

# Deep Learning

Vazgen Mikayelyan

YSU, Krisp

December 10, 2019

## 1 Generative Adversarial Networks

- Image-to-image translation is a class of vision and graphics problems where the goal is to learn the mapping between an input image and an output image using a training set of aligned image pairs.

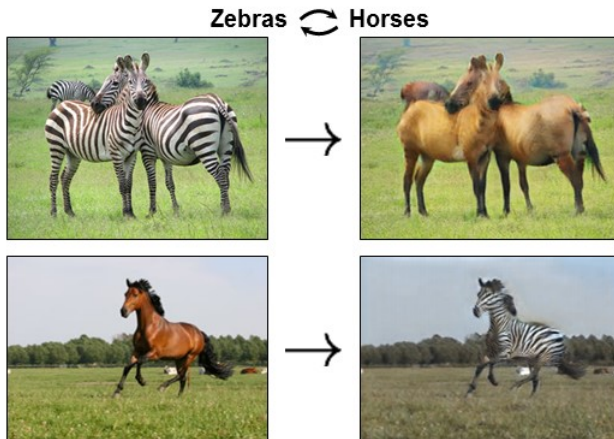
- Image-to-image translation is a class of vision and graphics problems where the goal is to learn the mapping between an input image and an output image using a training set of aligned image pairs.
- However, for many tasks, paired training data will not be available.

- Image-to-image translation is a class of vision and graphics problems where the goal is to learn the mapping between an input image and an output image using a training set of aligned image pairs.
- However, for many tasks, paired training data will not be available.
- So our task is to learn how to translate an image from a source domain  $X$  to a target domain  $Y$  in the absence of paired examples.

- Image-to-image translation is a class of vision and graphics problems where the goal is to learn the mapping between an input image and an output image using a training set of aligned image pairs.
- However, for many tasks, paired training data will not be available.
- So our task is to learn how to translate an image from a source domain  $X$  to a target domain  $Y$  in the absence of paired examples.
- Our goal is to learn a mapping  $G : X \rightarrow Y$  such that the distribution of images from  $G(X)$  is indistinguishable from the distribution  $Y$ .

- Image-to-image translation is a class of vision and graphics problems where the goal is to learn the mapping between an input image and an output image using a training set of aligned image pairs.
- However, for many tasks, paired training data will not be available.
- So our task is to learn how to translate an image from a source domain  $X$  to a target domain  $Y$  in the absence of paired examples.
- Our goal is to learn a mapping  $G : X \rightarrow Y$  such that the distribution of images from  $G(X)$  is indistinguishable from the distribution  $Y$ .
- We will couple it with an inverse mapping  $F : Y \rightarrow X$  and introduce a cycle consistency loss to enforce  $F(G(X)) \approx X$ .

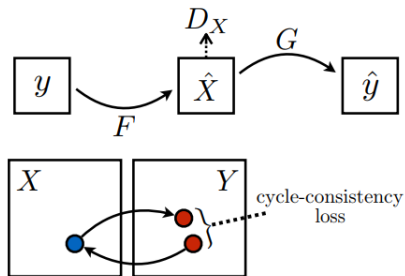
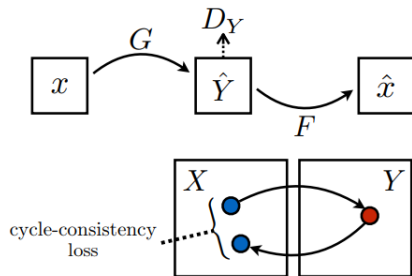
# CycleGANs







# CycleGANs



# CycleGANs

Our loss function will be the following

$$L(G, F, D_X, D_Y) = L_{GAN}(G, D_Y, X, Y) + L_{GAN}(F, D_X, Y, X) + \lambda L_{cyc}(G, F),$$

# CycleGANs

Our loss function will be the following

$$L(G, F, D_X, D_Y) = L_{GAN}(G, D_Y, X, Y) + L_{GAN}(F, D_X, Y, X) + \lambda L_{cyc}(G, F),$$

where

$$\begin{aligned} & L_{GAN}(G, D_Y, X, Y) \\ &= \mathbb{E}_{y \sim p_{data}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{data}(x)} [\log (1 - D_Y(G(x)))] \end{aligned}$$

# CycleGANs

Our loss function will be the following

$$L(G, F, D_X, D_Y) = L_{GAN}(G, D_Y, X, Y) + L_{GAN}(F, D_X, Y, X) + \lambda L_{cyc}(G, F),$$

where

$$\begin{aligned} & L_{GAN}(G, D_Y, X, Y) \\ &= \mathbb{E}_{y \sim p_{data}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{data}(x)} [\log (1 - D_Y(G(x)))] \end{aligned}$$

$$\begin{aligned} & L_{GAN}(F, D_X, Y, X) \\ &= \mathbb{E}_{x \sim p_{data}(x)} [\log D_X(x)] + \mathbb{E}_{y \sim p_{data}(y)} [\log (1 - D_X(F(y)))] \end{aligned}$$

# CycleGANs

Our loss function will be the following

$$L(G, F, D_X, D_Y) = L_{GAN}(G, D_Y, X, Y) + L_{GAN}(F, D_X, Y, X) + \lambda L_{cyc}(G, F),$$

where

$$\begin{aligned} & L_{GAN}(G, D_Y, X, Y) \\ &= \mathbb{E}_{y \sim p_{data}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{data}(x)} [\log (1 - D_Y(G(x)))] \end{aligned}$$

$$\begin{aligned} & L_{GAN}(F, D_X, Y, X) \\ &= \mathbb{E}_{x \sim p_{data}(x)} [\log D_X(x)] + \mathbb{E}_{y \sim p_{data}(y)} [\log (1 - D_X(F(y)))] \end{aligned}$$

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

# CycleGANs

Our loss function will be the following

$$L(G, F, D_X, D_Y) = L_{GAN}(G, D_Y, X, Y) + L_{GAN}(F, D_X, Y, X) + \lambda L_{cyc}(G, F),$$

where

$$\begin{aligned} & L_{GAN}(G, D_Y, X, Y) \\ &= \mathbb{E}_{y \sim p_{data}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{data}(x)} [\log (1 - D_Y(G(x)))] \end{aligned}$$

$$\begin{aligned} & L_{GAN}(F, D_X, Y, X) \\ &= \mathbb{E}_{x \sim p_{data}(x)} [\log D_X(x)] + \mathbb{E}_{y \sim p_{data}(y)} [\log (1 - D_X(F(y)))] \end{aligned}$$

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

We aim to solve

$$G^*, F^* = \arg \min_{G, F} \max_{D_X, D_Y} L(G, F, D_X, D_Y)$$

# WGANs

Let  $\mathcal{X}$  be a compact metric set (e.g.  $[0, 1]^d$ ) and let  $\Sigma$  denote the set of all Borel subsets of  $\mathcal{X}$ .



Let  $\mathcal{X}$  be a compact metric set (e.g  $[0, 1]^d$ ) and let  $\Sigma$  denote the set of all Borel subsets of  $\mathcal{X}$ .

- Total Variation (TV) distance

$$\delta(\mathbb{P}_r, \mathbb{P}_g) = \sup_{A \in \Sigma} |\mathbb{P}_r(A) - \mathbb{P}_g(A)|.$$

Let  $\mathcal{X}$  be a compact metric set (e.g  $[0, 1]^d$ ) and let  $\Sigma$  denote the set of all Borel subsets of  $\mathcal{X}$ .

- Total Variation (TV) distance

$$\delta(\mathbb{P}_r, \mathbb{P}_g) = \sup_{A \in \Sigma} |\mathbb{P}_r(A) - \mathbb{P}_g(A)|.$$

- Kullback-Leibler (KL) divergence

$$KL(\mathbb{P}_r || \mathbb{P}_g) = \int_{\mathbb{R}} \log \left( \frac{p_r(x)}{p_g(x)} \right) p_r(x) d\mu(x).$$

Let  $\mathcal{X}$  be a compact metric set (e.g  $[0, 1]^d$ ) and let  $\Sigma$  denote the set of all Borel subsets of  $\mathcal{X}$ .

- Total Variation (TV) distance

$$\delta(\mathbb{P}_r, \mathbb{P}_g) = \sup_{A \in \Sigma} |\mathbb{P}_r(A) - \mathbb{P}_g(A)|.$$

- Kullback-Leibler (KL) divergence

$$KL(\mathbb{P}_r \| \mathbb{P}_g) = \int_{\mathbb{R}} \log \left( \frac{p_r(x)}{p_g(x)} \right) p_r(x) d\mu(x).$$

- Jensen-Shannon (JS) distance

$$JS(\mathbb{P}_r \| \mathbb{P}_g) = \frac{1}{2} (KL(\mathbb{P}_r \| \mathbb{P}_m) + KL(\mathbb{P}_g \| \mathbb{P}_m))$$

$$\text{where } \mathbb{P}_m = \frac{\mathbb{P}_r + \mathbb{P}_g}{2}.$$

- The Earth-Mover (EM) distance or Wasserstein-1

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} (\|x - y\|)$$

where  $\Pi(\mathbb{P}_r, \mathbb{P}_g)$  denotes the set of all joint distributions  $\gamma(x, y)$  whose marginals are respectively  $\mathbb{P}_r$  and  $\mathbb{P}_g$ .

# Example

Let  $Z \sim U[0, 1]$ ,  $\mathbb{P}_0$  be the distribution of points  $(0, Z) \in \mathbb{R}^2$  and  $g_\theta(z) = (\theta, z)$ , then

# Example

Let  $Z \sim U[0, 1]$ ,  $\mathbb{P}_0$  be the distribution of points  $(0, Z) \in \mathbb{R}^2$  and  $g_\theta(z) = (\theta, z)$ , then

- $\delta(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} 1, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$

# Example

Let  $Z \sim U[0, 1]$ ,  $\mathbb{P}_0$  be the distribution of points  $(0, Z) \in \mathbb{R}^2$  and  $g_\theta(z) = (\theta, z)$ , then

- $\delta(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} 1, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$
- $KL(\mathbb{P}_0 \| \mathbb{P}_\theta) = \begin{cases} \infty, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$

# Example

Let  $Z \sim U[0, 1]$ ,  $\mathbb{P}_0$  be the distribution of points  $(0, Z) \in \mathbb{R}^2$  and  $g_\theta(z) = (\theta, z)$ , then

- $\delta(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} 1, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$
- $KL(\mathbb{P}_0 \parallel \mathbb{P}_\theta) = \begin{cases} \infty, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$
- $JS(\mathbb{P}_0 \parallel \mathbb{P}_\theta) = \begin{cases} \log 2, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$



# Example

Let  $Z \sim U[0, 1]$ ,  $\mathbb{P}_0$  be the distribution of points  $(0, Z) \in \mathbb{R}^2$  and  $g_\theta(z) = (\theta, z)$ , then

- $\delta(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} 1, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$
- $KL(\mathbb{P}_0 \| \mathbb{P}_\theta) = \begin{cases} \infty, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$
- $JS(\mathbb{P}_0 \| \mathbb{P}_\theta) = \begin{cases} \log 2, & \text{if } \theta \neq 0 \\ 0, & \text{if } \theta = 0, \end{cases}$
- $W(\mathbb{P}_0, \mathbb{P}_\theta) = |\theta|$

Note that

Note that

- All distances other than EM are not continuous.

Note that

- All distances other than EM are not continuous.
- When  $\theta_t \rightarrow 0$ , the sequence  $(\mathbb{P}_{\theta_t})_{t \in \mathbb{N}}$  converges to  $\mathbb{P}_0$  under the EM distance, but does not converge at all under either the JS, KL, reverse KL or TV divergences.

Note that

- All distances other than EM are not continuous.
- When  $\theta_t \rightarrow 0$ , the sequence  $(\mathbb{P}_{\theta_t})_{t \in \mathbb{N}}$  converges to  $\mathbb{P}_0$  under the EM distance, but does not converge at all under either the JS, KL, reverse KL or TV divergences.
- Only EM distance has informative gradient.

## Definition 1

Let  $X$  and  $Y$  be normed vector spaces. A function  $f : X \rightarrow Y$  is called

- $K$ -Lipschitz if there exists a real constant  $K > 0$  such that, for all  $x_1$  and  $x_2$  in  $X$

$$\|f(x_1) - f(x_2)\| \leq K\|x_1 - x_2\|.$$

- local Lipschitz if for every  $x \in X$  there exists a neighbourhood  $U$  of  $x$  such that  $f$  is Lipschitz on  $U$ .

## Definition 1

Let  $X$  and  $Y$  be normed vector spaces. A function  $f : X \rightarrow Y$  is called

- $K$ -Lipschitz if there exists a real constant  $K > 0$  such that, for all  $x_1$  and  $x_2$  in  $X$

$$\|f(x_1) - f(x_2)\| \leq K\|x_1 - x_2\|.$$

- local Lipschitz if for every  $x \in X$  there exists a neighbourhood  $U$  of  $x$  such that  $f$  is Lipschitz on  $U$ .

## Theorem 1

If function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  has bounded gradient, then  $f$  is a Lipschitz function.

**Theorem 1.** *Let  $\mathbb{P}_r$  be a fixed distribution over  $\mathcal{X}$ . Let  $Z$  be a random variable (e.g Gaussian) over another space  $\mathcal{Z}$ . Let  $g : \mathcal{Z} \times \mathbb{R}^d \rightarrow \mathcal{X}$  be a function, that will be denoted  $g_\theta(z)$  with  $z$  the first coordinate and  $\theta$  the second. Let  $\mathbb{P}_\theta$  denote the distribution of  $g_\theta(Z)$ . Then,*

- 1. If  $g$  is continuous in  $\theta$ , so is  $W(\mathbb{P}_r, \mathbb{P}_\theta)$ .*
- 2. If  $g$  is locally Lipschitz and satisfies regularity assumption 1, then  $W(\mathbb{P}_r, \mathbb{P}_\theta)$  is continuous everywhere, and differentiable almost everywhere.*
- 3. Statements 1-2 are false for the Jensen-Shannon divergence  $JS(\mathbb{P}_r, \mathbb{P}_\theta)$  and all the KLs.*



**Theorem 1.** *Let  $\mathbb{P}_r$  be a fixed distribution over  $\mathcal{X}$ . Let  $Z$  be a random variable (e.g. Gaussian) over another space  $\mathcal{Z}$ . Let  $g : \mathcal{Z} \times \mathbb{R}^d \rightarrow \mathcal{X}$  be a function, that will be denoted  $g_\theta(z)$  with  $z$  the first coordinate and  $\theta$  the second. Let  $\mathbb{P}_\theta$  denote the distribution of  $g_\theta(Z)$ . Then,*

- 1. If  $g$  is continuous in  $\theta$ , so is  $W(\mathbb{P}_r, \mathbb{P}_\theta)$ .*
- 2. If  $g$  is locally Lipschitz and satisfies regularity assumption 1, then  $W(\mathbb{P}_r, \mathbb{P}_\theta)$  is continuous everywhere, and differentiable almost everywhere.*
- 3. Statements 1-2 are false for the Jensen-Shannon divergence  $JS(\mathbb{P}_r, \mathbb{P}_\theta)$  and all the KLs.*

**Corollary 1.** *Let  $g_\theta$  be any feedforward neural network<sup>4</sup> parameterized by  $\theta$ , and  $p(z)$  a prior over  $z$  such that  $\mathbb{E}_{z \sim p(z)}[\|z\|] < \infty$  (e.g. Gaussian, uniform, etc.). Then assumption 1 is satisfied and therefore  $W(\mathbb{P}_r, \mathbb{P}_\theta)$  is continuous everywhere and differentiable almost everywhere.*

**Theorem 2.** *Let  $\mathbb{P}$  be a distribution on a compact space  $\mathcal{X}$  and  $(\mathbb{P}_n)_{n \in \mathbb{N}}$  be a sequence of distributions on  $\mathcal{X}$ . Then, considering all limits as  $n \rightarrow \infty$ ,*

1. *The following statements are equivalent*

- $\delta(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$  with  $\delta$  the total variation distance.
- $JS(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$  with  $JS$  the Jensen-Shannon divergence.

2. *The following statements are equivalent*

- $W(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$ .
- $\mathbb{P}_n \xrightarrow{\mathcal{D}} \mathbb{P}$  where  $\xrightarrow{\mathcal{D}}$  represents convergence in distribution for random variables.

3.  $KL(\mathbb{P}_n \| \mathbb{P}) \rightarrow 0$  or  $KL(\mathbb{P} \| \mathbb{P}_n) \rightarrow 0$  imply the statements in (1).

4. *The statements in (1) imply the statements in (2).*

## Summary 1

*Wasserstein (or EM) loss for neural networks is continuous and differentiable almost everywhere. Moreover, Convergence in KL implies convergence in TV and JS which implies convergence in EM.*

## Kantorovich-Rubinstein duality 2

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} (\mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)])$$

where the supremum is over all the 1-Lipschitz functions  $f : \mathcal{X} \rightarrow \mathbb{R}$ .

## Kantorovich-Rubinstein duality 2

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} (\mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)])$$

where the supremum is over all the 1-Lipschitz functions  $f : \mathcal{X} \rightarrow \mathbb{R}$ .

Note that if we replace  $\|f\|_L \leq 1$  for  $\|f\|_L \leq K$  (consider  $K$ -Lipschitz for some constant  $K$ ) then we end up with  $K \cdot W(\mathbb{P}_r, \mathbb{P}_g)$ .

## Kantorovich-Rubinstein duality 2

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} (\mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)])$$

where the supremum is over all the 1-Lipschitz functions  $f : \mathcal{X} \rightarrow \mathbb{R}$ .

Note that if we replace  $\|f\|_L \leq 1$  for  $\|f\|_L \leq K$  (consider  $K$ -Lipschitz for some constant  $K$ ) then we end up with  $K \cdot W(\mathbb{P}_r, \mathbb{P}_g)$ .

Therefore, if we have a parameterized family of functions  $\{f_w\}_{w \in \mathcal{W}}$  that are all  $K$ -Lipschitz for some  $K$ , we could consider solving the problem

$$\max_{w \in \mathcal{W}} (\mathbb{E}_{x \sim \mathbb{P}_r} [f_w(x)] - \mathbb{E}_{z \sim p(z)} [f_w(g_\theta(z))])$$

for estimating  $W(\mathbb{P}_r, \mathbb{P}_\theta)$ .