

Comprehensive Wasserstein GAN (WGAN) Cheat Sheet

AI Assistant

October 21, 2024

1 Key Concepts

- **Wasserstein Distance:**
 - Also known as Earth Mover's Distance
 - Measures the distance between two probability distributions
 - Provides a meaningful and smooth gradient everywhere
- **Critic vs Discriminator:**
 - WGAN uses a critic instead of a discriminator
 - Critic outputs a real number instead of a probability
 - Critic is trained to approximate the Wasserstein distance
- **Lipschitz Continuity:**
 - Critic must satisfy the Lipschitz continuity condition
 - Ensures the Wasserstein distance is well-defined
- **Weight Clipping:**
 - Original method to enforce Lipschitz continuity
 - Clip weights to a fixed range after each gradient update
- **Gradient Penalty (WGAN-GP):**
 - Improved method to enforce Lipschitz continuity
 - Adds a penalty term to the critic's loss

2 Mathematical Formulation

2.1 Wasserstein Distance

$$W(P_r, P_g) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P_r}[f(x)] - \mathbb{E}_{x \sim P_g}[f(x)]] \quad (1)$$

Where:

- P_r is the real data distribution
- P_g is the generated data distribution
- f is a 1-Lipschitz function

2.2 WGAN Objective

$$\min_G \max_C [\mathbb{E}_{x \sim P_r} [C(x)] - \mathbb{E}_{z \sim p(z)} [C(G(z))]] \quad (2)$$

Where:

- G is the generator
- C is the critic
- z is random noise

2.3 Gradient Penalty (WGAN-GP)

$$L = \mathbb{E}_{x \sim P_r} [C(x)] - \mathbb{E}_{x \sim P_g} [C(x)] + \lambda \cdot \mathbb{E}_{x \sim P_x} [(\|\nabla_x C(x)\|_2 - 1)^2] \quad (3)$$

Where:

- λ is the penalty coefficient (typically 10)
- P_x is the distribution of interpolated samples

3 WGAN Algorithm

Algorithm 1 WGAN Training

Require: α , the learning rate

Require: c , the clipping parameter

Require: m , the batch size

Require: n_{critic} , the number of critic iterations per generator iteration

```
1: for number of training iterations do
2:   for  $t = 1, \dots, n_{critic}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim P_r$  a batch from the real data
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples
5:      $\tilde{x}^{(i)} \leftarrow G_\theta(z^{(i)})$  for  $i = 1, \dots, m$ 
6:      $L^{(i)} \leftarrow C_w(\tilde{x}^{(i)}) - C_w(x^{(i)})$  for  $i = 1, \dots, m$ 
7:      $w \leftarrow w + \alpha \cdot \text{RMSPProp}(w, \nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)})$ 
8:      $w \leftarrow \text{clip}(w, -c, c)$ 
9:   end for
10:  Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples
11:   $\tilde{x}^{(i)} \leftarrow G_\theta(z^{(i)})$  for  $i = 1, \dots, m$ 
12:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSPProp}(\theta, -\nabla_\theta \frac{1}{m} \sum_{i=1}^m C_w(\tilde{x}^{(i)}))$ 
13: end for
```

4 PyTorch Implementation

4.1 Generator and Critic Classes

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4
5 class Generator(nn.Module):
6     def __init__(self, z_dim, img_dim):
7         super(Generator, self).__init__()
```

```

8         self.gen = nn.Sequential(
9             nn.Linear(z_dim, 128),
10            nn.LeakyReLU(0.2),
11            nn.Linear(128, 256),
12            nn.BatchNorm1d(256),
13            nn.LeakyReLU(0.2),
14            nn.Linear(256, img_dim),
15            nn.Tanh()
16        )
17
18    def forward(self, z):
19        return self.gen(z)
20
21 class Critic(nn.Module):
22     def __init__(self, img_dim):
23         super(Critic, self).__init__()
24         self.critic = nn.Sequential(
25             nn.Linear(img_dim, 128),
26             nn.LeakyReLU(0.2),
27             nn.Linear(128, 64),
28             nn.LeakyReLU(0.2),
29             nn.Linear(64, 1)
30         )
31
32     def forward(self, img):
33         return self.critic(img)

```

4.2 Training Loop

```

1  # Hyperparameters
2  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
3  learning_rate = 5e-5
4  batch_size = 64
5  img_dim = 28 * 28  # for MNIST
6  z_dim = 100
7  num_epochs = 50
8  critic_iterations = 5
9  weight_clip = 0.01
10
11 # Initialize generator and critic
12 generator = Generator(z_dim, img_dim).to(device)
13 critic = Critic(img_dim).to(device)
14
15 # Optimizers
16 opt_gen = optim.RMSprop(generator.parameters(), lr=learning_rate)
17 opt_critic = optim.RMSprop(critic.parameters(), lr=learning_rate)
18
19 # Training loop
20 for epoch in range(num_epochs):
21     for batch_idx, (real, _) in enumerate(loader):
22         real = real.view(-1, 784).to(device)
23         batch_size = real.shape[0]
24
25         # Train Critic
26         for _ in range(critic_iterations):
27             noise = torch.randn(batch_size, z_dim).to(device)
28             fake = generator(noise)

```

```

29         critic_real = critic(real).reshape(-1)
30         critic_fake = critic(fake).reshape(-1)
31         loss_critic = -(torch.mean(critic_real) - torch.mean(critic_fake))
32         critic.zero_grad()
33         loss_critic.backward(retain_graph=True)
34         opt_critic.step()
35
36         # Clip critic weights
37         for p in critic.parameters():
38             p.data.clamp_(-weight_clip, weight_clip)
39
40         # Train Generator
41         output = critic(fake).reshape(-1)
42         loss_gen = -torch.mean(output)
43         generator.zero_grad()
44         loss_gen.backward()
45         opt_gen.step()
46
47         if batch_idx % 100 == 0:
48             print(
49                 f"Epoch [{epoch}/{num_epochs}] Batch {batch_idx}/{len(loader)} \
50                 Loss D: {loss_critic:.4f}, loss G: {loss_gen:.4f}"
51             )

```

5 Key Differences from Standard GAN

- Uses Wasserstein distance instead of Jensen-Shannon divergence
- Critic outputs unbounded real numbers, not probabilities
- No log in the loss function
- Weight clipping or gradient penalty to enforce Lipschitz continuity
- More stable training and meaningful loss

6 Advantages of WGAN

- Improved stability during training
- Meaningful loss metric that correlates with sample quality
- Reduced mode collapse
- Better gradient flow for the generator
- Less sensitive to architecture choices and hyperparameters

7 Practical Tips

- Use RMSProp optimizer instead of Adam for more stable training
- Train the critic to optimality before each generator update
- Monitor the Wasserstein estimate during training
- Use gradient penalty (WGAN-GP) for better performance than weight clipping

- Experiment with different architectures for generator and critic
- Start with a lower learning rate (e.g., 5e-5) and adjust as needed
- Ensure proper normalization of input data (e.g., scale to [-1, 1])
- Use BatchNorm in the generator but not in the critic
- Experiment with different z_{dim} values for the noise input