

- GANs have become very popular for learning deep generative models

- Informally, the main idea is:

- Two competing neural network models

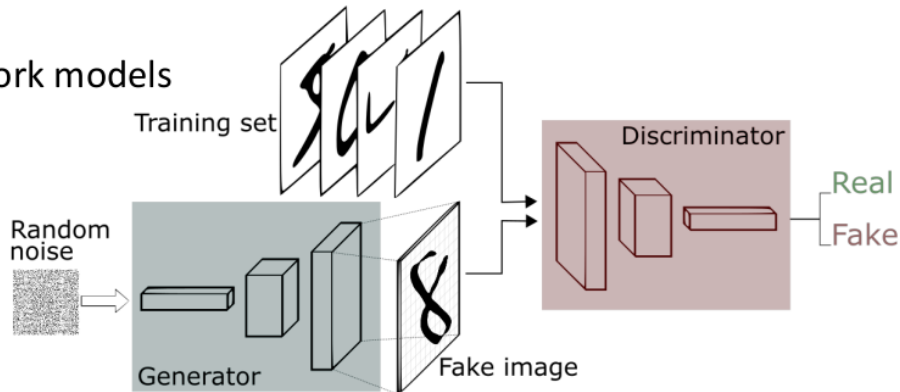
- Generator: takes noise as input and generates ("fake") samples

- Discriminator: receives samples from both generator and training data

and has to distinguish between the two → classify input as "real" or "fake"

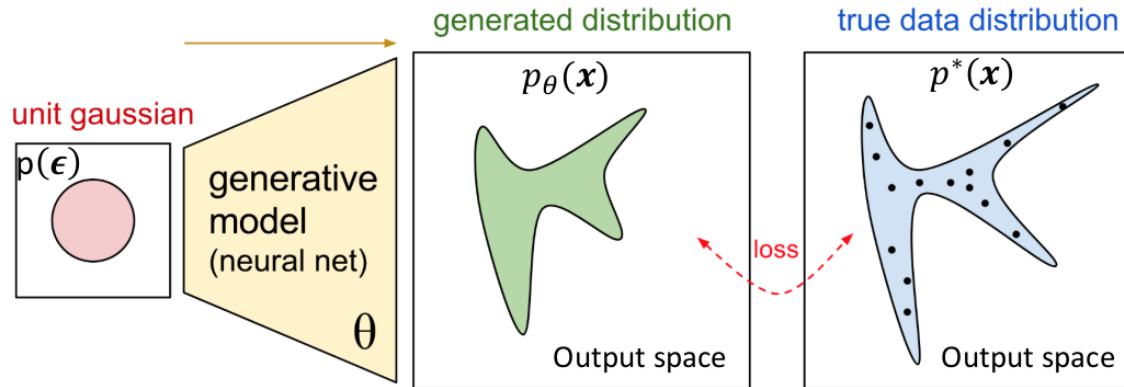
- Goal: Train the generator in such a way that the discriminator can **not** distinguish between real and "fake" samples

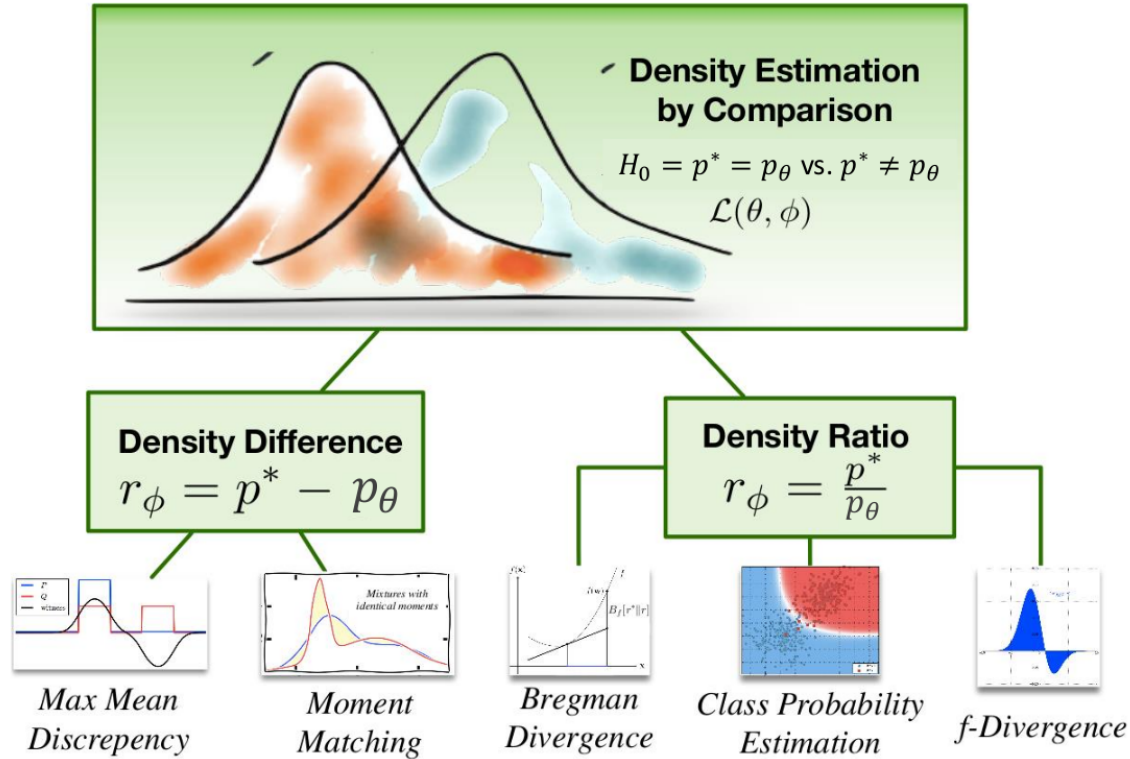
➤ In this case, the generator generates realistic examples



*From "Mining Massive Datasets - Graphs Deep Generative Models.", Prof. Dr. Stephan Günnemann et.al. TUM, 2019.*

- What can we do if we can **easily draw samples** from the model but we cannot evaluate the density?
- Idea: We can use **any** method that compares two **sets of samples**.
  - This process is called ***density estimation by comparison***.
- We **test** the **hypothesis** that the true data distribution  $p^*(x)$  and the model distribution  $p_\theta(x)$  are **equal**





From "Mining Massive Datasets - Graphs Deep Generative Models.", Prof. Dr. Stephan Günnemann et.al. TUM, 2019.

- Density ratio  $r^*(\mathbf{x}) = p^*(\mathbf{x}) / p_\theta(\mathbf{x})$ 
  - In the best case always 1, i.e. the two distributions are indistinguishable
  - However, we cannot compute ratio in closed form/easily
- Idea: Approximate the true density ratio  $r^*(\mathbf{x})$  by  $r_\phi(\mathbf{x})$ 
  - Finding the approximation  $r_\phi(\mathbf{x})$  often means solving again a learning problem
- Thus, we get the following general principle for learning
  - **Optimize Ratio loss:**  
approximate the true density ratio  $r^*(\mathbf{x})$  (i.e. learning  $\phi$ )
  - **Optimize Generative loss:**  
drive the density ratio towards 1 (i.e. learning  $\theta$ )
  - Essentially a bi-level optimization problem, which is usually just solved alternatingly

- Density ratio  $r^*(\mathbf{x}) = p^*(\mathbf{x}) / p_\theta(\mathbf{x})$ 
  - In the best case always 1, i.e. the two distributions are indistinguishable
  - However, we cannot compute ratio in closed form/easily
- Idea: Approximate the true density ratio  $r^*(\mathbf{x})$  by  $r_\phi(\mathbf{x})$ 
  - Finding the approximation  $r_\phi(\mathbf{x})$  often means solving again a learning problem
- Thus, we get the following general principle for learning
  - **Optimize Ratio loss:**  
approximate the true density ratio  $r^*(\mathbf{x})$  (i.e. learning  $\phi$ )
  - **Optimize Generative loss:**  
drive the density ratio towards 1 (i.e. learning  $\theta$ )
  - Essentially a bi-level optimization problem, which is usually just solved alternatingly

- Let  $Y$  denote a random variable which assigns label  $Y = 1$  to samples from the true data distribution;  $Y = 0$  to those from the generator distribution
- Then  $p^*(\mathbf{x}) = p(\mathbf{x} \mid Y = 1)$  and  $p_\theta(\mathbf{x}) = p(\mathbf{x} \mid Y = 0)$

- Denote  $P(Y = 1) = \pi$ . From Bayes we have:

$$r^*(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p_\theta(\mathbf{x})} = \frac{p(Y = 1 \mid \mathbf{x})}{p(Y = 0 \mid \mathbf{x})} \frac{1 - \pi}{\pi}$$

- Apparently density ratio estimation is equal to **class-probability estimation**
  - Simply speaking: we can consider a classifier for  $\mathbf{x}$  (predicting labels  $Y=0$  or  $1$ )
- Specify a scoring function or a **discriminator**  $D_\phi(\mathbf{x}) = p(Y = 1 \mid \mathbf{x})$
- e.g. logistic regression or a neural network

*From "Mining Massive Datasets - Graphs Deep Generative Models.", Prof. Dr. Stephan Günnemann et.al. TUM, 2019.*

- Specify a scoring function or a **discriminator**  $D_\phi(\mathbf{x}) = p(Y = 1 | \mathbf{x})$ 
  - e.g. logistic regression or a neural network
- For learning  $D_\phi(\mathbf{x})$ , we need a loss function, e.g., the cross-entropy loss:

$$\begin{aligned}\mathcal{L}_{\theta, \phi} &= \mathbb{E}_{p(\mathbf{x}|y)p(y)}[-y \log[D_\phi(\mathbf{x})] - (1 - y) \log[1 - D_\phi(\mathbf{x})]] \\ &= \pi \mathbb{E}_{p^*(\mathbf{x})}[-\log D_\phi(\mathbf{x})] + (1 - \pi) \mathbb{E}_{p(\mathbf{z})}[-\log[1 - D_\phi(f_\theta(\mathbf{z}))]]\end{aligned}$$

*From "Mining Massive Datasets - Graphs Deep Generative Models.", Prof. Dr. Stephan Günnemann et.al. TUM, 2019.*

1. Solving 
$$\phi^*(\theta) = \underset{\phi}{\operatorname{argmin}} \mathcal{L}_{\theta, \phi}$$
 leads to the "best" discriminator for a given generative model ( $\theta$ )
    - That is, we well approximate 
$$r^*(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p_{\theta}(\mathbf{x})} = \frac{p(Y=1 | \mathbf{x})}{p(Y=0 | \mathbf{x})} \approx \frac{D_{\phi^*(\theta)}(\mathbf{x})}{1 - D_{\phi^*(\theta)}(\mathbf{x})}$$
      - here w.l.o.g. we set  $\pi = 0.5$
  2. We aim to drive the density ratio  $r^*(\mathbf{x})$  towards 1
    - aim:  $p(Y = 1 | \mathbf{x}) = p(Y = 0 | \mathbf{x})$
    - That is, find generative model ( $\theta$ ) such that (even) the "best" discriminator cannot distinguish the classes
- Solving 
$$\theta^* = \underset{\theta}{\operatorname{argmax}} \mathcal{L}_{\theta, \phi^*(\theta)}$$



- Generator and discriminator play a **minimax game**:

$$\min_{\theta} \max_{\phi} \pi \mathbb{E}_{p^*(x)} [\log D_{\phi}(x)] + (1 - \pi) \mathbb{E}_{p(z)} [\log[1 - D_{\phi}(f_{\theta}(z))]]$$

- Discriminator: aims to distinguish between (samples from)  $p^*(x)$  and  $p_{\theta}(x)$

➤ Maximization

// minimization of cross-entropy

- Generator: aims to generate samples that are indistinguishable

➤ Minimization

// maximization of (lowest) cross-entropy

- This bilevel problem is typically tackled via alternating optimization.

- **Ratio loss** (discriminator loss) optimization:

$$\min_{\phi} \pi \mathbb{E}_{p^*(x)} [-\log D_{\phi}(x)] + (1 - \pi) \mathbb{E}_{p(z)} [-\log[1 - D_{\phi}(f_{\theta}(z))]]$$

- **Generative loss** optimization:

$$\min_{\theta} \mathbb{E}_{p(z)} [\log[1 - D_{\phi}(f_{\theta}(z))]]$$

*From "Mining Massive Datasets - Graphs Deep Generative Models.", Prof. Dr. Stephan Günnemann et.al. TUM, 2019.*

# Wasserstein Generative Adversarial Networks

---

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values  $\alpha = 0.00005$ ,  $c = 0.01$ ,  $m = 64$ ,  $n_{\text{critic}} = 5$ .

---

**Require:** :  $\alpha$ , the learning rate.  $c$ , the clipping parameter.  $m$ , the batch size.  $n_{\text{critic}}$ , the number of iterations of the critic per generator iteration.

**Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters.

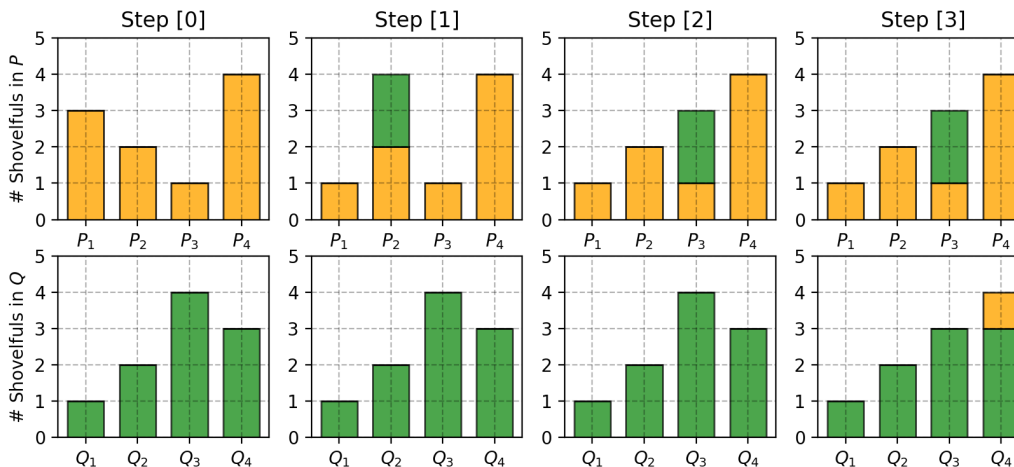
```

1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSPProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSPProp}(\theta, g_\theta)$ 
12: end while
    
```

---

# Wasserstein distance

Wasserstein Distance is a measure of the distance between two probability distributions. It is also called Earth Mover's distance. It can be interpreted as the minimum energy cost of moving and transforming a pile of dirt in the shape of one probability distribution to the shape of the other distribution.



$$P_1 = 3, P_2 = 2, P_3 = 1, P_4 = 4$$
$$Q_1 = 1, Q_2 = 2, Q_3 = 4, Q_4 = 3$$

If we label the cost to pay to make  $P_i$  and  $Q_i$  match  $\delta_i$ , we would have  $\delta_{i+1} = \delta_i + P_i - Q_i$  and in example:

$$\begin{aligned}\delta_0 &= 0 \\ \delta_1 &= 0 + 3 - 1 = 2 \\ \delta_2 &= 2 + 2 - 2 = 2 \\ \delta_3 &= 2 + 1 - 4 = -1 \\ \delta_4 &= -1 + 4 - 3 = 0\end{aligned}\tag{1}$$

When dealing with the continuous probability domain, the distance formula becomes:

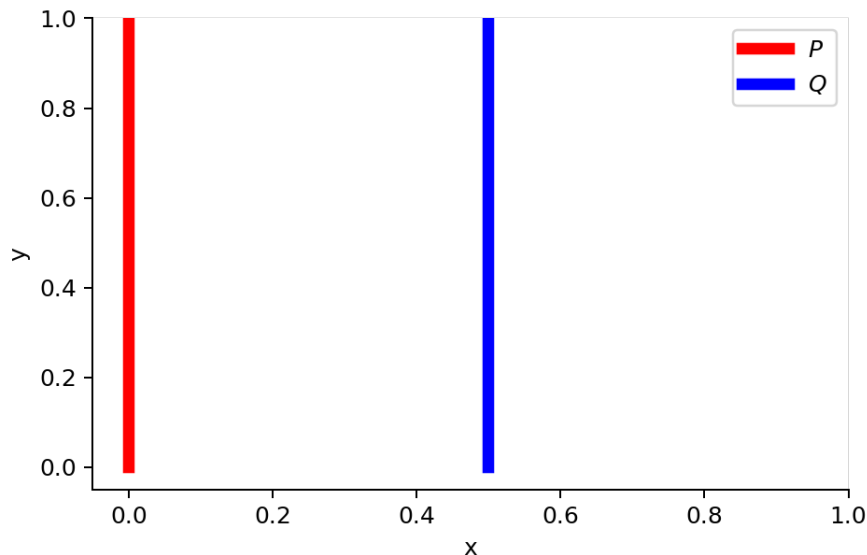
$$W(p^*(x), p(z)) = \inf_{\gamma \sim \Pi(p^*(x), p(z))} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]\tag{2}$$

*Image source: <https://arxiv.org/pdf/1904.08994.pdf>.*

# Why Wasserstein is better than JS or KL divergence?

Suppose we have two probability distributions,  $P$  and  $Q$ :

$$\begin{aligned} \forall (x, y) \in P, x = 0 \text{ and } y \sim U(0, 1) \\ \forall (x, y) \in Q, x = \theta, 0 \leq \theta \leq 1 \text{ and } y \sim U(0, 1) \end{aligned} \quad (3)$$



$$\begin{aligned}D_{KL}(P\|Q) &= \sum_{x=0, y \sim U(0,1)} 1 \cdot \log \frac{1}{0} = +\infty \\D_{KL}(Q\|P) &= \sum_{x=\theta, y \sim U(0,1)} 1 \cdot \log \frac{1}{0} = +\infty \\D_{JS}(P, Q) &= \frac{1}{2} \left( \sum_{x=0, y \sim U(0,1)} 1 \cdot \log \frac{1}{1/2} + \sum_{x=0, y \sim U(0,1)} 1 \cdot \log \frac{1}{1/2} \right) = \log 2 \\W(P, Q) &= |\theta|\end{aligned}\tag{4}$$

But when  $\theta = 0$ , two distributions are fully overlapped:

$$\begin{aligned}D_{KL}(P\|Q) &= D_{KL}(Q\|P) = D_{JS}(P, Q) = 0 \\W(P, Q) &= 0 = |\theta|\end{aligned}\tag{5}$$

From "<https://arxiv.org/pdf/1904.08994.pdf>".

# The differences in implementation for the WGAN

- 1 Use a linear activation function in the output layer of the critic model (instead of sigmoid).
- 2 Use -1 labels for real images and 1 labels for fake images (instead of 1 and 0).
- 3 Use Wasserstein loss to train the critic and generator models.
- 4 Constrain critic model weights to a limited range after each mini batch update (e.g.  $[-0.01, 0.01]$ ).
- 5 Update the critic model more times than the generator each iteration (e.g. 5).
- 6 Use the RMSProp version of gradient descent with a small learning rate and no momentum (e.g. 0.00005).

*From "<https://arxiv.org/pdf/1904.08994.pdf>".*