

# Bài 44 - Model Wasserstein GAN (WGAN)

25 Jul 2020 - phamdinhhkhanh

## Menu

- 1. Giới thiệu chung
- 2. Cross Entropy
- 3. Kullback Leibler divergence
  - 3.1. Ví dụ về độ đo Kullback Leibler divergence
- 4. Jensen-Shannon
  - 4.1. Hội tụ của GAN và khoảng cách Jensen-Shannon
- 5. Độ đo Earth-Mover (EM) hoặc Wasserstein
  - 5.1. Ví dụ về Wasserstein
- 6. Hạn chế của DCGAN
- 7. Wasserstein GAN
  - 7.1. Khái niệm liên tục Lipschitz (Lipschitz continuity)
  - 7.2. Wasserstein GAN
- 8. Thực hành Wasserstein GAN
  - 8.1. Load dữ liệu
  - 8.2. Discriminator
  - 8.3. Generator
  - 8.4. WGAN-GP model
  - 8.5. GAN monitoring
  - 8.6. Huấn luyện
- 9. Tổng kết
- 10. Tài liệu

## 1. Giới thiệu chung

Huấn luyện model GAN là một việc không dễ dàng bởi trong quá trình huấn luyện model thường xuyên không thể hội tụ đến điểm cân bằng. Một trong những cách chúng ta có thể áp dụng để khắc phục sự không hội tụ đó là thay đổi hàm loss function. Trong bài viết này chúng ta sẽ cùng tìm hiểu phương pháp huấn luyện Wasserstein GAN (<https://arxiv.org/abs/1701.07875>) và phạt WGAN-Gradient. Nếu đọc bài báo gốc bạn sẽ thấy hàm loss function của Wasserstein GAN có một chút khó hiểu bởi các lý thuyết khá hàn lâm và không dễ tiếp cận với người bắt đầu, nhưng trên thực tế việc áp dụng phương pháp này lại khá đơn giản. Nhằm đưa ra một lời giải thích một cách dễ hiểu kèm theo code mẫu, bài viết này mình sẽ làm rõ hơn về phương pháp Wasserstein GAN thông qua cách tiếp cận nắm vững lý thuyết. Qua đó bạn đọc sẽ hiểu được đặc trưng của các dạng hàm loss function khác nhau, ưu nhược điểm và trường hợp sử dụng của chúng. Đây sẽ là những kiến thức cơ bản nhưng là tiền đề quan trọng để chúng ta đi xa hơn để hiểu những nghiên cứu của thuật toán GAN SOTA và thậm chí là những lớp mô hình deep learning khác.

Nhưng trước khi bắt đầu bài viết này các bạn hãy ôn lại Bài 43 - Model GAN (<https://phamdinhhkhanh.github.io/2020/07/13/GAN.html>) để hiểu rõ kiến trúc và nguyên lý hoạt động của GAN. Sau khi đã hiểu GAN hoạt động như thế nào? Cách thức xây dựng và huấn luyện một mô hình GAN ra sao thì chúng ta sẽ quay lại đọc tiếp bài này. WGAN sẽ tập trung vào việc thay đổi hàm loss function. Do đó, đầu tiên chúng ta sẽ cùng phân tích các dạng loss function được sử dụng chung trong các model GAN và đơn giản nhất là hàm cross entropy bên dưới.

## 2. Cross Entropy

Chúng ta đã biết trong hồi qui logistic và softmax sử dụng loss function là hàm *cross entropy* để đo lường **tương quan** giữa 2 phân phối xác suất thực tế  $\mathbf{p} = (p_1, \dots, p_C)$  và phân phối xác suất dự báo  $\mathbf{q} = (q_1, \dots, q_C)$  như sau:

$$D_{CE}(\mathbf{p} \parallel \mathbf{q}) \triangleq \mathbf{H}(\mathbf{p}, \mathbf{q}) = - \sum_{i=1}^C p_i \log(q_i)$$

Trong đó  $\sum_{i=1}^C p_i = \sum_{i=1}^C q_i = 1$  và  $C$  là số lượng classes.  $\mathbf{H}(\mathbf{p}, \mathbf{q})$  chính là ký hiệu của hàm *cross entropy* của hai phân phối  $\mathbf{p}$  và  $\mathbf{q}$ .

Mục tiêu tối ưu hóa hàm loss function cũng đồng nghĩa với mục tiêu tìm ra mô hình dự báo phân phối  $\mathbf{q}$  *sát nhất* với phân phối  $\mathbf{p}$ . Hai nói cách khác, khi phân phối  $\mathbf{p}$  và  $\mathbf{q}$  càng tương quan thì giá trị hàm *cross entropy* càng nhỏ và mô hình dự báo càng chuẩn xác. Trái lại khi  $\mathbf{p}$  và  $\mathbf{q}$  không tương quan thì giá trị của *cross entropy* càng lớn.

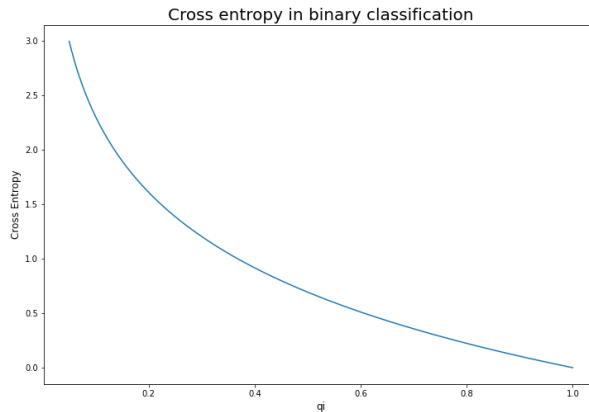
Thật vậy, để dễ hình dung chúng ta lấy ví dụ trong trường hợp phân loại nhị phân, nhãn  $p_i$  sẽ nhận một trong hai giá trị 0, 1. Ta có:  $D_{CE}(\mathbf{p} \parallel \mathbf{q}) = -\log(q_i)$  (trường hợp  $p_i = 0$  thì đóng góp vào cross entropy bằng 0, trường hợp  $p_i = 1$  ta thu được công thức như trên). Khi đó đồ thị của hàm *cross entropy* biểu diễn theo  $q_i$  như sau:

Top

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  q = np.linspace(0.05, 1, 200)
5  d_ce = -np.log(q)
6
7  plt.figure(figsize=(12, 8))
8  plt.plot(q, d_ce)
9  plt.ylabel('Cross Entropy', fontsize=12)
10 plt.xlabel('qi', fontsize=12)
11 plt.title('Cross entropy in binary classification', fontsize=20)
12 plt.show()

```



Ta có thể thấy giá trị nhỏ nhất của hàm *cross entropy* đạt được khi  $q_i = 1$ . Tức là phân phối xác suất của  $\mathbf{p}$  và  $\mathbf{q}$  là trùng nhau. Trong trường hợp tổng quát cho nhiều nhãn ta cũng có thể chứng minh được sự **tương quan** giữa  $\mathbf{p}$  và  $\mathbf{q}$  tại điểm cực tiểu của *cross entropy*. Bài toán tối ưu *cross entropy* có dạng như sau:

$$(\mathbf{p}^*, \mathbf{q}^*) = \arg \min_{(\mathbf{p}, \mathbf{q})} \mathbf{H}(\mathbf{p}, \mathbf{q}) = \arg \min_{(\mathbf{p}, \mathbf{q})} - \sum_{i=1}^C p_i \log(q_i)$$

Thỏa mãn ĐK:  $\sum_{i=1}^C p_i = \sum_{i=1}^C q_i = 1$ .

Đây là bài toán tối ưu lồi với điều kiện ràng buộc tuyến tính nên các giải quen thuộc là thông qua tối ưu hàm lagrange:

$$\mathcal{L}(\mathbf{p}, \mathbf{q}, \lambda) = - \sum_{i=1}^C p_i \log(q_i) + \lambda \left( \sum_{i=1}^C q_i - 1 \right)$$

Nếu coi phân phối  $\mathbf{p}$  là cố định và tìm phân phối  $\mathbf{q}$  theo  $\mathbf{p}$  để tối thiểu hóa hàm lagrange. Các phương trình đạo hàm bậc nhất của hàm lagrange như sau:

$$\begin{cases} \nabla_{q_i} \mathcal{L}(\mathbf{p}, \mathbf{q}, \lambda) &= -\frac{p_i}{q_i} + \lambda, \forall i = \overline{1, C} \\ \nabla_{\lambda} \mathcal{L}(\mathbf{p}, \mathbf{q}, \lambda) &= \sum_{i=1}^C q_i - 1 \end{cases}$$

Điều kiện cần để bài toán có cực trị đó là đạo hàm bậc nhất có nghiệm bằng 0. Tức là:  $\frac{p_1}{q_1} = \frac{p_2}{q_2} = \dots = \frac{p_C}{q_C} = \lambda$ . Suy ra cực trị đạt được khi  $\mathbf{p}$  và  $\mathbf{q}$  có cùng phân phối xác suất.

Mặc khác  $\nabla_{q_i}^2 \mathcal{L}(\mathbf{p}, \mathbf{q}, \lambda) = \frac{p_i}{q_i^2} \geq 0, \forall i = \overline{1, C}$  nên hàm mục tiêu là một hàm lồi. Do đó đây chính là điều kiện đủ để cực trị  $(\mathbf{p}^*, \mathbf{q}^*)$  đạt được chính là điểm cực tiểu.

Chứng minh bài toán không quá khó phải không nào?

Như vậy tìm cực tiểu của *cross entropy* chính là tìm các giá trị để phân phối dự báo sát với phân phối thực tế nhất. Đó cũng là lý do *cross entropy* thường được lựa chọn làm hàm mục tiêu trong các bài toán liên quan tới dự báo phân phối xác suất.

### 3. Kullback Leibler divergence

*cross entropy* là một độ đo được sử dụng để đo lường tương quan giữa hai phân phối xác suất. Tuy nhiên *cross entropy* có một số hạn chế đó là giá trị của nó có thể lớn tùy ý. Chẳng hạn như khi  $p_i = q_i = \frac{1}{n}, i \in \overline{1, n}$ . Khi  $n \rightarrow +\infty$  thì giá trị của *cross entropy* sẽ là :

$$\lim_{n \rightarrow +\infty} D_{CE}(\mathbf{p} \parallel \mathbf{q}) = \lim_{n \rightarrow +\infty} - \sum_{i=1}^n \frac{1}{n} \log\left(\frac{1}{n}\right) = \lim_{n \rightarrow +\infty} \log(n) = +\infty$$

Do đó giá trị của *cross entropy* có thể đưa ra giá trị *tương quan* giữa hai phân phối  $\mathbf{p}$  và  $\mathbf{q}$  là một số rất lớn mặc dù chúng thể hiện cùng một phân phối xác suất. Tồn tại một thước đo khác giúp tìm ra khác biệt giữa hai phân phối xác suất  $P(x), Q(x)$  và trả về giá trị là 0 trong trường hợp hai phân phối này là trùng nhau. Đó chính là độ đo *Kullback leiber divergence* như sau:

Top

$$\begin{aligned}
 D_{KL}(P||Q) &= \sum_{i=1}^n P(x) \log\left(\frac{P(x)}{Q(x)}\right) \\
 &= \sum_{i=1}^n P(x) [\log(P(x)) - \log(Q(x))]
 \end{aligned} \tag{1}$$

Công thức trên là xét với trường hợp phân phối của  $x$  rời rạc. Nếu phân phối của  $x$  liên tục thì tổng sẽ được chuyển sang tích phân:

$$D_{KL}(p||q) = \int_{\chi} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

Với  $\chi$  là miền xác định của  $x$ . Các ký hiệu chữ thường  $p, q$  đại diện cho *hàm mật độ xác suất (pdf)* của biến liên tục  $x$ . Nếu bạn vẫn chưa biết *hàm mật độ xác suất* là gì hãy xem lại Apenddix 1 - Lý thuyết phân phối và kiểm định thống kê ([https://phamdinhhkhanh.github.io/2019/05/10/Hypothesis\\_Statistic.html#21-kh%C3%A1i-ni%E1%BB%87m-bi%E1%BA%BFn-ng%E1%BA%ABu-nhi%C3%AAn](https://phamdinhhkhanh.github.io/2019/05/10/Hypothesis_Statistic.html#21-kh%C3%A1i-ni%E1%BB%87m-bi%E1%BA%BFn-ng%E1%BA%ABu-nhi%C3%AAn)).

Hàm  $D_{KL}(p||q)$  rõ ràng không xác định tại một số trường hợp đặc biệt khi  $p(x) = 0$  hoặc  $q(x) = 0$ . Khi đó chúng ta phải qui ước giá trị của nó tại những điểm đặc biệt này. Cụ thể:

Tại  $p(x) = 0$  thì  $p(x) \log(p(x)) = 0$ , nó không có đóng góp gì vào độ lớn của hàm số  $D_{KL}(p||q)$ .

Trường hợp  $q(x) = 0$ ,  $p(x)$  hữu hạn thì ta coi  $\log(q(x)) \rightarrow -\infty$ , suy ra:  $D_{KL}(p||q) = +\infty$  Ngoài ra ta nhận thấy một tính chất đặc biệt đó là phương trình (1) tương đương với:

$$D_{KL}(P||Q) = \mathbf{H}(P, Q) - \mathbf{H}(P)$$

Trong đó  $\mathbf{H}(P) = -P(x) \log P(x)$  chính là hàm *entropy* đo lường giá trị trung bình của  $\log(P(x))$  theo phân phối  $P(x)$ . Đây là hàm số thường được sử dụng để đánh giá mức độ tinh khiết của chỉ một phân phối xác suất (tinh khiết được hiểu là output trả ra chỉ thuộc về một nhóm) và thường được dùng trong các thuật toán về decision tree.

Hàm *cross entropy*  $\mathbf{H}(P, Q)$  đo lường giá trị trung bình của  $\log(Q(x))$  theo phân phối  $P(x)$ .

Như vậy, khoảng cách  $D_{KL}(P||Q)$  là một thước đo giá trị khác biệt giữa hai phân phối  $P(x)$  và  $Q(x)$ . Giá trị của  $D_{KL}$  sẽ bằng 0 khi  $P(x)$  và  $Q(x)$  có phân phối trùng nhau. Ta chứng minh điều này như sau:

**Tìm cực trị của  $D_{KL}(P||Q)$**

Giả sử giá trị của hai phân phối là:  $P = (p_1, \dots, p_C)$  và  $Q = (q_1, \dots, q_C)$ .

Trong trường hợp tồn tại  $q_j = 0$  thì theo như qui ước ở (2) ta có  $D_{KL}(p||q) = +\infty$ .

Xét hàm số:  $f(x) = -\log(x)$

Trên miền xác định  $x > 0$ . Ta có :

$$\begin{cases} \nabla_x f(x) &= -\frac{1}{x} \\ \nabla_x^2 f(x) &= \frac{1}{x^2} > 0 \end{cases}$$

Như vậy  $f(x)$  là một hàm lồi ngặt trên toàn miền  $x \in (0, 1]$ . Do đó nó thỏa mãn bất đẳng thức Jensen ([https://en.wikipedia.org/wiki/Jensen%27s\\_inequality](https://en.wikipedia.org/wiki/Jensen%27s_inequality)) đối với hàm lồi:

$$\lambda f(x_1) + (1 - \lambda)f(x_2) \geq f(\lambda x_1 + (1 - \lambda)x_2)$$

với mọi  $\lambda \in (0, 1]$

Trong trường hợp tổng quát ta có:

$$\sum_{i=1}^n \lambda_i f(x_i) \geq f\left(\sum_{i=1}^n \lambda_i x_i\right)$$

Với  $\sum_{i=1}^n \lambda_i = 1$  và  $\lambda_i \in (0, 1]$ ,  $\forall i = \overline{1, n}$

Áp dụng bất đẳng thức đối với hàm lồi ta có :

$$\begin{aligned}
 D_{KL}(P||Q) &= -\sum_{i=1}^C p_i \log\left(\frac{q_i}{p_i}\right) \\
 &= \sum_{i=1}^C p_i f\left(\frac{q_i}{p_i}\right) \geq f\left(\sum_{i=1}^C p_i \frac{q_i}{p_i}\right) \\
 &= \log\left(\sum_{i=1}^C q_i\right) = \log(1) = 0
 \end{aligned}$$

Như vậy giá trị của  $D_{KL}(P||Q)$  nhỏ nhất bằng 0 trên toàn miền  $(0, 1]$ . Kết hợp với trường hợp tồn tại  $q_i = 0$  thì  $D_{KL}(p||q) = +\infty$  ta suy ra giá trị nhỏ nhất là 0 đạt được trên toàn miền  $[0, 1]$ . Vì  $f(x)$  là một hàm lồi ngặt (*strictly convex*) nên đẳng thức đạt được khi:

$$\frac{q_1}{p_1} = \frac{q_2}{p_2} = \dots = \frac{q_C}{p_C} = \lambda$$

Tức là phân phối  $P$  và  $Q$  là đồng nhất.

Như vậy *Kullback Leibler divergence* là một hàm số sát hơn đo lường khoảng giữa hai phân phối. So với cross entropy thì giá trị của nó nhỏ hơn và đặc biệt khi hai phân phối trùng nhau thì khoảng cách *Kullback Leibler divergence* đạt giá trị cực tiểu và bằng 0. Top

### 3.1. Ví dụ về độ đo Kullback Leibler divergence

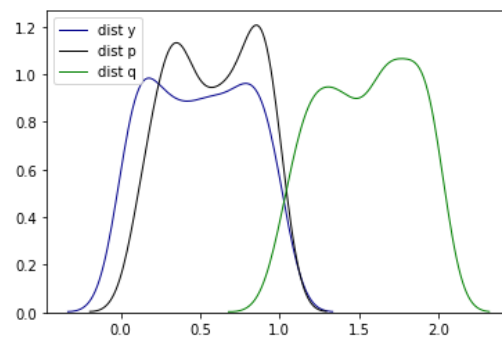
Để hiểu rõ hơn về đặc tính đo lường sự khác biệt giữa hai phân phối xác suất của *Kullback Leibler divergence* chúng ta cùng lấy ví dụ về ba phân phối xác suất  $p$ ,  $q$  và  $y$  cùng được tạo ra từ phân phối chuẩn theo công thức:

$$\text{pdf}(x) = \frac{\exp \frac{-(x-\mu)^2}{2\sigma^2}}{\sqrt{2\pi\sigma^2}}$$

Trong đó phân phối của  $y \sim \mathbf{N}(0, 1)$ ,  $p \sim \mathbf{N}(0.1, 1)$  và  $q \sim \mathbf{N}(2, 1)$ .

Ta nhận thấy cả ba phân phối này có cùng phương sai, tuy nhiên tâm của phân phối chính là giá trị kỳ vọng của các phân phối của  $p$  lại gần với  $y$  hơn so với  $q$  nên khả năng cao phân phối  $p$  sẽ giống  $y$  hơn phân phối  $q$ . Chúng ta cùng kiểm chứng thông qua mô phỏng các phân phối này:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 y = np.random.uniform(0, 1, 200)
5 p = np.random.uniform(0.1, 1, 200)
6 q = np.random.uniform(2, 1, 200)
7
8
9
10 sns.distplot(y,
11              hist = False,
12              kde = True,
13              color = 'darkblue',
14              kde_kws={'linewidth':1},
15              label = 'dist y'
16              )
17
18 sns.distplot(p,
19              hist = False,
20              kde = True,
21              color = 'black',
22              kde_kws={'linewidth':1},
23              label = 'dist p'
24              )
25
26 plt.show()
```



Đồng thời áp dụng công thức (1) để tính khoảng cách giữa các phân phối theo độ đo Kullback Leibler ta có:

```
1 def kl(p, q):
2     return np.sum(p*np.log(p/q))
3
4 kl_py = kl(p, y)
5 kl_qy = kl(q, y)
6
7 print('Kullback leibler (p, y): ', kl_py)
8 print('Kullback leibler (q, y): ', kl_qy)
9
10
11 Kullback leibler (p, y):  69.72425099387095
12 Kullback leibler (q, y):  459.4113963613252
```

Kết quả cho thấy khoảng cách  $D_{KL}(p||y) \leq D_{KL}(q||y)$ . Kết quả này cho thấy nhận định về độ đo *Kullback Leibler* đo lường khoảng cách giữa 2 phân phối là hoàn toàn phù hợp.

Mặc dù *Kullback Leibler* trông có vẻ khá hoàn hảo nhưng nó cũng tồn tại những hạn chế khiến cho nó không phù hợp với một phép đo khoảng cách đó là không có tính chất đối xứng và giá trị không bị chặn khiến cho chúng ta phải tìm đến một độ đo mới khắc phục được cả hai nhược điểm này. Đó chính là Jensen-Shannon.

## 4. Jensen-Shannon

Jensen-Shannon là độ đo được xây dựng dựa trên Kullback Leibler

$$D_{JS}(P\|Q) = D_{KL}(P\|\frac{P+Q}{2}) + D_{KL}(Q\|\frac{P+Q}{2})$$

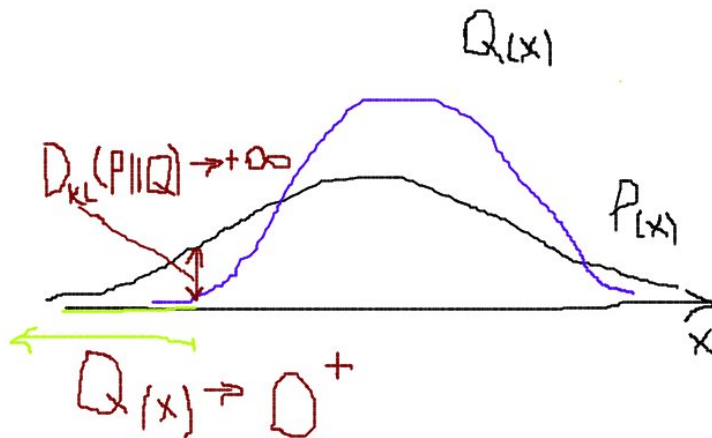
Độ đo Jensen-Shannon thay vì đo khoảng cách trực tiếp giữa  $P$  và  $Q$  thì nó đo khoảng cách này thông qua khoảng cách giữa  $P$  với  $\frac{P+Q}{2}$  và giữa  $Q$  với  $\frac{P+Q}{2}$ . Ở đây  $\frac{P+Q}{2}$  là phân phối trung bình của hai phân phối  $P$  và  $Q$ . Nếu  $P$  và  $Q$  càng giống nhau thì khoảng cách của nó tới phân phối trung bình càng nhỏ và khi chúng trùng nhau thì khoảng cách này bằng 0.

Một tính chất quan trọng của Jensen-Shannon so với Kullback Leibler đó là nó có tính chất đối xứng. Nghĩa là:  $D_{JS}(P\|Q) = D_{JS}(Q\|P)$ . Thật vậy, khai triển Jensen-Shannon ta có:

$$\begin{aligned} D_{JS}(P\|Q) &= D_{KL}(P\|\frac{P+Q}{2}) + D_{KL}(Q\|\frac{P+Q}{2}) \\ &= \sum_{i=1}^C P(x) \log(P(x)) - \sum_{i=1}^C P(x) \log(\frac{P(x)+Q(x)}{2}) + \\ &\quad \sum_{i=1}^C Q(x) \log(Q(x)) - \sum_{i=1}^C Q(x) \log(\frac{P(x)+Q(x)}{2}) \\ &= - \sum_{i=1}^C (P(x) + Q(x)) \log(\frac{P(x)+Q(x)}{2}) + \sum_{i=1}^C P(x) \log(P(x)) + \sum_{i=1}^C Q(x) \log(Q(x)) \\ &= 2 \cdot \mathbf{H}(\frac{P+Q}{2}) - \mathbf{H}(P) - \mathbf{H}(Q) \end{aligned}$$

Jensen-Shannon sẽ bằng hiệu giữa *entropy* của phân phối  $\frac{P+Q}{2}$  với các *entropy* của  $P$  và  $Q$ . Do đó nó có tính chất đối xứng và độ đo bằng 0 khi  $P = Q$ .

Từ công thức (1) ta nhận thấy độ đo Kullback-Leiber sẽ rất lớn tại những điểm  $x$  mà giá trị của  $Q(x)$  rất nhỏ và giá trị của  $P(x)$  khác 0. Bởi vì khi đó  $P(x)[\log(P(x)) - \log(Q(x))]\rightarrow +\infty$ . Do đó hàm loss function sẽ thường không hội tụ nếu xuất hiện những điểm  $x$  mà phân phối của  $Q(x) \rightarrow 0^+$ .



**Hình 1:** Giá trị của Kullback-Leibler tại những điểm có  $Q(x) \rightarrow 0^+$  sẽ khiến cho  $D_{KL}(P\|Q) \rightarrow +\infty$ .

Độ đo Jensen-Shannon sẽ khắc phục được hạn chế này vì giá trị của  $\log(\frac{P(x)}{P(x)+Q(x)})$  luôn giao động trong khoảng:

$$\log(\frac{P(x)}{P(x)+1}) \leq \log(\frac{P(x)}{P(x)+Q(x)}) \leq 1$$

Do đó

$$D_{KL}(P\|\frac{P+Q}{2}) = P(x) \log(\frac{P(x)}{P(x)+Q(x)}) \leq P(x)$$

Điều đó đảm bảo đóng góp vào khoảng cách  $D_{JS}(P\|Q)$  tại mọi điểm  $x$  luôn bị giới hạn trên và không làm giá trị khoảng cách trở nên quá lớn.

Như vậy ưu điểm chính của độ đo Jensen-Shannon so với Kullback-Leibler đó là có tính chất đối xứng và có giá trị hữu hạn. Do đó nó phù hợp để trở thành một độ đo thay thế cho Kullback-Leibler trong việc tìm khoảng cách giữa hai phân phối xác suất.

## 4.1. Hội tụ của GAN và khoảng cách Jensen-Shannon

Ta sẽ chứng minh rằng hội tụ của GAN sẽ chính là hội khoảng cách Jensen-Shannon. Ở bài 43

(<https://phamdinhhkhanh.github.io/2020/07/13/GAN.html>) ta đã biết hàm loss function của GAN là kết hợp giữa generator và discriminator có dạng như sau :

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_r(x)} [\log(D(x))] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

Ở đây  $x \sim p_r(x)$  là phân phối của dữ liệu thật và  $z \sim p_z(z)$  là phân phối của dữ liệu sinh.

Khi mô hình GAN hội tụ thì dữ liệu sinh ra giống với dữ liệu thật nhất. Giả sử nghiệm hội tụ generator của GAN là hàm  $G^*$ . Khi đó  $G^*(z) \rightarrow x$  và giá trị hội tụ của loss function sẽ là:

$$\max_D V(D, G^*) = \mathbb{E}_{x \sim p_r(x)} [\log(D(x))] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))]$$

Ở đây,  $x \sim p_g(x)$  là phân phối của dữ liệu từ mô hình generator. Đầu vào của generator là một véc tơ noise ngẫu nhiên nên có thể coi giá trị của nó là một hàm liên tục. Mặt khác khi  $x$  liên tục ta có tính chất :

$$\mathbb{E}_{x \sim p(x)} f(x) = \int_x p(x) f(x) dx$$

Do đó:

$$\begin{aligned} \max_D V(D, G^*) &= \mathbb{E}_{x \sim p_r(x)} [\log(D(x))] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))] \\ &= \int_x p_r(x) \log(D(x)) dx + \int_x p_g(x) \log(1 - D(x)) dx \\ &= \int_x p_r(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx \end{aligned}$$

Mục tiêu của chúng ta là tìm kiếm discriminator để tối đa hóa giá trị của hàm loss function nên  $D(x)$  được xem như biến của bài toán cực trị. Đặt  $D(x) = y$ .

Xét hàm  $V(y) = p_r(x) \log(y) + p_g(x) \log(1 - y)$

Hàm số trên có đạo hàm bậc nhất và bậc 2 :

$$\begin{cases} \nabla_y V(y) &= \frac{p_r(x)}{y} - \frac{p_g(x)}{1-y} \\ \nabla_y^2 V(y) &= -\frac{p_r(x)}{y^2} - \frac{p_g(x)}{(1-y)^2} < 0 \end{cases}$$

Như vậy  $V(y)$  là một hàm lõm. Do đó giá trị cực đại của nó chính là nghiệm của phương trình đạo hàm bậc nhất :

$$\begin{aligned} \frac{p_r(x)}{y^*} - \frac{p_g(x)}{1-y^*} &= 0 \\ \Leftrightarrow y^* &= \frac{p_r(x)}{p_r(x) + p_g(x)} \end{aligned}$$

Ta có :

$$\begin{aligned} V(y) &\leq V(y^*) \\ \Leftrightarrow \int_x p_r(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx &\leq \int_x p_r(x) \log\left(\frac{p_r(x)}{p_r(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_r(x) + p_g(x)}\right) dx \end{aligned}$$

Suy ra giá trị hội tụ:

$$\begin{aligned} V(D^*, G^*) &= \int_x p_r(x) \log\left(\frac{p_r(x)}{p_r(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{p_g(x)}{p_r(x) + p_g(x)}\right) dx \\ &= \int_x p_r(x) \log\left(\frac{2 p_r(x)}{p_r(x) + p_g(x)}\right) + p_g(x) \log\left(\frac{2 p_g(x)}{p_r(x) + p_g(x)}\right) dx \\ &\quad - \log(2) \int_x p_r(x) - \log(2) \int_x p_g(x) dx \\ &= D_{KL}(p_r(x) || \frac{p_g(x) + p_r(x)}{2}) + D_{KL}(p_g(x) || \frac{p_g(x) + p_r(x)}{2}) - 2 \log(2) \\ &= D_{JS}(p_r(x) || p_g(x)) - 2 \log(2) \end{aligned}$$

Tại điểm hội tụ, hàm loss function của GAN bằng giá trị khoảng cách Jensen-Shannon giữa phân phối của dữ liệu thật  $p_r(x)$  và dữ liệu sinh  $p_g(x)$  trừ đi  $2 \log(2)$ .

Mặt khác khi ảnh thật và sinh rất giống nhau thì  $D_{JS}(p_r(x) || p_g(x)) = 0$  và giá trị loss function của GAN hội tụ về  $-2 \log(2)$ . Và đồng thời nghiệm hội tụ của hai phân phối  $p_r, p_g$  của GAN trùng với Jensen-Shannon.

## 5. Độ đo Earth-Mover (EM) hoặc Wasserstein

Top

Giả sử  $\mathbf{p}_r$  là phân phối thực tế của dữ liệu gốc được gọi tắt là *phân phối thực* và  $\mathbf{p}_g$  là *phân phối sinh*, tức là phân phối của dữ liệu được sinh ra từ model GAN. Wasserstein hoặc EM là giá trị *cận dưới nhỏ nhất* (infimum) của khoảng cách giữa phép dịch chuyển  $\mathbf{p}_r$  sang  $\mathbf{p}_g$  theo công thức:

$$W(\mathbf{p}_r, \mathbf{p}_g) = \inf_{\gamma \in \Pi(\mathbf{p}_r, \mathbf{p}_g)} \mathbf{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

Giải thích một chút về công thức trên:

$\Pi(\mathbf{p}_r, \mathbf{p}_g)$  là *phân phối đồng thời* (join distribution) của hai phân phối  $\mathbf{p}_r$  và  $\mathbf{p}_g$  và mỗi một phân phối  $\mathbf{p}_r$  và  $\mathbf{p}_g$  được coi như là *phân phối biên* (margin distribution). Bạn đọc có thể tìm hiểu thêm về phân phối đồng thời và phân phối biên tại Appendix 1 - Lý thuyết phân phối và kiểm định thống kê ([https://phamdinhhkhanh.github.io/2019/05/10/Hypothesis\\_Statistic.html#22-ph%C3%A2n-ph%E1%BB%91i-x%C3%A1c-su%E1%BA%A5t-C4%91%E1%BB%93ng-th%E1%BB%9Di](https://phamdinhhkhanh.github.io/2019/05/10/Hypothesis_Statistic.html#22-ph%C3%A2n-ph%E1%BB%91i-x%C3%A1c-su%E1%BA%A5t-C4%91%E1%BB%93ng-th%E1%BB%9Di)).  $\gamma(x, y)$  chính là giá trị của phân phối đồng thời của  $(x, y)$ . Khoảng cách *wasserstein* là cận dưới lớn nhất của kỳ vọng khoảng cách giữa các điểm  $(x, y)$  theo phân phối chung  $\gamma$ . Khoảng cách được đo lường theo độ đo  $\|x - y\|$ .

Để dễ hình dung hơn về Wasserstein chúng ta tìm hiểu ví dụ bên dưới:

## 5.1. Ví dụ về Wasserstein

Giả sử dữ liệu gốc của chúng ta gồm các giá trị 1, 2, 3, 4, 5, 6 có phân phối nằm trên 3 cột 1, 2, 3. Từ dữ liệu gốc chúng ta có các phương án dịch chuyển các giá trị này sang một phân phối mới nằm trên 4 cột 7, 8, 9, 10 và thu được bảng giá trị phân phối xác suất đồng thời như bảng  $\gamma_1$  hình bên dưới:



**Hình 2:** Ví dụ về độ đo Wasserstein biến đổi từ phân phối thật  $\mathbf{p}_r$  (cột 1, 2, 3) bên trái sang phân phối sinh  $\mathbf{p}_g$  (cột 7, 8, 9, 10) ở giữa. Bên phải là bảng giá trị của phân phối đồng thời. Source WGAN - jonathan hui ([https://medium.com/@jonathan\\_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490](https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490)).

Xét phép biến đổi  $\gamma_1$  ta thấy ô có giá trị 1 ở đã được dịch chuyển từ cột 1 của phân phối thật sang cột 7 của phân phối sinh. Do đó khoảng cách  $\gamma(x, y) = \|x - y\| = 7 - 1 = 6$ . Tương tự ô giá trị 4 được dịch chuyển từ cột 2 sang 8 nên có giá khoảng cách là 6. Phương án dịch chuyển từ cột 1 sang cột 10 tương ứng với 2 ô giá trị 3 và 2 nên trên bảng phân phối xác suất đồng thời giá trị tương ứng với dòng 1, cột 10 được điền là 2. Khoảng cách Wasserstein bằng các khoảng cách cần dịch chuyển để biến đổi phân phối thật thành phân phối sinh nhân với giá trị *phân phối xác suất đồng thời*. Kết quả ta thu được:  $(6 + 6 + 6 + 6 + 2 \times 9) = 42$ .

Khi hai phân phối thật và phân phối sinh trùng nhau thì khoảng cách này bằng 0. Khoảng cách càng nhỏ thì 2 phân phối càng gần nhau và dữ liệu được sinh ra càng giống với dữ liệu thật.

## 6. Hạn chế của DCGAN

### Phân biệt giữa GAN và DCGAN

Mô hình GAN là tên gọi chung cho những mô hình sử dụng generator model để sinh ra dữ liệu giả và discriminator để phân biệt giữa dữ liệu thật và giả. Ở kiến trúc cơ bản đầu tiên của GAN chỉ sử dụng những fully connected layer đơn thuần.

DCGAN là *Deep Convolution GAN*, một kiến trúc cũng tương tự như GAN nhưng tập trung vào các mạng học sâu trong xử lý ảnh. Ngoài các fully connected layer của GAN thì trong kiến trúc này chúng ta sử dụng thêm các layer tích chập (convolutional layer) để trích lọc đặc trưng của ảnh. DCGAN có lẽ phù hợp hơn với dữ liệu đầu vào là hình ảnh/video, trong khi đó mô hình GAN tổng quát có thể áp dụng ở domain rộng hơn.

### Hạn chế của DCGAN

Trong một bài báo vào năm 2017, Arjovsky đã chỉ ra nhược điểm của DCGAN như sau:

- Quá trình huấn luyện discriminator sẽ tạo ra những thông tin tốt để cải thiện generator nhưng generator lại không thực hiện tốt vai trò của mình khiến cho giá trị gradient descent của nó thường bị triệt tiêu và generator không học được gì.

$$-\nabla_{\theta_g} \log(1 - D(G(z^{(i)}))) \rightarrow 0$$

- Một hàm loss function mới được Arjovsky đề xuất để thay thế hàm loss function cũ của generator nhằm giải quyết tình trạng triệt tiêu đạo hàm là:

$$\nabla_{\theta_g} \log(D(G(z^{(i)})))$$

Tuy nhiên đạo hàm này sẽ có một số giá trị rất lớn và dẫn tới mô hình hoạt động không ổn định.

Arjovsky đề xuất một phương pháp mới sẽ add thêm noise vào hình ảnh thật trước khi đưa vào huấn luyện discriminator. Tại sao lại cần add thêm noise? Đó là bởi do ảnh hưởng của noise nên discriminator sẽ hoạt động không tốt quá và cũng không kém quá dẫn tới  $D(G(z^{(i)}))$  không hội tụ về 0. Do đó huấn luyện mô hình generator sẽ trở nên ổn định hơn. Tuy nhiên generator sẽ học cách sinh ra ảnh giả giống như ảnh thật + noise. Để ảnh giống thật hơn thì ta cần điều chỉnh giá trị noise nhỏ hơn.



Những điểm này có thể được khắc phục trong mô hình Wasserstein GAN.

## 7. Wasserstein GAN

### 7.1. Khái niệm liên tục Lipschitz (Lipschitz continuity)

Một khái niệm khá quan trọng ám chỉ các hàm số liên tục, khả vi và có độ lớn (về giá trị tuyệt đối) đạo hàm bị giới hạn đó là liên tục Lipschitz. Giả sử  $f: \mathbb{R} \rightarrow \mathbb{R}$ . Khi đó:

Top

$$\left| \frac{f(x_1) - f(x_2)}{x_1 - x_2} \right| \leq K, \forall x_1, x_2 \in \mathbb{R}$$

Với  $K$  là một giá trị hằng số thì hàm  $f$  được gọi là *liên tục K-lipchitz*. Ví dụ về một hàm *liên tục lipchitz* phổ biến đó chính là hàm  $\sin(x)$  vì giá trị *độ lớn* đạo hàm của nó bị chặn trên bởi 1. Hàm  $x^2$  không phải là hàm *liên tục lipchitz* vì đạo hàm của nó là  $2x$  không bị chặn trên. Tương tự như vậy  $\log(x)$  cũng không *liên tục lipchitz* vì đạo hàm của nó là  $\frac{1}{x}$  không bị chặn tại  $x = 0$ .

Do tính chất đạo hàm bị chặn nên khi sử dụng loss function là hàm *liên tục K-lipchitz* có thể khắc phục được hiện tượng đạo hàm bùng nổ (explosion) dẫn tới thiếu ổn định trong huấn luyện. Cụ thể hơn chúng ta cùng tìm hiểu qua Wasserstein GAN.

## 7.2. Wasserstein GAN

Model Wasserstein GAN sẽ áp dụng khoảng cách wasserstein để tìm ra phân phối gần nhất giữa 2 phân phối thật và giả. Việc giải trực tiếp bài toán khoảng cách Wasserstein là khá khó. Do đó áp dụng đối ngẫu Kantorovic-rubinstein chúng ta chuyển về bài toán đối ngẫu:

$$W(p_r, p_\theta) = \sup_{\|f\|_L \leq 1} \mathbf{E}_{x \sim p_r}[f(x)] - \mathbf{E}_{x \sim p_\theta}[f(x)]$$

Trong đó  $\|f\|_L$  là một hàm *liên tục 1-lipchitz* có độ lớn của đạo hàm bị chặn trên bởi 1. Khoảng cách wasserstein chính là *supremum* (cận trên nhỏ nhất) của chênh lệch kỳ vọng giữa phân phối của dữ liệu thật và dữ liệu sinh được đo lường thông qua hàm  $f(x)$  là một hàm *liên tục 1-lipchitz* thỏa mãn:

$$|f(x_1) - f(x_2)| \leq |x_1 - x_2|, \forall x_1, x_2 \in \mathbb{R}$$

Để hiểu hơn về phát biểu của bài toán đối ngẫu của khoảng cách Wasserstein các bạn xem thêm qua wasserstein duality (<https://vincentherrmann.github.io/blog/wasserstein/>). Tuy nhiên tôi không khuyến khích các bạn có nền tảng toán học yếu tìm hiểu sâu về mảng này mà chúng ta nên công nhận những gì đã được chứng minh.

Như vậy để áp dụng Wasserstein GAN, chúng ta chỉ cần tìm hàm *liên tục 1-lipchitz*. Kiến trúc mạng của discriminator sẽ được giữ nguyên và chúng ta chỉ bỏ qua hàm sigmoid ở cuối. Như vậy hàm dự báo output chính là một linear projection và đó là một hàm *liên tục 1-lipchitz*. Kết quả dự báo sau cùng sẽ đưa ra một điểm số scalar thay vì xác suất. Điểm này có thể được hiểu là điểm đánh giá chất lượng hình ảnh được sinh ra theo mức độ giống với ảnh thật. Tính chất của discriminator đã thay đổi từ phân loại ảnh thật/fake sang chấm điểm chất lượng ảnh nên để phù hợp với mục tiêu thì chúng ta thay đổi tên của discriminator thành critic.

Để hiểu rõ hơn sự khác biệt giữa kiến trúc GAN và WGAN chúng ta cùng theo dõi hình bên dưới:

### Kiến trúc GAN



### Kiến trúc WGAN



Tổng kết lại ta có sự khác biệt giữa GAN và WGAN đó là:

1. Bỏ hàm sigmoid ở critic model và thay vào đó là linear projection.
2. Ở model GAN sẽ thay đổi từ mô hình phân loại sang mô hình đánh giá. Do đó xác suất được chuyển sang điểm số có tác dụng đánh giá chất lượng ảnh tạo ra thay cho xác suất. Điểm này càng lớn thì ảnh càng giống với thật và điểm này càng nhỏ thì ảnh sẽ khác với thật. Nhân của mô hình cũng được thay đổi từ 0, 1 sang  $-1, 1$  để phù hợp hơn với mục tiêu là chấm điểm.
3. Critic sẽ được huấn luyện nhiều lượt hơn so với generator và quá trình huấn luyện sẽ được thực hiện xen kẽ giữa critic và generator.
4. Quá trình cập nhật gradient descent sẽ được thực hiện theo phương pháp RMSProp.
5. Ràng buộc độ lớn weights của mô hình về một khoảng giới hạn sau mỗi mini-batch.
6. Sử dụng RMSProp để cập nhật gradient descent với momentum = 0.

Quá trình huấn luyện GAN sẽ thực hiện xen kẽ những bước như sau:

---

**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:      $g_w \leftarrow \nabla_w [\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
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:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$ 
12: end while

```

---

## 8. Thực hành Wasserstein GAN

Chúng ta sẽ cùng thực hành Wasserstein GAN trên bộ dữ liệu MNIST. Code mẫu của bài này được tham khảo từ keras - wgan ([https://keras.io/examples/generative/wgan\\_gp/](https://keras.io/examples/generative/wgan_gp/)) và được mình điều chỉnh lại một chút. Mình sẽ giải thích chi tiết các bước huấn luyện và xử lý dữ liệu của model Wasserstein GAN.

Top



## 8.1. Load dữ liệu

Dữ liệu trong bài được lấy từ bộ dữ liệu MNIST bao gồm 60.000 ảnh huấn luyện và 10.000 ảnh test. Để quá trình hội tụ nhanh hơn thì giá trị cường độ pixels của ảnh được resize về khoảng [-1, 1] như sau:

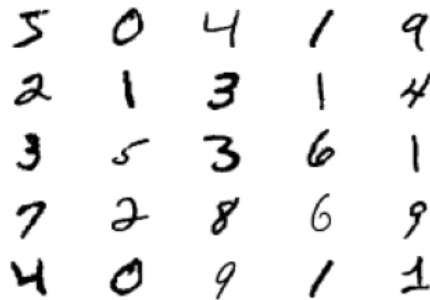
```

1      # Load dữ liệu mnist
2      from tensorflow.keras.datasets.mnist import load_data
3      import numpy as np
4
5      (trainX, trainy), (testX, testy) = load_data()
6      # Chuẩn hóa dữ liệu về khoảng [-1, 1]
7      trainX = (trainX - 127.5)/127.5
8      testX = (testX - 127.5)/127.5
9      trainX = np.expand_dims(trainX, axis = -1)
10     testX = np.expand_dims(testX, axis = -1)
11     print('Train', trainX.shape, trainy.shape)
12     print('Test', testX.shape, testy.shape)

1      Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/mnist.npz
2      11493376/11490434 [=====] - 0s 0us/step
3      Train (60000, 28, 28, 1) (60000,)
4      Test (10000, 28, 28, 1) (10000,)

1      import matplotlib.pyplot as plt
2
3      # plot images from the training dataset
4      def _plot(X):
5          for i in range(25):
6              # define subplot
7              plt.subplot(5, 5, 1 + i)
8              # turn off axis
9              plt.axis('off')
10             # plot raw pixel data
11             plt.imshow(X[i].reshape(28, 28), cmap='gray_r')
12     plt.show()
13
14     _plot(trainX[:25, :])

```



## 8.2. Discriminator

Chúng ta sẽ thiết kế Discriminator là một mạng CNN với layer Conv2D kết hợp xen kẽ với layer BatchNormalization. Stride và padding được kết hợp sao cho kích thước output shape của các layers thay đổi như sau: (32, 32) -> Conv\_s2 -> (16, 16) -> Conv\_s2 -> (8, 8) -> Conv\_s2 -> (4, 4) -> Conv\_s2 -> (2, 2). Để hạn chế overfitting thì layer Dropout được áp dụng cuối mỗi block CNN. Một chú ý quan trọng đó là chúng ta bỏ sigmoid activation ở cuối cùng mà thay vào đó sử dụng linear activation.

```
1  from tensorflow.keras.layers import Conv2D, BatchNormalization, Dropout, LeakyReLU, \
2  Input, ZeroPadding2D, Flatten, Dense, \
3  UpSampling2D, Reshape, Cropping2D, Activation
4
5  def conv_block(
6      x,
7      filters,
8      activation,
9      kernel_size=(3, 3),
10     strides=(1, 1),
11     padding="same",
12     use_bias=True,
13     use_bn=False,
14     use_dropout=False,
15     drop_value=0.5
16 ):
17     x = Conv2D(
18         filters, kernel_size, strides=strides, padding=padding, use_bias=use_bias
19     )(x)
20     if use_bn:
21         x = BatchNormalization()(x)
22     x = activation(x)
23     if use_dropout:
24         x = Dropout(drop_value)(x)
25     return x
```

```

1  from tensorflow.keras.models import Model
2
3  def get_discriminator_model():
4      img_input = Input(shape=IMG_SHAPE)
5      # Zero pad input để chuyển về kích thước (32, 32, 1).
6      x = ZeroPadding2D((2, 2))(img_input) # --> (32, 32)
7      x = conv_block(
8          x,
9          64,
10         kernel_size=(5, 5),
11         strides=(2, 2),
12         use_bn=False,
13         use_bias=True,
14         activation=LeakyReLU(0.2),
15         use_dropout=False,
16         drop_value=0.3,
17     ) # --> (16, 16)
18     x = conv_block(
19         x,
20         128,
21         kernel_size=(5, 5),
22         strides=(2, 2),
23         use_bn=False,
24         activation=LeakyReLU(0.2),
25         use_bias=True,
26         use_dropout=True,
27         drop_value=0.3,
28     ) # --> (8, 8)
29     x = conv_block(
30         x,
31         256,
32         kernel_size=(5, 5),
33         strides=(2, 2),
34         use_bn=False,
35         activation=LeakyReLU(0.2),
36         use_bias=True,
37         use_dropout=True,
38         drop_value=0.3,
39     ) # --> (4, 4)
40     x = conv_block(
41         x,
42         512,
43         kernel_size=(5, 5),
44         strides=(2, 2),
45         use_bn=False,
46         activation=LeakyReLU(0.2),
47         use_bias=True,
48         use_dropout=False,
49         drop_value=0.3,
50     ) # --> (2, 2)
51
52     x = Flatten()(x)
53     x = Dropout(0.2)(x)
54     # Bỏ sigmoid activation function và thay bằng linear activation
55     x = Dense(1, activation='linear')(x)
56
57     d_model = Model(img_input, x, name="discriminator")
58     return d_model
59
60
61 IMG_SHAPE = (28, 28, 1)
62 d_model = get_discriminator_model()
63 d_model.summary()

```

```

1  Model: "discriminator"
2
3  Layer (type)                Output Shape                Param #
4  =====
5  input_1 (InputLayer)        [(None, 28, 28, 1)]        0
6  .....
7
8  .....
9  dense (Dense)                (None, 1)                   2049
10 =====
11 Total params: 4,305,409
12 Trainable params: 4,305,409
13 Non-trainable params: 0
14

```

Top

## 8.3. Generator

Đầu vào của Generator là một véc tơ noise ngẫu nhiên có kích thước là 128. Sau đó đi qua các fully connected layer (hay còn gọi là dense layer) để tăng chiều dữ liệu và reshape về kích thước tensor3D với shape là (4, 4). Sử dụng layer Upsampling2D để tăng kích thước ảnh lên gấp đôi sau mỗi layer. Ngoài sử dụng Upsampling2D, chúng ta cũng có thể sử dụng tích chập chuyển vị Conv2DTranspose (<https://phamdinhhkhanh.github.io/2020/06/10/ImageSegmentation.html#7-t%C3%ADch-ch%E1%BA%ADp-chuy%E1%BB%83n-v%E1%BB%8B-transposed-convolution>) cũng có tác dụng tăng kích thước layers. Chuỗi kích thước layers của chúng ta sẽ theo thứ tự: (4, 4) -> Upsampling\_s2 -> (8, 8) -> Upsampling\_s2 -> (16, 16) -> Upsampling\_s2 -> (32, 32). Output shape cuối cùng sẽ bằng với kích thước của ảnh gốc.

Ngoài ra hàm activation là hàm tanh được áp dụng phía sau layer cuối cùng để chuẩn hóa giá trị output về khoảng  $[-1, 1]$ .

```

1  from tensorflow.keras.layers import Conv2D, BatchNormalization, Dropout, LeakyReLU, Input, ZeroPadding2D
2  def upsample_block(
3      x,
4      filters,
5      activation,
6      kernel_size=(3, 3),
7      strides=(1, 1),
8      up_size=(2, 2),
9      padding="same",
10     use_bn=False,
11     use_bias=True,
12     use_dropout=False,
13     drop_value=0.3,
14 ):
15     x = UpSampling2D(up_size)(x)
16     x = Conv2D(
17         filters, kernel_size, strides=strides, padding=padding, use_bias=use_bias
18     )(x)
19
20     if use_bn:
21         x = BatchNormalization()(x)
22
23     if activation:
24         x = activation(x)
25     if use_dropout:
26         x = Dropout(drop_value)(x)
27     return x
28
29 def get_generator_model():
30     noise = Input(shape=(noise_dim,))
31     x = Dense(4 * 4 * 256, use_bias=False)(noise)
32     x = BatchNormalization()(x)
33     x = LeakyReLU(0.2)(x)
34
35     x = Reshape((4, 4, 256))(x)
36     x = upsample_block(
37         x,
38         128,
39         LeakyReLU(0.2),
40         strides=(1, 1),
41         use_bias=False,
42         use_bn=True,
43         padding="same",
44         use_dropout=False,
45     )
46     x = upsample_block(
47         x,
48         64,
49         LeakyReLU(0.2),
50         strides=(1, 1),
51         use_bias=False,
52         use_bn=True,
53         padding="same",
54         use_dropout=False,
55     )
56     x = upsample_block(
57         x, 1, Activation("tanh"), strides=(1, 1), use_bias=False, use_bn=True
58     )
59     # Crop ảnh về kích thước (28, 28, 1)
60     x = Cropping2D((2, 2))(x)
61
62     g_model = Model(noise, x, name="generator")
63     return g_model
64
65 noise_dim = 128
66 g_model = get_generator_model()
67 g_model.summary()

```

1	Model: "generator"
2	
3	Layer (type)
4	Output Shape
5	Param #
6	=====
7	input_2 (InputLayer)
8	[(None, 128)]
9	0
10	=====
11	...
12	
13	cropping2d (Cropping2D)
14	(None, 28, 28, 1)
15	0
16	=====
17	Total params: 910,660
18	Trainable params: 902,082
19	Non-trainable params: 8,578
20	=====

## 8.4. WGAN-GP model

Ở bước này chúng ta sẽ khởi tạo mô hình WGAN-GP. Để huấn luyện mô hình WGAN-GP thì chúng ta sẽ cần khai báo generator, discriminator và các hàm loss function của generator và discriminator. Hàm quan trọng nhất của WGAN-GP là hàm `train_step()` sẽ có tác dụng huấn luyện trên từng batch. Ở hàm này chúng ta sẽ thực hiện các bước:

1. Lấy mẫu ngẫu nhiên ảnh real và ảnh fake. Trong đó ảnh real được lựa chọn từ mô hình và ảnh fake được tạo ra từ generator từ một véc tơ ngẫu nhiên.
2. Tính toán phạt cho gradient dựa trên chênh lệch giữa ảnh real và ảnh fake. Gradient được lấy trong bối cảnh của hàm `tf.GradientTape()`.
3. Nhân phạt của gradient với một giá trị hệ số alpha.
4. Tính loss function cho discriminator
5. Thêm phạt gradient vào loss function của discriminator
6. Cập nhật gradient theo loss function của discriminator và generator thông qua hàm `tape.apply_gradients()`

```

1  import tensorflow as tf
2
3  class WGAN(Model):
4      def __init__(
5          self,
6          discriminator,
7          generator,
8          latent_dim,
9          discriminator_extra_steps=3,
10         gp_weight=10.0,
11     ):
12         super(WGAN, self).__init__()
13         self.discriminator = discriminator
14         self.generator = generator
15         self.latent_dim = latent_dim
16         self.d_steps = discriminator_extra_steps
17         self.gp_weight = gp_weight
18
19     def compile(self, d_optimizer, g_optimizer, d_loss_fn, g_loss_fn):
20         super(WGAN, self).compile()
21         self.d_optimizer = d_optimizer
22         self.g_optimizer = g_optimizer
23         self.d_loss_fn = d_loss_fn
24         self.g_loss_fn = g_loss_fn
25
26     def gradient_penalty(self, batch_size, real_images, fake_images):
27         """
28         Tính toán phạt cho gradient
29         hàm loss này được tính toán trên ảnh interpolated và được thêm vào discriminator loss
30         """
31         # tạo interpolated image
32         alpha = tf.random.normal([batch_size, 1, 1, 1], 0.0, 1.0)
33         diff = fake_images - real_images
34         interpolated = real_images + alpha * diff
35
36         with tf.GradientTape() as gp_tape:
37             gp_tape.watch(interpolated)
38             # 1. dự đoán discriminator output cho interpolated image.
39             pred = self.discriminator(interpolated, training=True)
40
41             # 2. Tính gradients cho interpolated image.
42             grads = gp_tape.gradient(pred, [interpolated])[0]
43             # 3. Tính norm của gradients
44             norm = tf.sqrt(tf.reduce_sum(tf.square(grads), axis=[1, 2, 3]))
45             gp = tf.reduce_mean((norm - 1.0) ** 2)
46             return gp
47
48     def train_step(self, real_images):
49         if isinstance(real_images, tuple):
50             real_images = real_images[0]
51         batch_size = tf.shape(real_images)[0]
52
53         # Với mỗi một batch chúng ta sẽ thực hiện những bước sau
54         # following steps as laid out in the original paper.
55         # 1. Lấy mẫu ngẫu nhiên ảnh real và ảnh fake. Trong đó ảnh real được lựa chọn từ mô hình và ảnh
56         # 2. Tính toán phạt cho gradient dựa trên chênh lệch giữa ảnh real và ảnh fake. Gradient được l
57         # 3. Nhân phạt của gradient với một giá trị hệ số alpha.
58         # 4. Tính loss function cho discriminator
59         # 5. Thêm phạt gradient vào loss function của discriminator
60         # 6. Cập nhật gradient theo loss function của discriminator và generator thông qua hàm `tape.ap
61         # Chúng ta sẽ huấn luyện discriminator với d_steps trước. và sau đó chúng ta mới huấn luyện gene
62         # Như vậy cùng một batch thì discriminator được huấn luyện gấp d_steps lần so với generator.
63
64         for i in range(self.d_steps):
65             # Khởi tạo latent vector
66             random_latent_vectors = tf.random.normal(
67                 shape=(batch_size, self.latent_dim)
68             )
69
70             # Chúng ta phải thực thi trong tf.GradientTap() để lấy được giá trị gradient.
71             with tf.GradientTape() as tape:
72                 # Tạo ảnh fake từ latent vector
73                 fake_images = self.generator(random_latent_vectors, training=True)
74                 # Dự báo điểm cho ảnh fake
75                 fake_logits = self.discriminator(fake_images, training=True)
76                 # Dự báo điểm cho ảnh real
77                 real_logits = self.discriminator(real_images, training=True)
78
79                 # Tính discriminator loss từ ảnh real và ảnh fake
80                 d_cost = self.d_loss_fn(real_img=real_logits, fake_img=fake_logits)
81                 # Tính phạt gradient
82                 gp = self.gradient_penalty(batch_size, real_images, fake_images)
83                 # Thêm phạt gradient cho discriminator loss ban đầu

```

Top

```

84         d_loss = d_cost + gp * self.gp_weight
85
86         # Lấy gradient của discriminator loss
87         d_gradient = tape.gradient(d_loss, self.discriminator.trainable_variables)
88
89         # Cập nhật weights của discriminator sử dụng discriminator optimizer
90         self.d_optimizer.apply_gradients(
91             zip(d_gradient, self.discriminator.trainable_variables)
92         )
93
94         # Huấn luyện generator
95         # Khởi tạo véc tơ latent
96         random_latent_vectors = tf.random.normal(shape=(batch_size, self.latent_dim))
97         with tf.GradientTape() as tape:
98             # Tạo ảnh fake từ generator
99             generated_images = self.generator(random_latent_vectors, training=True)
100            # Get the discriminator logits for fake images
101            # Dự báo điểm cho ảnh fake từ discriminator
102            gen_img_logits = self.discriminator(generated_images, training=True)
103            # Tính generator loss
104            g_loss = self.g_loss_fn(gen_img_logits)
105
106            # Lấy gradient cho generator loss
107            gen_gradient = tape.gradient(g_loss, self.generator.trainable_variables)
108            # Cập nhật hệ số của generator sử dụng generator optimizer
109            self.g_optimizer.apply_gradients(
110                zip(gen_gradient, self.generator.trainable_variables)
111            )
112        return {"d_loss": d_loss, "g_loss": g_loss}

```

## 8.5. GAN monitoring

GANMonitoring là một Callback được gọi mỗi khi kết thúc một epoch để kiểm tra kết quả dự báo của một số ảnh sau mỗi epoch huấn luyện của mô hình.

```

1  from tensorflow.keras.callbacks import Callback
2  class GANMonitor(Callback):
3      def __init__(self, num_img=6, latent_dim=128):
4          self.num_img = num_img
5          self.latent_dim = latent_dim
6
7      def on_epoch_end(self, epoch, logs=None):
8          random_latent_vectors = tf.random.normal(shape=(self.num_img, self.latent_dim))
9          generated_images = self.model.generator(random_latent_vectors)
10         generated_images = (generated_images * 127.5) + 127.5
11
12         for i in range(self.num_img):
13             img = generated_images[i].numpy()
14             img = tf.keras.preprocessing.image.array_to_img(img)
15             img.save("generated_img_{i}_{epoch}.png".format(i=i, epoch=epoch))

```

## 8.6. Huấn luyện

```

1  from tensorflow.keras.optimizers import Adam
2  # Khởi tạo optimizer
3  # learning_rate=0.0002, beta_1=0.5 được khuyến nghị
4  generator_optimizer = Adam(
5      learning_rate=0.0002, beta_1=0.5, beta_2=0.9
6  )
7  discriminator_optimizer = Adam(
8      learning_rate=0.0002, beta_1=0.5, beta_2=0.9
9  )
10
11 # Define the loss functions to be used for discriminator
12 # This should be (fake_loss - real_loss)
13 # We will add the gradient penalty later to this loss function
14
15 # Xác định loss functions được sử dụng cho discriminator. = (fake_loss - real_loss).
16 # Ở hàm train_step của WGAN chúng ta phải thêm phạt gradient penalty cho hàm này.
17 def discriminator_loss(real_img, fake_img):
18     real_loss = tf.reduce_mean(real_img)
19     fake_loss = tf.reduce_mean(fake_img)
20     return fake_loss - real_loss
21
22
23 # Xác định loss function cho generator.
24 def generator_loss(fake_img):
25     return -tf.reduce_mean(fake_img)
26
27
28 # Khởi tạo tham số và mô hình
29
30 epochs = 20
31
32 # Callbacks
33 cbk = GANMonitor(num_img=3, latent_dim=noise_dim)
34
35 # wgan model
36 wgan = WGAN(
37     discriminator=d_model,
38     generator=g_model,
39     latent_dim=noise_dim,
40     discriminator_extra_steps=3,
41 )
42
43 # Compile model
44 wgan.compile(
45     d_optimizer=discriminator_optimizer,
46     g_optimizer=generator_optimizer,
47     g_loss_fn=generator_loss,
48     d_loss_fn=discriminator_loss,
49 )
50
51 BATCH_SIZE = 512
52 # Huấn luyện
53 wgan.fit(trainX, batch_size=BATCH_SIZE, epochs=epochs, callbacks=[cbk])

```

Epoch 20/20  
118/118 [=====] - 455s 4s/step - d\_loss: -3.5191 - g\_loss: 6.5974

Lưu lại mô hình discriminator và generator.

```

1  from google.colab import drive
2  import os
3
4  drive.mount('/content/gdrive/')
5  os.chdir('gdrive/My Drive/Colab Notebooks/GAN')

```

```

1  !mkdir wgan
2  g_model.save('wgan/generator.h5')
3  d_model.save('wgan/discriminator.h5')

```

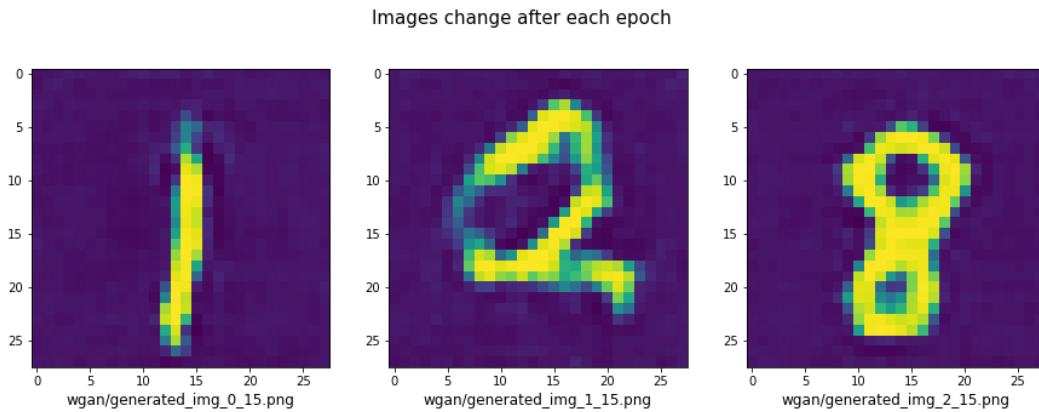
Tiếp theo ta sẽ đọc một vài kết quả huấn luyện từ epoch 0, 1 và 2.



```

1 import matplotlib.pyplot as plt
2
3 images = ["wgan/generated_img_0_15.png", "wgan/generated_img_1_15.png", "wgan/generated_img_2_15.png"]
4 fg, ax = plt.subplots(1, 3, figsize=(15, 5))
5 fg.suptitle('Images change after each epoch', fontsize=15)
6
7 for i in np.arange(3):
8     X = plt.imread(images[i])
9     ax[i].imshow(X)
10    ax[i].set_xlabel(images[i], fontsize=12)

```



## 9. Tổng kết

Như vậy qua bài này chúng ta đã biết được các kiến thức rất cơ bản và quan trọng về tiêu chuẩn đo lường khoảng cách chính giữa hai phân phối xác suất đó là các độ đo: Kullback Leibler, Jensen-Shannon kèm theo ưu, nhược điểm của chúng. Đồng thời chúng ta cũng biết được rằng giá trị hội tụ GAN chính là một hàm của khoảng cách Jensen-Shannon. Ý tưởng về việc áp dụng một hàm loss function liên tục 1-*Lipschitz* có đạo hàm bị chặn trong WGAN để tạo ra một mô hình huấn luyện ổn định hơn và việc chuyển mô hình từ phân loại sang chấm điểm chất lượng ảnh kèm theo code hướng dẫn.

Mặc dù chỉ thay đổi về hàm loss function và dường như là giữ nguyên kiến trúc nhưng WGAN đã tạo ra một lớp mô hình huấn luyện ổn định, hiệu quả và đưa ra chất lượng ảnh tốt hơn so với các mô hình được huấn luyện theo DCGAN.

## 10. Tài liệu

1. GAN wasserstein and GAN wasserstein GP - jonathan\_hui ([https://medium.com/@jonathan\\_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490](https://medium.com/@jonathan_hui/gan-wasserstein-gan-wgan-gp-6a1a2aa1b490))
2. GAN spectral normalization - jonathan\_hui ([https://medium.com/@jonathan\\_hui/gan-spectral-normalization-893b6a4e8f53](https://medium.com/@jonathan_hui/gan-spectral-normalization-893b6a4e8f53))
3. wasserstein GAN - vincentherrmann (<https://vincentherrmann.github.io/blog/wasserstein/>)
4. How to implement wasserstein loss for GAN - machinelearningmastery (<https://machinelearningmastery.com/how-to-implement-wasserstein-loss-for-generative-adversarial-networks/>)
5. How to code a wasserstein GAN from scratch - machinelearningmastery (<https://machinelearningmastery.com/how-to-code-a-wasserstein-generative-adversarial-network-wgan-from-scratch/>)