

The story about WGAN

Ngày 23 tháng 11 năm 2019

GAN - What is the optimal value for D ?

Loss function of GAN

$$L = \mathbb{E}_{x \sim P_r} [\log D(x)] + \mathbb{E}_{x \sim P_g} [\log (1 - D(x))]$$

$$f(D(x)) = P_r(x) \log D(x) + P_g(x) \log (1 - D(x))$$

$$\frac{\partial f(D(x))}{\partial D(x)} = P_r(x) * \frac{1}{D(x)} - P_g(x) \frac{1}{1 - D(x)} = 0$$

$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)}$$

Nếu hình ảnh là một hình ảnh thật thì Discriminator trả về xác suất là 0.5. Mặt khác, nếu hình ảnh là giả thì Discriminator trả về 0

The main problem : optimal discriminator

Loss function of GAN

$$L = \mathbb{E}_{x \sim P_r} [\log D(x)] + \mathbb{E}_{x \sim P_g} [\log (1 - D(x))]$$

$$\begin{aligned} L &= \mathbb{E}_{x \sim P_r} [\log D(x)] + \mathbb{E}_{x \sim P_g} [\log (1 - D(x))] \\ &= \mathbb{E}_{x \sim P_r} \left[\log \frac{P_r(x)}{P_r(x) + P_g(x)} \right] + \mathbb{E}_{x \sim P_g} \left[\log \left(1 - \frac{P_r(x)}{P_r(x) + P_g(x)} \right) \right] \\ &= \mathbb{E}_{x \sim P_r} \left[\log \frac{P_r(x)}{2 \times \frac{1}{2} (P_r(x) + P_g(x))} \right] + \mathbb{E}_{x \sim P_g} \left[\log \left(\frac{P_g(x)}{2 \times \frac{1}{2} (P_r(x) + P_g(x))} \right) \right] \\ &= \mathbb{E}_{x \sim P_r} \left[\log \frac{P_r(x)}{\frac{1}{2} (P_r(x) + P_g(x))} \right] + \mathbb{E}_{x \sim P_g} \left[\log \left(\frac{P_g(x)}{\frac{1}{2} (P_r(x) + P_g(x))} \right) \right] - 2 \log 2 \\ &= \mathbb{E}_{x \sim P_r} \left[\log \frac{P_r(x)}{P_{average}(x)} \right] + \mathbb{E}_{x \sim P_g} \left[\log \left(\frac{P_g(x)}{P_{average}(x)} \right) \right] - 2 \log 2 \\ &= 2JS(P_r | P_g) - 2 \log 2 \end{aligned}$$

Như kết quả trên, nếu ta tiếp tục training Discriminator thì giá trị của hàm loss sẽ bằng 0, và JS-divergence sẽ bằng $\log 2$

Different distance

- The Total Variation (TV) distance is :

$$\delta(P_r, P_g) = \sup_{A \in \Sigma} |P_r(A) - P_g(A)|$$

- The Kullback-Leibler (KL) divergence is

$$KL(P_r \| P_g) = \int \log \left(\frac{P_r(x)}{P_g(x)} \right) P_r(x) d\mu(x)$$

- The Jenson-Shannon (JS) divergence

$$JS(P_r, P_g) = KL(P_r \| P_m) + KL(P_g \| P_m)$$

- The Earth Mover (EM) or Wasserstein distance

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

Different distance

Example

Giả sử ta có hai distributions, \mathbb{P} và \mathbb{Q} :

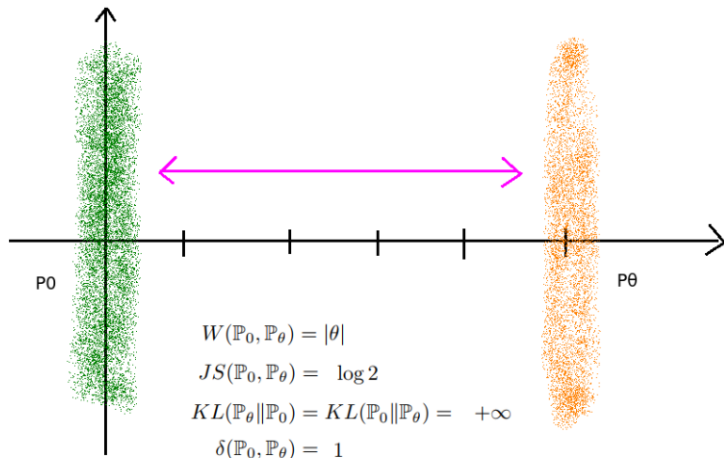
$$\forall (x, y) \in \mathbb{P}, x = 0 \text{ và } y \sim U(0, 1)$$

$$\forall (x, y) \in \mathbb{Q}, x = \theta, 0 \leq \theta \leq 1 \text{ và } y \sim U(0, 1)$$

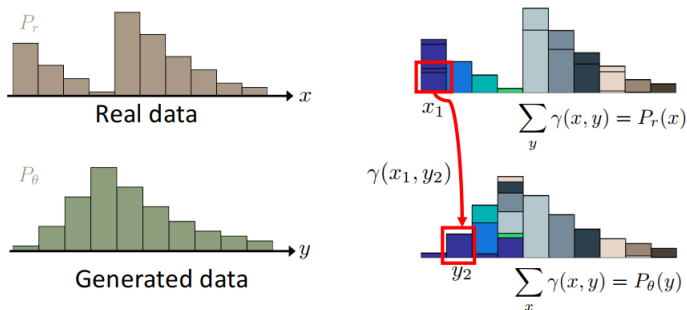
- $W(\mathbb{P}_0, \mathbb{P}_\theta) = |\theta|$
- $JS(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} 2 \log 2 & \text{nếu } \theta \neq 0 \\ 0 & \text{nếu } \theta = 0 \end{cases}$
- $KL(\mathbb{P}_0 \parallel \mathbb{P}_\theta) = KL(\mathbb{P}_\theta \parallel \mathbb{P}_0) = \begin{cases} +\infty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

Chỉ có duy nhất Wasserstein metric mới có a smooth measure, phù hợp để sử dụng gradient.

Different distance



Earth-Mover (EM) distance/ Wasserstein Metric



Earth-Mover (EM) distance/ Wasserstein Metric

$$\text{EMD}(P_r, P_\theta) = \inf_{\gamma \in \Pi} \sum_{x, y} \|x - y\| \gamma(x, y) = \inf_{\gamma \in \Pi} E_{(x, y) \sim \gamma} \|x - y\|$$

Earth-Mover (EM) distance/ Wasserstein Metric



$$\text{EMD}(P_r, P_\theta) = \inf_{\gamma \in \Pi} \sum_{x,y} \|x - y\| \gamma(x, y) = \inf_{\gamma \in \Pi} E_{(x,y) \sim \gamma} \|x - y\|$$

Đặt

$$\begin{cases} \Gamma = \gamma(x, y) \\ D = \|x - y\| \end{cases}$$

Ta viết lại :

$$\text{EMD}(P_r, P_\theta) = \inf_{\gamma \in \Pi} \langle D, \Gamma \rangle_F$$

Để tính được phương án di chuyển tối ưu Γ , ta dùng Linear Programming

Linear Programming

Our Problem

Objective function : $\text{EMD}(P_r, P_\theta) = \inf_{\gamma \in \Pi(P_r, P_\theta)} \langle D, \gamma \rangle_F$

Constraint :

$$\sum_y \gamma(x, y) = P_r(x)$$

$$\sum_x \gamma(x, y) = P_\theta(y)$$

$$\forall x, y \quad \gamma(x, y) \geq 0$$

Linear Programming

Objective function : minimize $z = c^T x$

Constraint :

$$Ax = b$$

$$x \geq 0$$

Objective function :

Minimize $z = c^T x$

$$\text{EMD}(P_r, P_\theta) = \inf_{\gamma \in \Pi(P_r, P_\theta)} \langle D, \Gamma \rangle_F$$

$$\begin{array}{cc} c = \text{vec}(D) & x = \text{vec}(\Gamma) \\ \left[\begin{array}{c} \|x_1 - y_1\| \\ \|x_1 - y_2\| \\ \vdots \\ \|x_2 - y_1\| \\ \|x_2 - y_2\| \\ \vdots \\ \|x_n - y_1\| \\ \|x_n - y_2\| \\ \vdots \end{array} \right] & \left[\begin{array}{c} \gamma(x_1, y_1) \\ \gamma(x_1, y_2) \\ \vdots \\ \gamma(x_2, y_1) \\ \gamma(x_2, y_2) \\ \vdots \\ \gamma(x_n, y_1) \\ \gamma(x_n, y_2) \\ \vdots \end{array} \right] \end{array}$$

Linear Programming

Điều kiện :

$$Ax = b$$

$$\sum_y \gamma(x, y) = P_r(x)$$

$$\sum_x \gamma(x, y) = P_\theta(y)$$

$$\begin{array}{c}
 A \\
 \left[\begin{array}{ccccccccc}
 1 & 1 & \cdots & 0 & 0 & \cdots & 0 & 0 & \cdots \\
 0 & 0 & \cdots & 1 & 1 & \cdots & 0 & 0 & \cdots \\
 \vdots & \vdots & \cdots & \vdots & \vdots & \cdots & \vdots & \vdots & \cdots \\
 0 & 0 & \cdots & 0 & 0 & \cdots & 1 & 1 & \cdots \\
 \hline
 1 & 0 & \cdots & 1 & 0 & \cdots & 1 & 0 & \cdots \\
 0 & 1 & \cdots & 0 & 1 & \cdots & 0 & 1 & \cdots \\
 \vdots & \vdots & \cdots & \vdots & \vdots & \cdots & \vdots & \vdots & \cdots \\
 0 & 0 & \cdots & 0 & 0 & \cdots & 0 & 0 & \cdots
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 x = \text{vec}(\Gamma) \\
 \left[\begin{array}{c}
 \gamma(x_1, y_1) \\
 \gamma(x_1, y_2) \\
 \vdots \\
 \hline
 \gamma(x_2, y_1) \\
 \gamma(x_2, y_2) \\
 \vdots \\
 \hline
 \gamma(x_n, y_1) \\
 \gamma(x_n, y_2) \\
 \vdots
 \end{array} \right]
 \end{array}
 =
 \begin{array}{c}
 b = \left[\begin{array}{c}
 P_r \\
 P_\theta \\
 P_r(x_1) \\
 P_r(x_2) \\
 \vdots \\
 P_r(x_n) \\
 P_\theta(y_1) \\
 \vdots \\
 P_\theta(y_2) \\
 P_\theta(y_n)
 \end{array} \right]
 \end{array}$$

Dual form

Primary form

minimize $z = c^T x$
so that $Ax = b$
and $x \geq 0$

Dual form

maximize $\tilde{z} = b^T y$
so that $A^T y \leq c$

Dual form

Objective function

$$\tilde{z} = b^T y$$

$$b = \begin{bmatrix} P_r \\ P_\theta \\ P_r(x_1) \\ P_r(x_2) \\ \vdots \\ P_r(x_n) \\ P_\theta(y_1) \\ P_\theta(y_n) \\ \vdots \\ P_\theta(y_n) \end{bmatrix} \quad y = \begin{bmatrix} f \\ g \\ f(x_1) \\ f(x_1) \\ \vdots \\ f(x_1) \\ g(x_1) \\ g(x_1) \\ \vdots \\ g(x_1) \end{bmatrix}$$
$$\Rightarrow \text{EMD}(P_r, P_\theta) = f^T P_r + g^T P_\theta$$

Dual form

Điều kiện

$$A^T y \leq c$$

$$\begin{array}{c}
 A^T \\
 \left[\begin{array}{cccc|cccc}
 1 & 0 & \dots & 0 & 1 & 0 & \dots & 0 \\
 1 & 0 & \dots & 0 & 0 & 1 & \dots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 0 & 1 & \dots & 0 & 1 & 0 & \dots & 0 \\
 0 & 1 & \dots & 0 & 0 & 1 & \dots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 0 & 0 & \dots & 1 & 1 & 0 & \dots & 0 \\
 0 & 0 & \dots & 1 & 0 & 1 & \vdots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots
 \end{array} \right]
 \end{array}$$

$$y = \begin{bmatrix} f \\ g \end{bmatrix}$$

$$\begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \\ g(x_1) \\ g(x_2) \\ \vdots \\ g(x_n) \end{bmatrix}$$

\leq

$$\begin{array}{c}
 c = \text{vec}(D) \\
 \left[\begin{array}{c} \|x_1 - x_1\| \\ \|x_1 - x_2\| \\ \vdots \\ \|x_2 - x_1\| \\ \|x_2 - x_2\| \\ \vdots \\ \|x_n - x_1\| \\ \|x_n - x_2\| \\ \vdots \end{array} \right]
 \end{array}$$

$$\Rightarrow \forall i, j \quad f(x_i) + g(x_j) \leq \|x_i - x_j\|$$

Dual form

Điều kiện

$$\forall i, j \quad f(x_i) + g(x_j) \leq \|x_i - x_j\|$$

If $i = j$:

$$f(x_i) + g(x_i) \leq \|x_i - x_i\| = 0$$

Ta muốn maximize : $\text{EMD}(P_r, P_\theta) = f^T P_r + g^T P_\theta$

$$\Rightarrow f(x_i) = -g(x_i)$$

Từ (1) và (2), ta được :

$$\begin{cases} f(x) - f(x) \leq \|x - x\| \\ f(x) - f(x) \geq -\|x - x\| \end{cases}$$

$$\Rightarrow -1 \leq \frac{f(x) - f(x)}{\|x - x\|} \leq 1$$

$$\Rightarrow \|f\|_{L \leq 1}$$

Vậy dạng Dual form của EMD là :

$$\text{EMD}(P_r, P_\theta) = \sup_{\|f\|_{L \leq 1}} \mathbb{E}_{x \sim P_r} f(x) - \mathbb{E}_{x \sim P_\theta} f(x)$$

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: : α , the learning rate. c , the clipping parameter. m , the batch size. n_{critic} , the number of iterations of the critic per generator iteration.

Require: : w_0 , initial critic parameters. θ_0 , initial generator's parameters.

```

1: while  $\theta$  has not converged do                                     ← Train until convergence
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
    
```

Critic Training

Generator Training

Ta nhận thấy được quá trình training của WGAN khác với GAN thông thường :

- The critics sẽ được cập nhật nhiều lần
- Trong quá trình loss, WGAN không sử dụng hàm logarit (không sử dụng cross entropy)
- Sử dụng weight clipping
- Không sử dụng các momentum

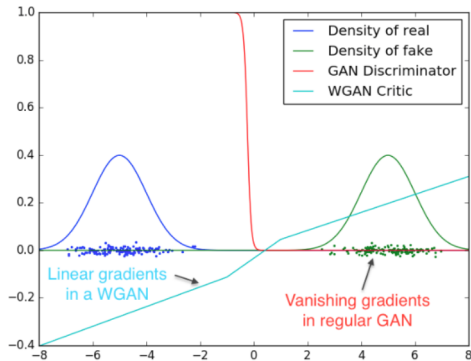
- Trong GAN, Discriminator cố gắng maximizes :

$$\frac{1}{m} \left(\sum_{i=1}^m \log D(x^{(i)}) + \sum_{i=1}^m \log (1 - D(g_{\theta}(z^{(i)}))) \right)$$

trong đó, $D(x)$ là một xác suất $p \in (0, 1)$ Trong WGAN, nó không đưa ra một xác suất cụ thể, mà nó đưa ra khoảng cách EMD của hai distribution

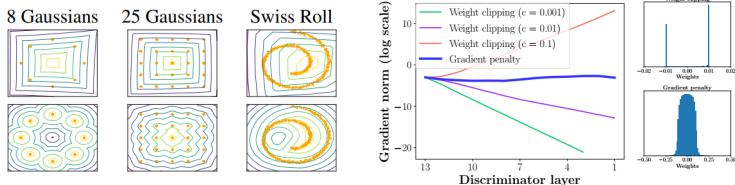
- Trong paper gốc của GAN, các độ đo khoảng cách sẽ dựa vào JS divergence
Trong WGAN, nó được thay thế bằng Wasserstein distance

Result of WGAN



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

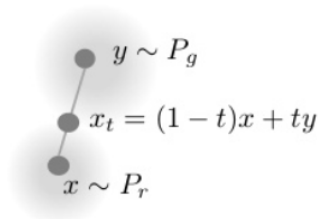
Problem in WGAN



- Sau khi WGAN được công bố, rất nhiều developer tìm ra được vấn đề của WGAN là tốc độ hội tụ của nó chậm và chất lượng hình ảnh được tạo ra không được tốt
- Nếu như ta chọn giá trị của weight clipping quá nhỏ thì sẽ bị dẫn đến hiện tượng vanishing gradient, ngược lại nếu quá lớn thì sẽ dẫn đến hiện tượng exploding gradient
- Việc sử dụng weight clipping sẽ làm giảm khả năng học được các distribution phức tạp của tập dataset, chỉ học được các function đơn giản

Gradient penalty

Optimal critic has gradient with norm 1 almost everywhere under P_r và P_g



$$\nabla f^*(x_t) = \frac{y - x_t}{\|y - x_t\|}$$
$$\Rightarrow \|\nabla f^*(x_t)\| = 1$$

Loss function

$$L = \underbrace{\mathbb{E}_{\tilde{x} \sim P_g} [D(\tilde{x})] - \mathbb{E}_{x \sim P_r} [D(x)]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{x} \sim P_{\hat{x}}} \left[(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2 \right]}_{\text{Gradient penalty}}$$

Algorithm 1 WGAN with gradient penalty. We use default values of $\lambda = 10$, $n_{\text{critic}} = 5$, $\alpha = 0.0001$, $\beta_1 = 0$, $\beta_2 = 0.9$.

Require: The gradient penalty coefficient λ , the number of critic iterations per generator iteration n_{critic} , the batch size m , Adam hyperparameters α, β_1, β_2 .

Require: initial critic parameters w_0 , initial generator parameters θ_0 .

```

1: while  $\theta$  has not converged do
2:   for  $t = 1, \dots, n_{\text{critic}}$  do
3:     for  $i = 1, \dots, m$  do
4:       Sample real data  $x \sim \mathbb{P}_r$ , latent variable  $z \sim p(z)$ , a random number  $\epsilon \sim U[0, 1]$ .
5:        $\tilde{x} \leftarrow G_\theta(z)$ 
6:        $\hat{x} \leftarrow \epsilon x + (1 - \epsilon)\tilde{x}$ 
7:        $L^{(i)} \leftarrow D_w(\tilde{x}) - D_w(x) + \lambda(\|\nabla_{\hat{x}} D_w(\hat{x})\|_2 - 1)^2$ 
8:     end for
9:      $w \leftarrow \text{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, \alpha, \beta_1, \beta_2)$ 
10:   end for
11:   Sample a batch of latent variables  $\{z^{(i)}\}_{i=1}^m \sim p(z)$ .
12:    $\theta \leftarrow \text{Adam}(\nabla_\theta \frac{1}{m} \sum_{i=1}^m -D_w(G_\theta(z)), \theta, \alpha, \beta_1, \beta_2)$ 
13: end while

```

Ở thuật toán này ta đã có thể sử dụng được momentum để cập nhật được model. Mỗi khi cập nhật critics, gradient penalty được thêm vào hàm loss