

# PART 4

## Deep Unsupervised Learning

- So far, **supervised learning**

- Discriminative methods:

$$\{\underline{x}_i, \underline{y}_i\}_{i=1}^N \rightarrow p(y|x)$$

- Generative methods:

$$\{\underline{x}_i, \underline{y}_i\}_{i=1}^N \rightarrow P(y), p(x|y)$$

↑              ↘  
class sizes      class conditional

*logistic regression  
softmax regression  
FNN  
CNN*

\* Latent / Hidden Variables:

Intelligence

\* Observed Variables :

Math  
grade

science  
grade

Literature  
grade

\* Latent Variable models.

Market Sentiment

unlabelled data

Gaussian Mixture Model

cluster sizes

cluster characteristics

- Next, **unsupervised learning**:

- *Finite mixture models* for clustering [Skipped]

$$\{\underline{x_i}\}_{i=1}^N \rightarrow P(z), p(x|z)$$

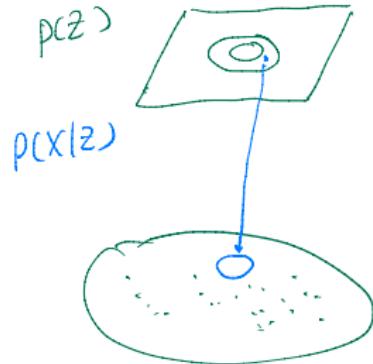
- *Variational autoencoder* for data generation and representation learning

$$\{\underline{x_i}\}_{i=1}^N, p(z) \rightarrow p(x|z) \quad q(z|x) \text{ used in inference}$$

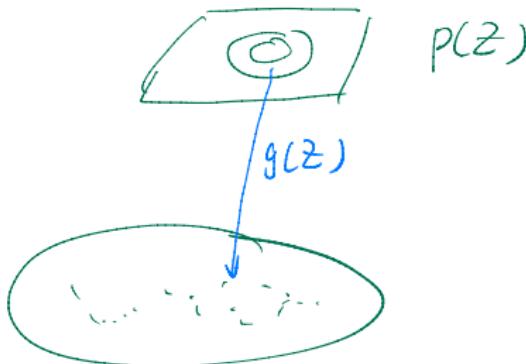
- *Generative adversarial networks* for data generation

$$\{\underline{x_i}\}_{i=1}^N, p(z) \rightarrow x = g(z)$$

VAE



GAN



## Unsupervised Learning

- Unknown true distribution  $P(\mathbf{x})$ .

$$P(\mathbf{x}) \xrightarrow{\text{sampling}} \mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N \xrightarrow{\text{learning}} Q(\mathbf{x})$$

- Objective:

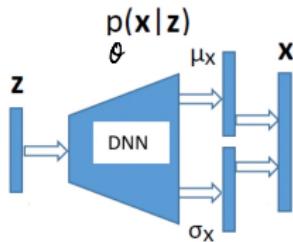
- **Minimizing KL:**  $KL(P||Q)$
- Same as **minimizing cross entropy:**  $H(P, Q)$
- Approximating the cross entropy using data:

$$\begin{aligned} H(P, Q) &= - \int P(\mathbf{x}) \log Q(\mathbf{x}) d\mathbf{x} \\ &\approx - \frac{1}{N} \sum_{i=1}^N \log Q(\mathbf{x}_i) \\ &= - \frac{1}{N} \log Q(\mathcal{D}) \end{aligned}$$

- Same as **maximizing likelihood:**  $\log Q(\mathcal{D})$ .

0	4	1	9	2	1
5	3	6	1	7	2
0	9	1	1	2	4
8	6	9	0	5	6
8	1	9	3	9	8
0	7	4	9	8	0
4	6	0	4	5	6

$$P(z) \Rightarrow$$



X

$$\log P_{\theta}(x) = \frac{1}{N} \sum_{i=1}^N \log \underline{P_{\theta}(x^i)}$$

$$p(z), p(x|z) \Rightarrow p(x) = \int p(z)p(x|z)dz$$

$$\log P_{\theta}(x^i) = \int p(z)p(x^i|z)dz \approx \frac{1}{L} \sum_{j=1}^L p(x^i|z^j)$$

$$z^1 \dots z^L \sim p(z)$$

$$x^i = (x_{i1}, x_{i2}, \dots, x_{iD})$$

$$D = 256 \times 256$$

$$= 65,536$$

$$P_0(x^i | z^e) = C \prod_{j=1}^D e^{-\frac{(x_{ij} - \mu_{ij}^e)^2}{2\sigma^2}}$$

$$= C \times \underbrace{0.99 \times 0.99 \times \dots \times 0.99}_{65K} \approx$$

Analogy:

$X$ : Image of person flying small airplane

$Z$ : income

$$z^l \sim p(z) \quad p(x|z) \quad \begin{cases} \text{negligible} & \checkmark \\ \text{not small} & \end{cases}$$

$$X \rightarrow g(z|x)$$

Distribution of income  
among those flying  
small airplane

$$z^l \sim g(z|x)$$

$$p(x|z^l) \quad \begin{cases} \text{negligible} & \\ \text{not small} & \checkmark \end{cases}$$

$$\text{Want: } z^\ell \sim p(z) \quad \frac{1}{L} \sum_{\ell=1}^L p(x^\ell | z^\ell)$$

$$\text{Do: } z^\ell \sim q(z|x) \quad \frac{1}{L} \sum_{\ell=1}^L p(x^\ell | z^\ell)$$

Bias

$$\begin{aligned}
 & E_{q(z|x)} \left[ \log \frac{p(x,z)}{q(z|x)} \right] = E_{q(z|x)} \left[ \log \frac{p(x,z)}{q(z|x)} \right] \\
 &= \int q(z|x) \log \frac{p(x,z)}{q(z|x)} dz = E_{q(z|x)} \left[ \log \frac{p(z)p(x|z)}{q(z|x)} \right] \\
 &\leq \log \int q(z|x) \frac{p(x,z)}{q(z|x)} dz = E_{q(z|x)} [\log p(x|z)] \\
 &\quad \text{Jensen's inequality} - E_{q(z|x)} \left[ \log \frac{q(z|x)}{p(z)} \right] \\
 &= \log p(x) \\
 &\quad \text{KL}(q(z|x) || p(z))
 \end{aligned}$$

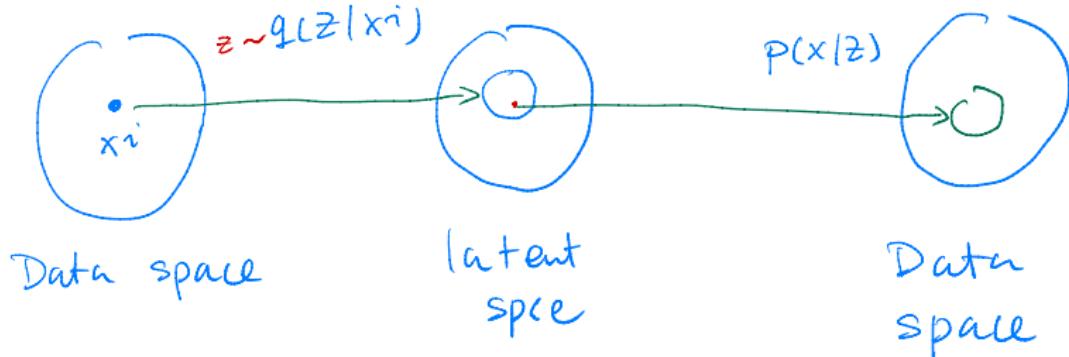
$$\log p(x) \geq E_{q(z|x)} [\log p(x|z)] - \text{KL}(q(z|x) || p(z))$$

Variational lower bound

Evidence lower Bound (ELBO)

## reconstruction error

$$\mathcal{L}(\mathbf{x}^{(i)}, \theta, \phi) = E_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[ \log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] - D_{KL}[q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})]$$



\* How well  $p(x|z)$  match  $x^i$ :  $\log p(x^i|z)$

\* How well  $x^i$  is reconstructed from  $q(z|x^i)$ .

## Representation learning

NLP: BERT, Word2Vec, ... ELMO

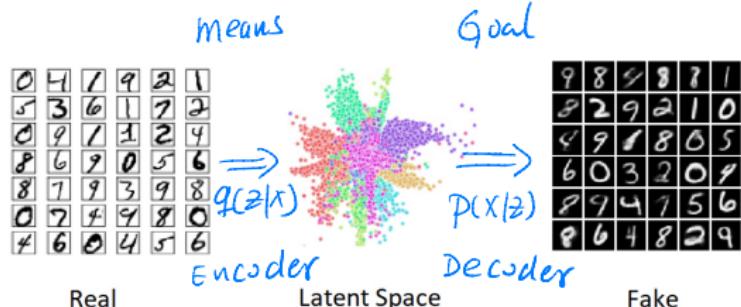
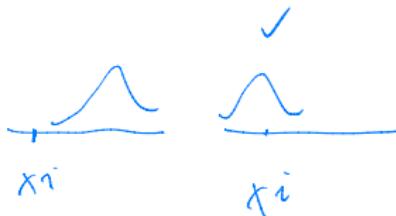
Images: AE  
DAE  
VAE

Contrastive learning = SimCLR

Self-Supervised learning, DINO

2022-04-06

VAE



$$\mathcal{L}(x^{(i)}, \theta, \phi) = E_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)}|z)] - D_{KL}[q_\phi(z|x^{(i)}) || p_\theta(z)]$$

How to  
Compute

(1)

(2)

?

Depends on  $\phi$

Depends on  $\phi$ ?

No

$$\mathcal{L}_1 = E_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] \approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)}|\underline{\mathbf{z}}^{(i,l)})$$

where  $\mathbf{z}^{(i,l)} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$ .

Reparameterization : scalar  $x, z$

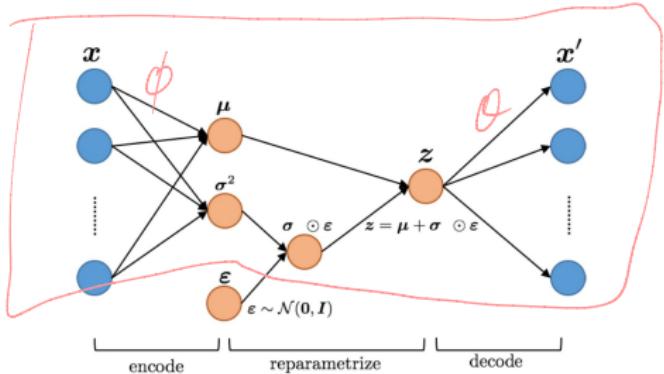
$$q(z|x) = N\left(\frac{x^2+2}{\mu}, \frac{x^4}{\sigma^2}\right) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(z-(x^2+2))^2}{2\sigma^2}}$$

$$\frac{dz}{dx} = ? \quad \frac{dq}{dx} = .. \text{OK}$$

$$z = \frac{x^2+2}{\mu} + \frac{x^2}{\sigma} \varepsilon \quad \varepsilon \sim N(0, 1)$$

$$\frac{dz}{dx} = 2x + 2x\varepsilon$$

$$\underline{\mathbf{z} = \mu_z(\mathbf{x}, \phi) + \sigma_z(\mathbf{x}, \phi) \odot \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})}$$



$$\begin{aligned}\mathcal{L}_1 &= E_{\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x}^{(i)})} \left[ \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}) \right] \\ &= \frac{1}{L} \sum_{l=1}^L \log p_\theta(x^{(i)} | z^{(i,l)})\end{aligned}$$

$\uparrow \quad \uparrow$   
 $\phi \quad \phi$

$$\begin{aligned}\varepsilon^l &\sim N(0, \mathbf{I}) \\ l &= 1, \dots, L\end{aligned}$$

$$\begin{aligned}z^{(i,l)} &= \mu_z(x^{(i)}, \phi) \\ &+ \sigma_z(x^{(i)}, \phi) \odot \varepsilon^l\end{aligned}$$

$$\mathcal{L}_2 = -\mathcal{D}_{KL}[q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})]. \quad \max \quad 1 + \log \sigma^2 - \sigma^2$$

$$= \frac{1}{2} \sum_{j=1}^J \left( 1 + \log((\sigma_j^{(i)})^2) - \underline{(\mu_j^{(i)})^2} - (\sigma_j^{(i)})^2 \right)$$

mean

$$\frac{q_{\phi}(z|x^{(i)})}{\begin{bmatrix} \mu_1^{(i)} \\ \vdots \\ \mu_J^{(i)} \end{bmatrix}} \stackrel{\max L_2}{\Rightarrow} \frac{p_{\theta}(z)}{\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}}$$

$$\sum \begin{bmatrix} \sigma_1^{(i)} & \dots & 0 \\ 0 & \ddots & \sigma_J^{(i)} \end{bmatrix} \stackrel{\max L_2}{\Rightarrow} \begin{bmatrix} 1 & \dots & 0 \\ 0 & \ddots & 1 \end{bmatrix}$$

LII: GAN

Train D

Real:  $\{x^i\}_{i=1}^m \quad Y_i = 1$

Fake:  $\{g(z^i)\}_{i=1}^m \quad Y_i = 0$

Binary Xentropy

$$\begin{cases} -\log \underline{P(Y_i=1 | x^i)} \\ = -\log \underline{D(x^i)} \end{cases}$$

$$\begin{cases} -\log \underline{P(Y_i=0 | g(z^i))} \\ = -\log [1 - \underline{P(Y_i=1 | g(z^i))}] \\ = -\log [1 - \underline{D(g(z^i))}] \end{cases}$$

D minimize:

$$J = -\frac{1}{m} \sum_{i=1}^m \log D(x^i) - \frac{1}{m} \sum_{i=1}^m \log (1 - D(g(z^i)))$$

D maximize

$$V(\theta_g, \theta_d) = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(g(z^i)))$$

D maximum: zero

Depends on G: ~~Yes~~ No

$$V(\theta_g, \theta_d) = \frac{1}{m} \sum_{i=1}^m \log D(x^{(i)}) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(g(z^{(i)})))$$

V(G, D)

$$V(G) = \max_D V(G, D)$$

min max Same

$$G: \min_G V(G) = \min_G \max_D V(G, D)$$

G

G

D

$$\min_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(g(z^{(i)})))$$

Theory

$$\min_{\theta_g} \frac{1}{m} \sum_{i=1}^m -\log D(g(z^{(i)}))$$

prob  
of  
false

{ max  
min }

practice

## The GAN training algorithm (Goodfellow et al 2014)

**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  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
    • Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{data}(x)$ .
    • Update the discriminator by ascending its stochastic gradient:
      
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))] .$$

  end for
  • Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
  • Update the generator by descending its stochastic gradient:
    
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))) .$$

end for
```

\* How is GAN training related to minimizing  $J_S(\Pr | P_g)$ ?

\* Why cannot we train the discriminator to optimal?

\* Why use the following loss for the generator?

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

$D$  maximizes:  $\frac{1}{m} \sum_{i=1}^m \log D(\mathbf{x}^{(i)}) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\underline{g(\mathbf{z}^{(i)})}))$

$$V(G, D) = E_{\mathbf{x} \sim p_r} [\log D(\mathbf{x})] + E_{\mathbf{x} \sim p_g} [\log (1 - D(\mathbf{x}))]$$
$$= \int \left[ \underbrace{\Pr_r(\mathbf{x})}_{\text{Fixed}} \log D(\mathbf{x}) + \underbrace{\Pr_g(\mathbf{x})}_{\text{Fixed}} \log (1 - D(\mathbf{x})) \right] d\mathbf{x}$$

optimal  $D$ :  $D^*(\mathbf{x}) = \frac{\Pr_r(\mathbf{x})}{\Pr_r(\mathbf{x}) + \Pr_g(\mathbf{x})}$

Thimbleback:  $\theta = P(H)$  Data: 3H, 7T

$$\ell(\theta | \text{Data}) = 0.3 \log \theta + 0.7 \log (1 - \theta)$$

MLE:  $\hat{\theta} = \frac{0.3}{0.3 + 0.7}$

$$V(G, D) = E_{x \sim p_r} [\log D(x)] + E_{x \sim p_g} [\log (1 - D(x))]$$

$$= \int \left[ \underbrace{p_r(x) \log D(x)}_{\text{Fixed}} + \underbrace{p_g(x) \log (1 - D(x))}_{\text{Fixed}} \right] dx$$

$z \quad g(z)$   
 $\| \longrightarrow \bar{x}$

optimal  $D$ :  $D^*(x) = \frac{p_r(x)}{p_r(x) + p_g(x)}$

$$G: \min_G \max_D V(G, D) = \min_G V(G, D^*) \quad p_a = \frac{1}{2}(p_r + p_g)$$

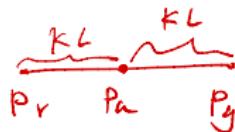
$$V(G, D^*) = E_{x \sim p_r} \left[ \log \frac{p_r(x)}{p_r(x) + p_g(x)} \right] + E_{x \sim p_g} \left[ \log \frac{p_g(x)}{p_r(x) + p_g(x)} \right]$$

$$= E_{x \sim p_r} \left[ \log \frac{p_r(x)}{p_a(x)} \right] + E_{x \sim p_g} \left[ \log \frac{p_g(x)}{p_a(x)} \right] - 2 \log 2$$

$$= 2 \underbrace{\left[ KL(p_r || p_a) + KL(p_g || p_a) \right]}_{JS(p_r || p_g)} - 2 \log 2$$

G: minimize  $JS(p_r || p_g)$

$$JS(p_r || p_g)$$



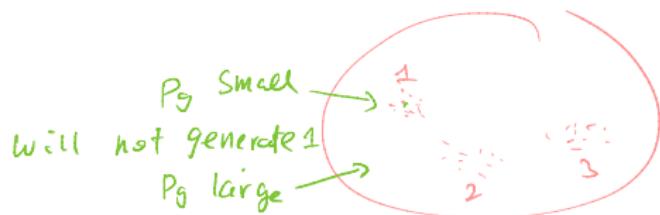
VAE : minimizes

$$KL(p_r || p_g) = \int p_r(x) \log \frac{p_r(x)}{p_g(x)} dx$$

J.S

Different regions in the data space

	Pr	Pg	$\log \frac{Pr(x)}{Pg(x)}$	$\log \frac{Pr(x)}{P_a(x)}$	$\frac{Pg(x)}{P_a(x)}$	$\log \frac{Pg(x)}{Pr(x)}$
Mode collapse	large	Small	large (unlikely)			
Unrealistic images	Small	large	Small (likely)			



Loss for G:

$$\min \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^i)))$$

$$\nabla_{\theta_g} \mathbb{E}_{x \sim P_g} [\log (1 - D(x))] = \nabla_{\theta_g} \int P_g(x) \frac{\log (1 - D(x))}{dx} dx$$

$$= \int \log (1 - D(x)) \nabla_{\theta_g} P_g(x) dx$$

↳ small

slow training

Depends  
on }  
 $\theta_g$  } No

At beginning:

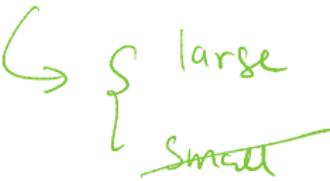
$$x \sim P_g(x): D(x) \left\{ \begin{array}{l} 1 \\ 0 \end{array} \right.$$

$$\min \frac{1}{m} \sum_{i=1}^m -\log D(G(z^i))$$

$$\min \frac{1}{m} \sum_{i=1}^m -\log D_G(z^i)$$

$$\nabla_{\theta_g} E_{x \sim g} [-\log D(x)] = - \nabla_{\theta_g} \int p_g(x) \log D(x) dx$$

$$= - \int \underbrace{\log D(x)}_{\text{large}} \nabla_{\theta_g} p_g(x) dx$$

If  $D$  is :  
close optimal  large  
small

$\Rightarrow$  unstable training

Do not train  $D$  to optimum.

## Difficulty with JS

$$JS(P_r \parallel P_g) = \frac{1}{2} KL(P_r \parallel P_a) + \frac{1}{2} KL(P_g \parallel P_a)$$

$$\text{Supp}(P_r) = \{x \mid P_r(x) > 0\}$$

$\approx 0$

$x \in \mathbb{R}^{dr}$   
 $dr = 286 \times 256$   
 $= 56k$

$$\text{Supp}(P_g) = \{x = g(z) \mid z \in \mathbb{R}^{dz}\}$$

$$dz \ll dr$$

$\approx 0$

$1000$        $56k$

$$P \left( \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \text{ look real} \right)$$

Volume of 2D surface  
 in 3D space = 0

$$\approx 0$$



$$\text{Supp}(P_r) \cap \text{Supp}(P_g) = \emptyset$$

$$\begin{aligned} JS(P_r \| P_g) &= \frac{1}{2} KL(P_r \| P_a) + \frac{1}{2} KL(P_g \| P_a) \\ &= \log 2 \end{aligned}$$

$$KL(P_r \| P_a) = \int_{\text{Supp}(P_r)} P_r(x) \log \frac{P_r(x)}{P_a(x)} dx$$

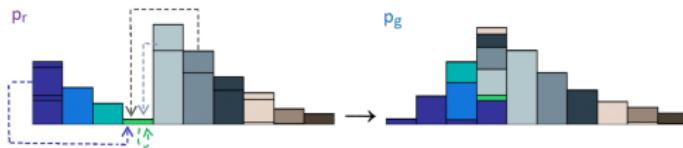
$$P_g(x) \leftarrow x \in \text{Supp}(P_r)$$

$$= \int_{x \in \text{Supp}(P_a)} \frac{P_r(x)}{P_a(x)} \log_2 dx$$

$$= \log 2$$

$$= \frac{1}{2} (P_r(x) + P_g(x))$$

# Earth Mover's Distance / Wasserstein Distance



$\gamma(x, y)$  : amount of dirt from  $x$  to  $y$

$$\textcircled{1} \quad \gamma(x, y) \geq 0$$

$$\textcircled{2} \quad \sum_y \gamma(x, y) = p_r(x)$$

amount of dirt  
at  $x$  at start

$$\textcircled{3} \quad \sum_x \gamma(x, y) = p_g(y)$$

$$\sum_{x,y} \gamma(x, y) = 1$$

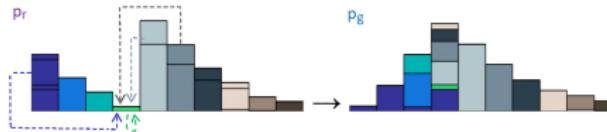
prob dist

Follow  $\gamma(x, y)$ , the total amount of work

$$\sum_{x, y} \frac{\gamma(x, y) \|x - y\|}{\text{Prob Dist}} = E_{(x, y) \sim p} \|x - y\|$$

$$E_{MD(p_r, p_g)} = \min_p E_{(x, y) \sim p} \|x - y\|$$

## Earth Mover's Distance or Wasserstein Distance



$$EMD(p_r, p_g) = \min_{\gamma} \mathbb{E}_{(x,y) \sim \gamma} \|x - y\|$$

$\gamma(x, y)$ : Amount  
to transport  
from  $x$  to  $y$

optimal transport problem

## Wasserstein GAN (WGAN)

$$\arg \min_{\theta} EMD(p_r, p_{\theta})$$

unable to  
optimize  
directly

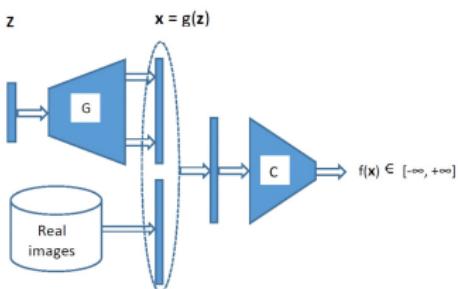
$$EMD(p_r, p_g) = \min_{\gamma} \mathbb{E}_{(x,y) \sim \gamma} \|x - y\|$$

## Kantorovich-Rubinstein Duality

$$EMD(p_r, p_\theta) = \max_{f \in \mathcal{F}_{L-L}} (\mathbb{E}_{x \sim p_r}[f(x)] - \mathbb{E}_{x \sim p_\theta}[f(x)])$$

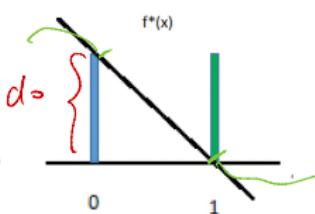
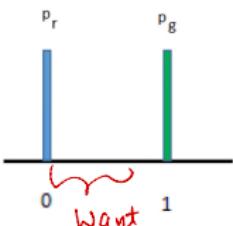
High scores  
for real  
Low scores  
for fake

$$\mathcal{F}_{L-L} \quad EMD_F = k EMD$$



restrict each weight to be from an interval  $[-c, c]$ .

to make  $f$  Lipschitz



Distance  $\leftarrow$  Height

Slope  $\leq \Delta$  everywhere?  
 $\hookrightarrow$  in  $[0, 1]$

$$EMD(p_r, p_\theta) = 1$$

$$EMD(p_r, p_\theta) = \max_{f \in \mathcal{F}_{1-L}} (\mathbb{E}_{x \sim p_r} [f(x)] - \mathbb{E}_{x \sim p_\theta} [f(x)])$$

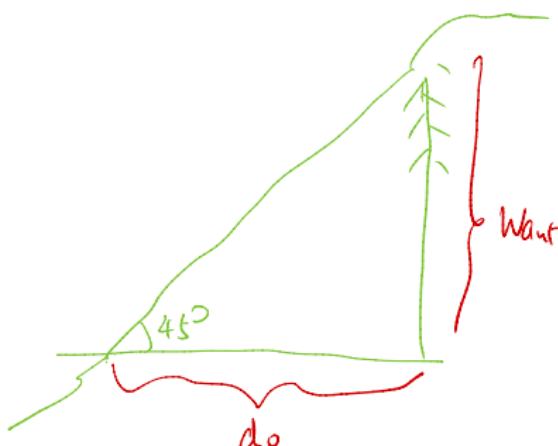
$$\geq \max (b - (a+b))$$

$$f(x) = ax + b$$

$$|a| \leq 1$$

$$= \max -a = 1$$

$$|a| \leq 1$$



Height  $\leftarrow$  Distance









v

