

对抗生成网络

[ref.GAN](#) [ref.WGAN](#) [ref.郑华滨知乎博客](#) [ref.WGAN-GP](#)

GAN的目标函数是

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

分析目标函数，其中对判别器 D ，其目标函数是： $E_r[\log D(x)] + E_g[\log(1 - D(G(z)))]$ ，对生成器 G ，其目标函数是： $E_g[\log(1 - D(G(z)))]$ 。

对判别器，一个样本的损失函数是

$$-P_r(x) \log D(x) - P_g(x) \log[1 - D(x)] \dots (2)$$

首先最大化 $D(x)$ ，对 $D(x)$ 求导，得到

$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)} \dots (3)$$

其含义是，最优判别器 $D^*(x)$ 可以准确给出输入样本中真实样本所占的比例。把(3)带入判别器的损失函数，得到

$$\mathbb{E}_{x \sim P_r} \log \frac{P_r(x)}{\frac{1}{2}[P_r(x) + P_g(x)]} + \mathbb{E}_{x \sim P_g} \log \frac{P_g(x)}{\frac{1}{2}[P_r(x) + P_g(x)]} =$$

即为2倍JS散度，

$$2JS(P_r \| P_g) = 2 \log 2 \dots (4)(b)$$

原始GAN的损失函数是JS散度，对JS散度，当两个分布是不相交时，JS散度恒为 $\log 2$ ，如图1，[明](#))

The Total Variation (TV) distance

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

The Kullback-Leibler (KL) divergence

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

The Jensen-Shannon (JS) divergence

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

其中 $\mathbb{P}_m = (\mathbb{P}_r + \mathbb{P}_g) / 2$

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

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

对图1的分布，三种不同距离分别为，后两种距离不连续

$$W(\mathbb{P}_0, \mathbb{P}_\theta) = |\theta|$$

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

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

由于推土机距离是计算下确界 \inf , 很难计算, 根据Kantorovich-Rubinstein duality(对偶), 具体计

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

距离计算变为求两个分布的期望差的上确界 \sup 。个人理解, 相较于下确界 \inf , \sup 的好处是, 我来进行优化。这是因为目标是最小化其上确界, 一般情况下的期望差必然更小

WGAN算法:

Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper use 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 number of samples per batch. 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
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:      $w \leftarrow w + \alpha \cdot \text{RMSPProp}(w, g_w)$ 
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

```

WGAN算法有一个混合真实样本和生成样本的步骤, 判别器损失和原始GAN相同, 生成器损失变为对WGAN-GP, WGAN对权值裁剪, WGAN-GP对梯度进行裁剪。原文中提到了一些梯度裁剪的好处如

Gradient norms of deep WGAN critics during training on toy datasets either explode or vanish when using weight clipping, but not when using a gradient penalty. (right) Weight clipping (top) pushes weights towards two values (the extremes of the clipping range), unlike gradient penalty (bottom).

WGAN-GP的算法如下:

Algorithm 1 WGAN with gradient penalty. We use default values of $\lambda = 10.00001$, $\beta_1 = 0$, $\beta_2 = 0.9$.

Require: The gradient penalty coefficient λ , the number of critic iterations per 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  $\mathbf{x} \sim \mathbb{P}_r$ , latent variable  $\mathbf{z} \sim p(\mathbf{z})$ , a random number
5:        $\tilde{\mathbf{x}} \leftarrow G_{\theta}(\mathbf{z})$ 
6:        $\hat{\mathbf{x}} \leftarrow \epsilon \mathbf{x} + (1 - \epsilon) \tilde{\mathbf{x}}$ 
7:        $L^{(i)} \leftarrow D_w(\tilde{\mathbf{x}}) - D_w(\mathbf{x}) + \lambda(\|\nabla_{\hat{\mathbf{x}}} D_w(\hat{\mathbf{x}})\|_2 - 1)^2$ 
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  $\{\mathbf{z}^{(i)}\}_{i=1}^m \sim p(\mathbf{z})$ .
12:    $\theta \leftarrow \text{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m -D_w(G_{\theta}(\mathbf{z})), \theta, \alpha, \beta_1, \beta_2)$ 
13: end while

```

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy import stats
4 mu = 0
5 sd = 0.3
6 t1 = np.linspace(-1, 1, 1000)
7 t2 = np.linspace(2, 4, 1000)
8 y1 = stats.norm(mu, sd).pdf(t1)
9 y2 = stats.norm(mu, sd).pdf(t1)
10
11 plt.figure(1)
12 plt.plot(t1, y1)
13 plt.plot(t2, y2)
14 plt.title("Figure 1")
15
16 plt.figure(2)
17 plt.plot(t1, y1)
18 plt.plot(t1, y1+0.5)
19 plt.title("Figure 2")

```

☞

```
Text(0.5, 1.0, 'Figure 2')
```

