



# Improved Wasserstein-GAN

(April 26, 2017)

Cheng Jianlin



# Outline

- **I. Introduction**
- **II. Optimal WGAN critic**
- **III. Improved Wasserstein GAN**
- **IV. Results**
- **V. Conclusion**



# • I. Introduction

## GAN

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

$$\begin{aligned} C(G) &= \max_D V(G, D) \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log(1 - D_G^*(G(\mathbf{z})))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \right] \end{aligned}$$

**Min:**  $C(G) = -\log(4) + 2 \cdot JSD(p_{\text{data}} \| p_g)$



# • I. Introduction

## Wasserstein GAN

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

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

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)] - \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})]$$

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



## • I. Introduction

# Wasserstein GAN

---

**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
```

---



## • II. Optimal WGAN critic

**Lemma 1.** Let  $\mathbb{P}_r$  and  $\mathbb{P}_g$  be two distributions in  $\mathcal{X}$ , a compact metric space. Then, there is a 1-Lipschitz function  $f^*$  which is the optimal solution of

$$\max_{\|f\|_L \leq 1} \mathbb{E}_{y \sim \mathbb{P}_r} [f(y)] - \mathbb{E}_{x \sim \mathbb{P}_g} [f(x)]$$

Let  $\pi$  be the optimal coupling between  $\mathbb{P}_r$  and  $\mathbb{P}_g$ , defined as the minimizer of:

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

Where  $\Pi(\mathbb{P}_r, \mathbb{P}_g)$  is the set of joint distributions  $\pi(x, y)$  whose marginals are  $\mathbb{P}_r$  and  $\mathbb{P}_g$ , respectively. Then, if  $f^*$  is differentiable<sup>2</sup> and  $x_t = tx + (1 - t)y$  with  $0 \leq t \leq 1$ ,

$$\mathbb{P}_{(x,y) \sim \pi} \left[ \nabla f^*(x_t) = \frac{y - x_t}{\|y - x_t\|} \right] = 1$$



## • III. Improved Wasserstein GAN

### Gradient penalty

As we discussed in the previous section, the use of weight clipping in WGAN can result in undesirable behaviors. We consider an alternative method to enforce the Lipschitz constraint on the training objective: a differentiable function is 1-Lipschitz if and only if it has gradients with norm less than or equal to 1 everywhere, so we would like to directly constrain the gradient norm of our critic function with respect to its input.

Exactly enforcing this constraint is not easily tractable, so instead we enforce a soft version: at certain points sampled from a distribution over the input space  $\hat{x} \sim \mathbb{P}_{\hat{x}}$ , we evaluate the gradient of the critic  $\nabla_{\hat{x}} D(\hat{x})$  and penalize its squared distance from 1 in the critic loss function. Our new objective for the critic is:

$$L = \underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Our gradient penalty}} \quad (3)$$

In the limit of high  $\lambda$ , the optimal critic under our formulation is still the optimal critic under the true Kantorovich-Rubinstein dual. Therefore, given enough capacity on the critic the cost function



# • III. Improved Wasserstein GAN

---

**Algorithm 1** WGAN with gradient penalty. We use default values of  $\lambda = 10$ ,  $n_{\text{critic}} = 5$ ,  $\alpha = 0.0001$ ,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.9$ .

---

**Require:** The gradient penalty coefficient  $\lambda$ , the number of critic iterations per generator iteration  $n_{\text{critic}}$ , the batch size  $m$ , Adam hyperparameters  $\alpha, \beta_1, \beta_2$ .

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

```
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
```

---





# • III. Improved Wasserstein GAN

**Sampling along straight lines** The gradient term  $||\nabla_{\hat{x}} D(\hat{x})||_2$  is with respect to the points  $\hat{x}$ , not the parameters of  $D$ . We implicitly define the distribution of points  $\mathbb{P}_{\hat{x}}$  at which to penalize the gradient by taking straight lines between points in the data distribution  $\mathbb{P}_r$  and the generator distribution  $\mathbb{P}_g$ :

$$\epsilon \sim U[0, 1], x \sim \mathbb{P}_r, \tilde{x} \sim \mathbb{P}_g \quad (4)$$

$$\hat{x} = \epsilon x + (1 - \epsilon)\tilde{x} \quad (5)$$

Our motivation for penalizing the gradient over this distribution comes from the fact that the graph of the optimal critic consists of straight lines connecting points from  $\mathbb{P}_r$  and  $\mathbb{P}_g$  (see subsection 2.3). Given that enforcing the Lipschitz constraint everywhere is intractable, enforcing it only along these straight lines seems sufficient and experimentally results in good performance.

**Hyperparameters** Our proposed penalty term introduces one hyperparameter,  $\lambda$ , which controls the trade-off between optimizing the penalty term and the original objective. All experiments in this paper use  $\lambda = 10$ , which we found to work well across a variety of architectures and datasets ranging from toy tasks to large ImageNet CNNs.

**No critic batch normalization** Most prior GAN implementations (Radford et al., 2015; Salimans et al., 2016; Arjovsky et al., 2017) make use of batch normalization in both the generator and the discriminator to help stabilize training. However using batch normalization in the discriminator changes the form of the problem: instead of specifying a function mapping a single input to a single output, a discriminator with batch normalization specifies a function mapping from an entire batch of inputs to a batch of outputs (Salimans et al., 2016).

Our penalized training objective is no longer valid in this setting, since we penalize the norm of the critic's gradient with respect to each input independently, and not the entire batch. To resolve



## • IV. Results

### CIFAR-10

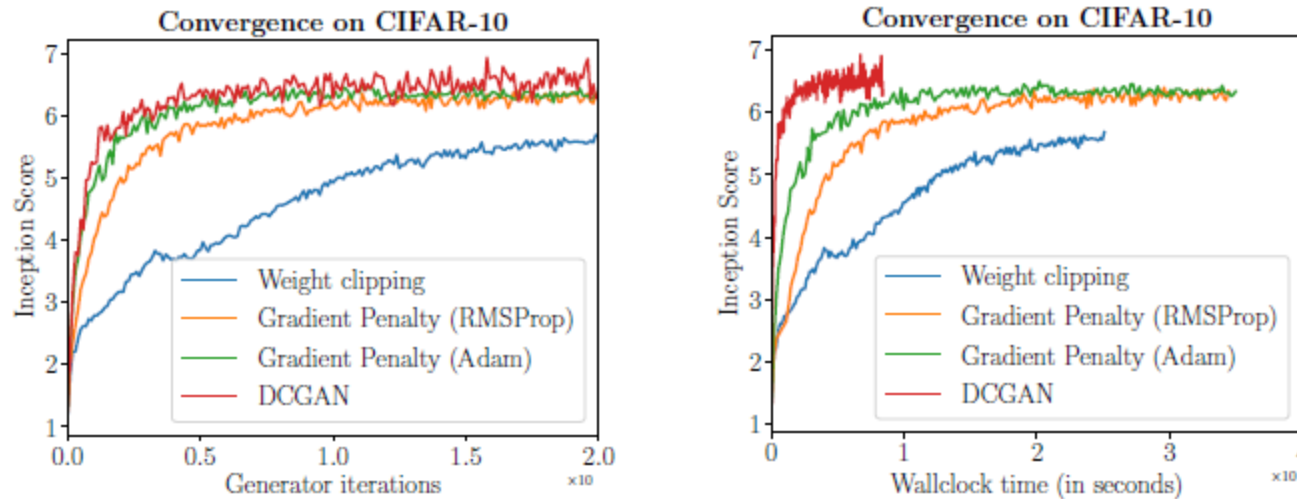


Figure 3: Plots of CIFAR-10 Inception score over generator iterations (left) or wall-clock time (right) for four models: WGAN with weight clipping, WGAN with gradient penalty and RMSProp (to control for the optimizer), WGAN with gradient penalty and Adam, and DCGAN. Even with the same learning rate, gradient penalty significantly outperforms weight clipping. DCGAN converges faster, but WGAN with gradient penalty achieves similar scores with improved stability.



## • IV. Results

**LSUN(Large-scale Scene Understanding) bedrooms dataset**



## • V. Conclusion

**Improve performance and successfully train difficult GAN architectures**

**we think our work opens the path for strong modeling performance on large-scale image datasets and language**

**it might stabilize training by encouraging the discriminator to learn smoother decision boundaries**



# Bibliography

- [1] Wasserstein GAN
- [2] Improved Training of Wasserstein GANs



# Thanks !

