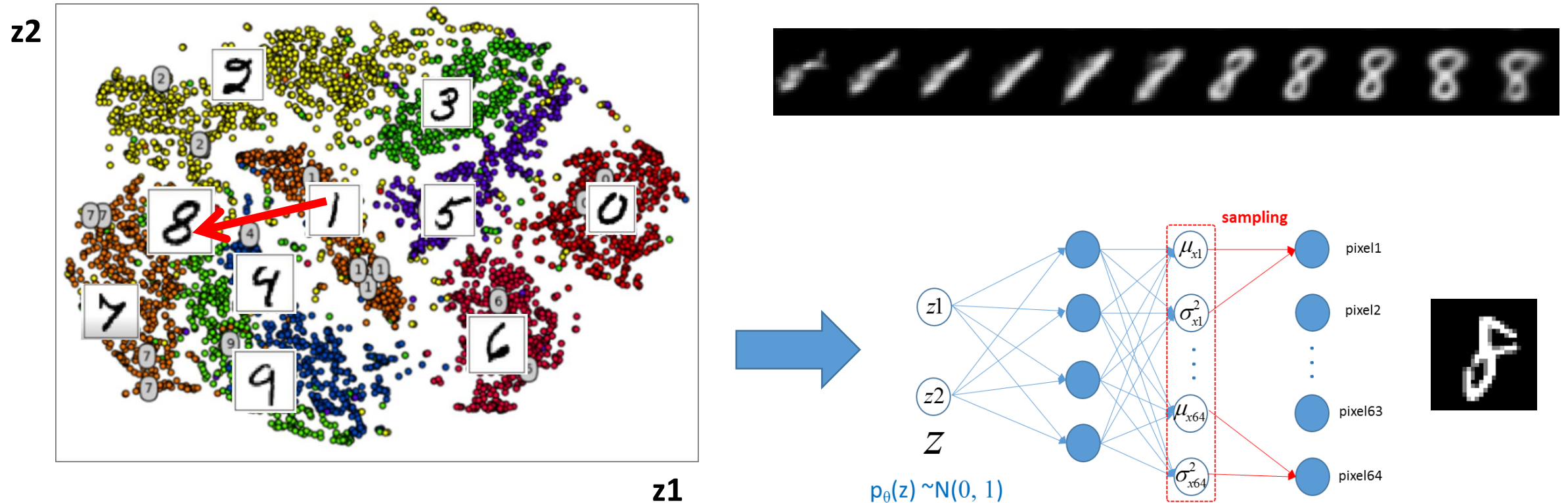- ❑ D: dimension of x

$$= \frac{1}{2} \sum_{j=1}^{J} \left(1 + \log((\sigma_{z_j}^{(i)})^2) - (\mu_{z_j}^{(i)})^2 - (\sigma_{z_j}^{(i)})^2\right)$$

$$= \sum_{j=1}^{D} \left( \frac{1}{2} \log((\sigma_{x_j}^{(i)})^2) + \frac{\left(x_j^{(i)} - \mu_{x_j}\right)^2}{2\sigma_{x_j}^2} \right)$$

$$\boxed{\min \mid x - \hat{x} \mid}$$

Auto-Encoding Variational Bayes (appendix B:
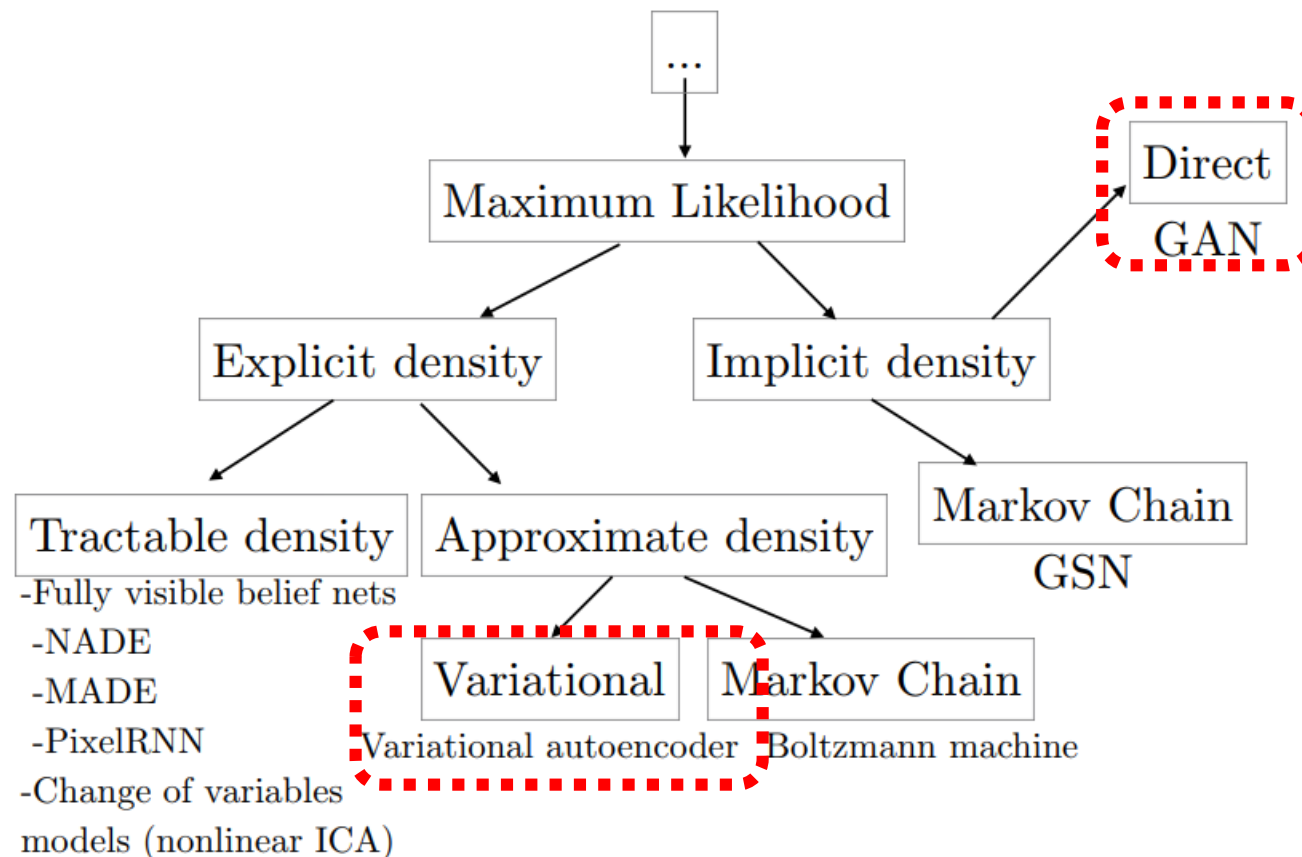derivation) https://arxiv.org/pdf/1312.6114.pdf

- ❑  A generative model based on a neural network (AutoEncoder)
- ❑  Its loss function is derived based on variational inference approach (Variational)
- ❑  The loss function calculates the error used to train Auto Encoder through backpropagation
  - That is the reason why it is called "Variational AutoEncoder" (VAE).
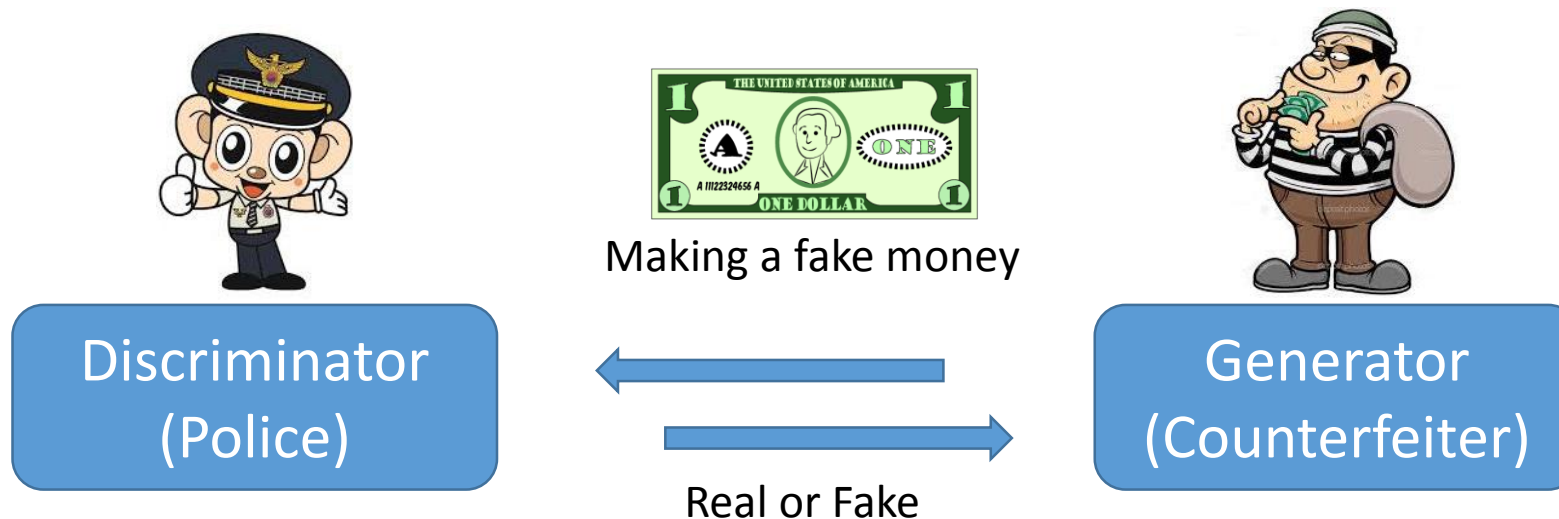
# Generative Adversarial Networks (GAN)

❑ In VAE, we design a latent space which maps to a data space.
❑ Then, a latent variable in the space is used to generate a data sample.
❑ However, actually we are interested in not the latent space but a sample itself.
❑ Then, why do we generate samples directly without the latent space estimation?

https://arxiv.org/pdf/1701.00160.pdf

- ❑ GAN: Generative Adversarial Network
- ❑ Based on game theory to train the system which directly generates a sample
- ❑ Adversarial:

*'GAN framework can naturally be analyzed with the tools of game theory, we call GANs "adversarial".'* - Ian Goodfellow

Making a fake money

| Discriminator (Police) | | Generator (Counterfeiter) |
|---|---|---|

Real or Fake

https://arxiv.org/pdf/1701.00160.pdf

# Theory: formulation of an optimization problem

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

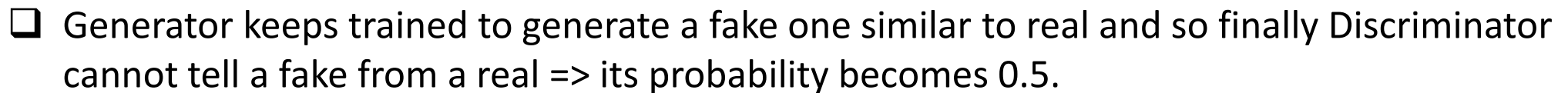- ❑ Expectation that discriminator (D) tells **real is real** (D successes)
- ❑ Training discriminator to maximize it

- ❑ Expectation that D tells **fake is real** (D fails)
- ❑ Training generator to minimize a fake

| notation | description |
|---|---|
| $x \sim p_{data}(x)$ | Real data sample |
| $z \sim p_z(z)$ | A random number from N(0, 1) |
| $G(z)$ | Fake data sample |
| $D(x)=1$ | Probability of discriminator (D) telling that given real data "x" is real |
| $D(G(z))=0$ | Probability of discriminator (D) telling that given fake data "G(z)" is fake |
| $1 - D(G(z))$ | Probability of discriminator (D) telling that given fake data "G(z)" is real |

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



Real data dist.

Fake data dist.

Discriminator

0.5

(a)          (b)          (c)          ...          (d)

❑ Generator keeps trained to generate a fake one similar to real and so finally Discriminator cannot tell a fake from a real => its probability becomes 0.5.

https://arxiv.org/abs/1406.2661

❑ For the fixed generator (G), the optimal discriminator value (D*) is

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

$$\max_D V(D) = \int_x p_{data}(x) \log(D(x)) dx + \int_z p_z(z) \log(1 - D(G(z))) dz$$

$$= \int_x p_{data}(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx$$

$$\frac{dV(D)}{dD} = \frac{p_{data}(x)}{D(x)} - \frac{p_g(x)}{1 - D(x)} = 0$$

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

$$= 0.5 \ (p_g = p_{data})$$

❑ With the optimum value of D*, lower bound of V(G) is

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

$$\min_{G} V(G) = E_{x \sim p_{data}(x)}[\log D^*(x)] + E_{z \sim p_g(x)}[\log(1 - D^*(x)))]$$

**Backup slide**

$$= -\log(4) + 2 \cdot JSD(p_{data} \parallel p_g)$$

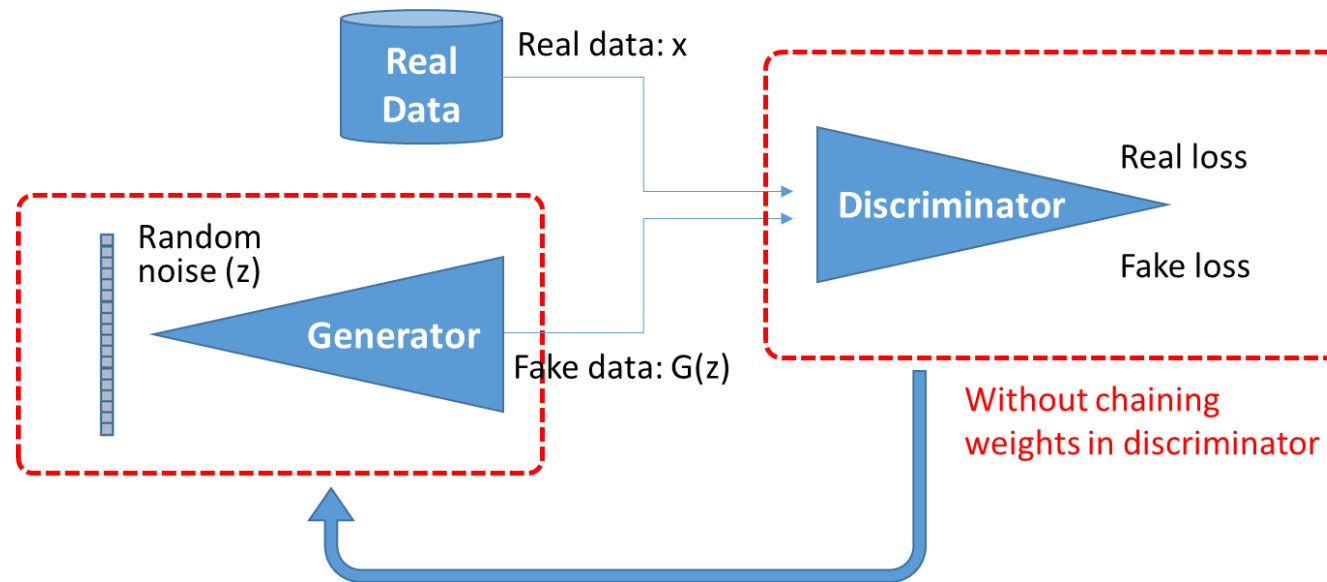This is the minimum value of V(G) when JSD=0 ($p_g=p_{data}$)

❑ JSD: Jensen Shannon divergence
  - A method of measuring the similarity between two probability distribution.
  - 0 ≤ JSD(p|q) ≤1

❑ Given the system below, we train it based on the objective function.
❑ The objective function is derived based on game theory.
  - Generator tries to make a real like fake data to deceive the discriminator
  - Discriminator tries not to be deceived by the generator
❑ In this manner, generator learns how to make a sample close to real data.
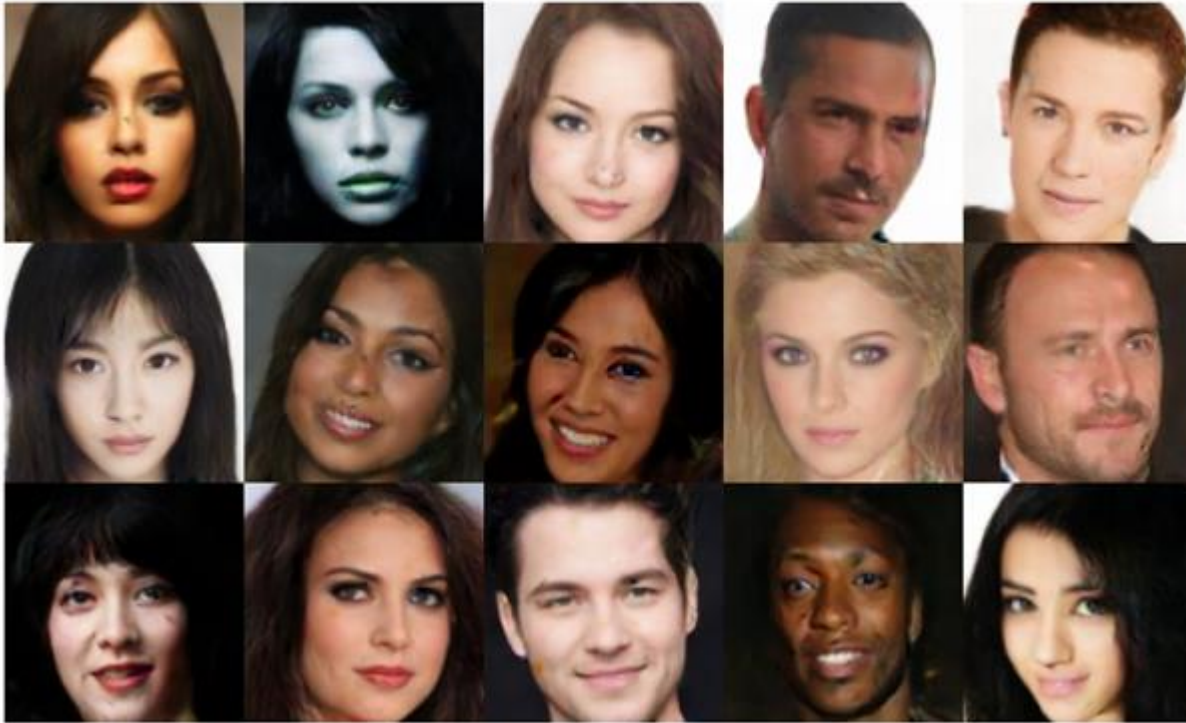❑ It is about how to define the objective function and whether it converges to an optimum solution.

# Applications

Cumulative number of named GAN papers by month

https://deephunt.in/the-gan-zoo-79597dc8c347

# High resolution image generation

❑ [https://arxiv.org/pdf/1703.10717.pdf](https://arxiv.org/pdf/1703.10717.pdf) (BGAN)

❑ https://arxiv.org/pdf/1710.10916.pdf (StackGAN)
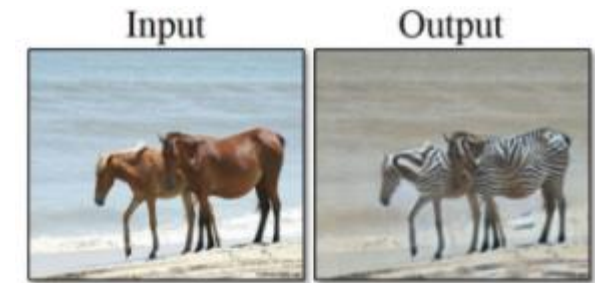
❑ https://github.com/junyanz/CycleGAN



winter Yosemite → summer Yosemite

summer Yosemite → winter Yosemite
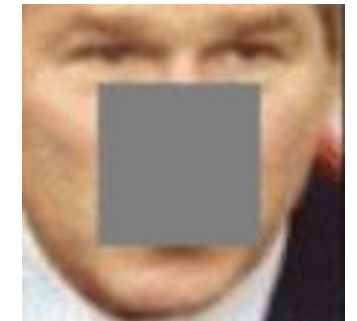


Input    Output

horse → zebra

zebra → horse

apple → orange

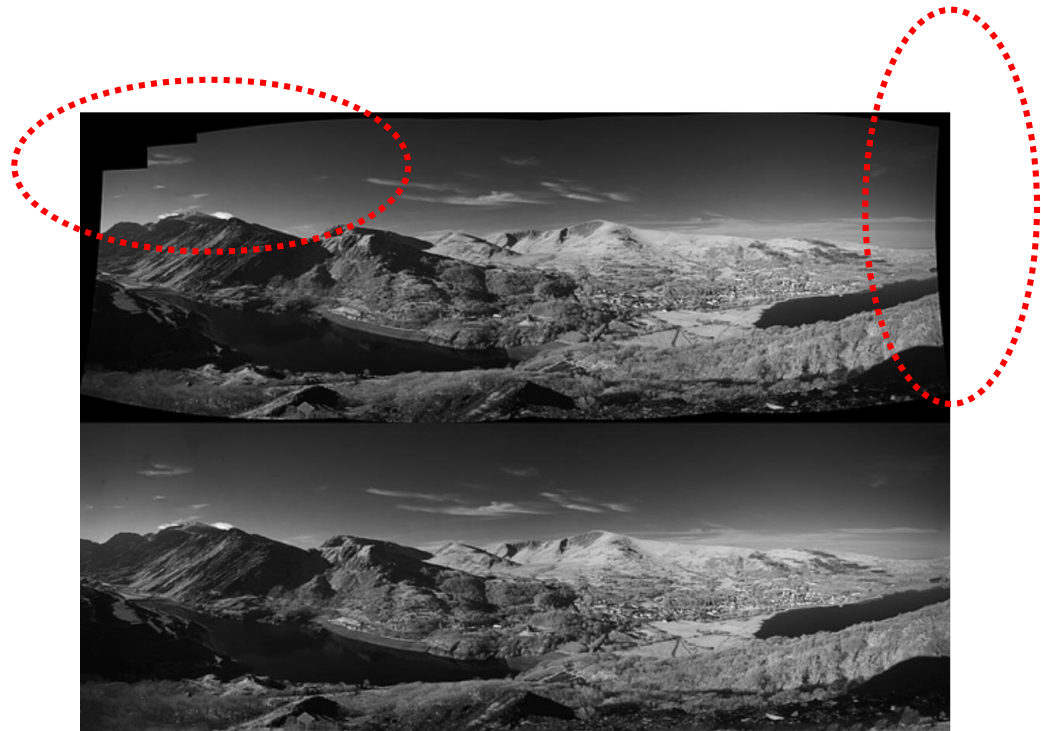orange → apple

❑ http://bamos.github.io/2016/08/09/deep-completion/



33

# Image2Vec

❑ https://arxiv.org/pdf/1511.06434.pdf



A->C = B->D

C-A = D-B

C-A+B=D



man with glasses − man without glasses + woman without glasses = woman with glasses

# Deep Feature Interpolation

❑ https://arxiv.org/pdf/1611.05507.pdf

# Backup Slides

$$\min_{G} V(G) = E_{x \sim p_{data}(x)} \left[ \log \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \right] + E_{x \sim p_g(x)} \left[ \log \frac{p_g(x)}{p_{data}(x) + p_g(x)} \right]$$

$$KL(P \| Q) = \sum P \log \frac{P}{Q}$$

$$JSD(P \| Q) = \frac{1}{2} KL(P \| M) + \frac{1}{2} KL(Q \| M)$$

$$V(G) = V(G) + \log(4) - \log(4)$$

$$= -\log(4) + E_{x \sim p_{data}(x)}[\log D(x)] + E_{x \sim p_q(x)}[\log(1 - D(x))] + \log(4)$$

$$= -\log(4) + \sum p_{data}(x) \log(D*(x)) + \log(2) + \sum p_g(x) \log(1 - D*(x)) + \log(2)$$

$$= -\log(4) + \sum p_{data}(x) \log \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} + \log(2) + \sum p_g(x) \log \frac{p_g(x)}{p_{data}(x) + p_g(x)} + \log(2)$$

$$= -\log(4) + KL\left( p_{data} \| \frac{p_{data} + p_g}{2} \right) + KL\left( p_g \| \frac{p_{data} + p_g}{2} \right)$$

$$= -\log(4) + 2JSD(p_{data} \| p_g)$$