# Exercise 5: Generative Adversarial Networks

## Erik Stolt

## What you did and how

Generative Adversarial Networks (GANs) work by incorporating two neural network architectures: a generator (G) and a discriminator (D). The generator learns to produce realistic images from an abstract input distribution, typically Gaussian noise. The discriminator learns to determine whether an image is real, meaning from the training dataset, or fake, meaning generated by G.

The discriminator's loss function is the binary cross-entropy (BCE) loss. Real images are labeled as 1 and fake images as 0. During training, the discriminator is presented with both real and fake images in each batch, along with their corresponding labels. It updates its weights to improve its ability to classify them correctly. The discriminator's total loss is computed as the average of the BCE loss on real and fake samples.

The generator is trained indirectly through the discriminator. To train the generator, fake images are generated and passed to the discriminator, but this time with the label 1, as if they were real. The generator's objective is to produce fake images that the discriminator classifies as real. When the discriminator is fooled by a fake image and predicts a value close to 1, the generator receives a lower loss. Therefore, the generator learns to generate fake images that closely resemble real images from the training data.

For this exercise, we defined a generator and a discriminator network, and trained them in an adversarial setup where the generator learned to produce increasingly realistic digit images while the discriminator tried to distinguish them from real MNIST digits. We monitored the training progress by generating sample images at different epochs by fixing the noise beforehand, which clearly showed how the image quality improved over time. In practice we first normalized the input images s.t. the pixel values transformed like $[0, 1] \rightarrow [-1, 1]$ which we also later denormalized for visualization according to image $=$ image $* 0.5 + 0.5$. Regarding the network structures for the generator and discriminator, we settled on two standard NN. The generator had tanh activation in the output layer to transform the values to the $[-1, 1]$ range. The discriminator instead used sigmoid activation since its task is a binary classification task. The networks were then trained for 30 epochs with a batch size of 32, monitoring the performances using Tensorboard.

## What results you obtained

Example print-out during training:

```
Epoch [28/30] | Loss D: 0.6806 | lr D: 5e-05 | Loss G: 0.7323 | lr G: 0.0001 | 135.27s
Epoch [29/30] | Loss D: 0.6682 | lr D: 5e-05 | Loss G: 0.7079 | lr G: 0.0001 | 122.10s
Epoch [30/30] | Loss D: 0.6654 | lr D: 5e-05 | Loss G: 0.7033 | lr G: 0.0001 | 126.16s
```

I also present the evolution of the digits and the loss as a function of epoch:
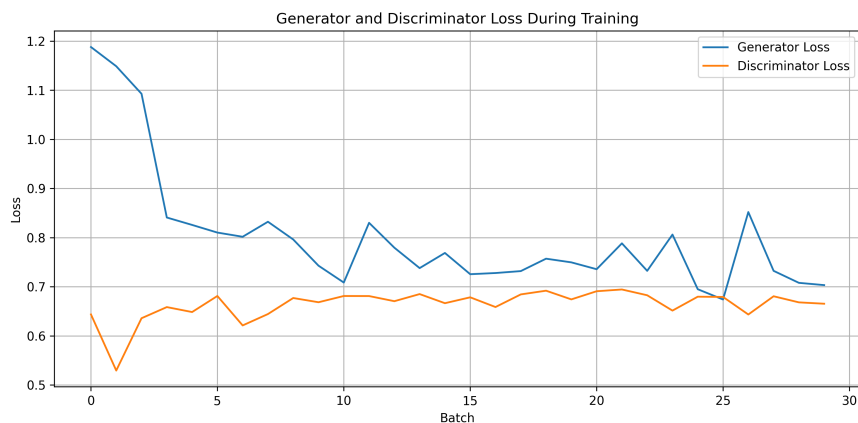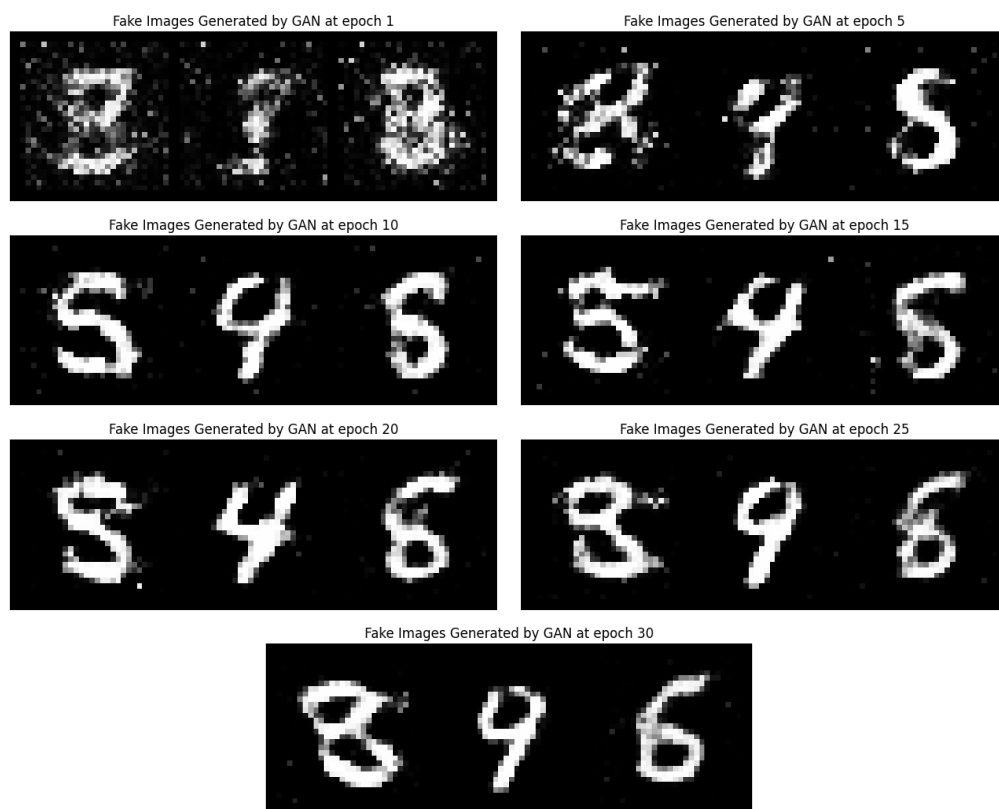
*Figure 1: Loss vs epoch for both G and D.*



*Figure 2: Evolution of G generating digits.*

We clearly see that the loss stays the same but the digits obviously get better. This is the result of a good trade-off between G and D improving at the same rate.

## What challenges you encountered and what could be improved

GANs are notoriously unstable which lead to the discriminant having zero loss, hence, no learning. This is fixed by setting different initial lrs for the two models. Furthermore, I also replaced the true labels ones and zeros with 0.9 and 0.1 (instead of 1 and 0) in order to trick D not to be overconfident. This works because it

penalizes the model in being overconfident (see Appendix 1). I also added dropout and batch normalization in order for G not to be deterministic as this also lead to model collapse. Batch normalization is a critical part of G but often left out in D since it can cause numerical intabilities [1]. To improve the model one could use WGANs to get even better results and prevent model collapse.

This exercise also taught me a lot about interpreting losses. In GANs, monitoring losses is not enough and even misleading. If one only were to monitor the loss, model collapse seems good since it happens when the discriminators loss is zero. However, we also have the two scenarios

- Discriminator loss may go to zero when it is dominating, but that means G is not learning.
- Generator loss may go up even when it is improving image quality because D is getting better.

This means that in adversarial setups, both networks are playing a game, so, their losses do not converge in the usual sense.

## Link to GitHub

Link: `https://github.com/erst6955/Advanced-Applied-Deep-Learning-in-Physics-And-Engineering`

## Appendix 1

BCE is given by

$$\text{BCE}(p, y) = -(y \cdot \log(p) + (1 - y) \cdot \log(1 - p)),$$

which means that if $p = 0.99$ and $y = 1$ we would get BCE $\approx 0.01$ which is very small loss, hence, the gradient would vanish. However, if we let $y \to 0.9$ we would instead get BCE $\approx 0.47$ which is higher! Therefore, more learning is possible with this small change. The same goes for label $y \to 0.1$. This is a common trick which works well when D is learning too strong.

## Code (for backup)

**~\OneDrive\Desktop\Exercise 5 - Generative Adversarial Networks\GANs for MNIST.py**

```python
 1
 2
 3   import torchvision
 4   from torchvision import transforms, datasets
 5   import torch
 6   import torch.nn as nn
 7   import torch.optim as optim
 8   from torch.utils.data import DataLoader
 9   from torch.utils.tensorboard import SummaryWriter
10   import matplotlib.pyplot as plt
11   import time
12   import random
13
14
15
16   # Hyperparameters
17   device = "cuda" if torch.cuda.is_available() else "cpu"
18   lr = 3e-4
19   batchSize = 32
20   logStep = 625
21   latent_dimension = 128
22   image_dimension = 784
23
24   myTransforms = transforms.Compose([ # we define a tranform that converts the image to
     tensor and normalizes it with mean and std of 0.5
25       transforms.ToTensor(),          # which will convert the image range from [0, 1] to
     [-1, 1]
26       transforms.Normalize((0.5,), (0.5,))
27   ])
28
29   dataset = datasets.MNIST(root="dataset/", transform=myTransforms, download=True)
30   loader = DataLoader(dataset, batch_size=batchSize, shuffle=True)
31
32
33   class Generator(nn.Module):
34       """
35       Generator Model which takes a random noise vector and generates a fake image. The
     input noise in latent space is an abstract space where each point corresponds to a
     different kind of image.
36       The noise vector doesn't directly mean anything in the beginning, but after training,
     it gets mapped by the generator to a specific kind of image — like a "7", or a "2", or a
     "3". We make
37       assumptions about how the input noise is distributed.
38       """
39       def __init__(self):
40           super().__init__()
41           self.gen = nn.Sequential( # simple NN
42               nn.Linear(latent_dimension, 256),
43               nn.BatchNorm1d(256),
44               nn.ReLU(),
45
46               nn.Linear(256, 512),
```

```python
47              nn.BatchNorm1d(512),
48              nn.ReLU(),
49
50              nn.Linear(512, 1024),
51              nn.BatchNorm1d(1024),
52              nn.ReLU(),
53
54              nn.Linear(1024, image_dimension),
55              nn.Tanh(), # tanh activation function to get the output in the range of [-1,
    1]
56          )
57
58      def forward(self, x):
59          return self.gen(x)
60
61  class Discriminator(nn.Module):
62      """
63      Discriminator Model which takes an image and outputs a probability of it being real or
    fake. Furthermore,
64      the discriminator is a binary classifier which uses the sigmoid activation and the BCE
    loss function.
65      """
66      def __init__(self):
67          super().__init__()
68          self.disc = nn.Sequential(
69              nn.Linear(image_dimension, 1024),
70              nn.LeakyReLU(0.2),
71              nn.Dropout(0.3),  # add this
72              nn.Linear(1024, 512),
73              nn.LeakyReLU(0.2),
74              nn.Dropout(0.3), # add this
75              nn.Linear(512, 256),
76              nn.LeakyReLU(0.2),
77              nn.Linear(256, 1),
78              nn.Sigmoid(),
79          )
80
81
82      def forward(self, x):
83          return self.disc(x)
84
85
86  def train_models(discriminator, generator, fixed_noise, num_examples, opt_discriminator,
    opt_generator, criterion, epochs, writer):
87      """
88      Trains the GAN models using the MNIST dataset.
89      The generator uses BCE where 1 is the label for fake images and 0 is the label for
    real images. This mean that it learns
90      to generate fake images that are similar to the real images in the dataset since the
    disciminator has real images labeled as
91      1 and fake images labeled as 0. I.e., the generator aims to minimize the loss by
    producing images that are classified as real by the discriminator.
92      The disciminator on the other hand is simply trained to identify the real and fake
    images.
93      """
```

```python
 94        step = 0 # for tensorboard logging
 95        gen_losses = []
 96        disc_losses = []
 97
 98        # ===================== Training Loop =======================================
 99        for epoch in range(epochs):
100
101            start_time = time.time()
102
103            for batch_idx, (real, _) in enumerate(loader):
104                real = real.view(-1, image_dimension).to(device) # flatten the image in order
       to pass it to the discriminator
105                batch_size = real.shape[0] # get the batch size
106
107                noise = torch.randn(batch_size, latent_dimension).to(device) # generate random
       noise from a normal distribution
108                fake = generator(noise) # generate fake images from the noise
109
110                # Discriminator Loss
111                disc_real = discriminator(real).view(-1) # get the discriminator output for
       real images
112                loss_real = criterion(disc_real, torch.full_like(disc_real, 0.9)) # pass
       through BCE where real images are labeled as 1
113
114                disc_fake = discriminator(fake.detach()).view(-1) # get the discriminator
       output for fake images
115                loss_fake = criterion(disc_fake, torch.full_like(disc_fake, 0.1)) # pass
       through BCE where fake images are labeled as 0
116
117                loss_discriminator = (loss_real + loss_fake) / 2 # average the loss for real
       and fake images
118
119                discriminator.zero_grad() # zero the gradients
120                loss_discriminator.backward(retain_graph=True) # backpropagate the loss
121                opt_discriminator.step() # update the discriminator weights
122
123                # ===== Generator Loss =====
124                # The generator tries to fool the discriminator, so we want the discriminator
       to think that the fake images are real
125                output = discriminator(fake).view(-1) # get the discriminator output for fake
       images
126                loss_generator = criterion(output, torch.ones_like(output)) # pass through BCE
       where fake images are labeled as 1
127                # The generator tries to maximize the probability of the discriminator being
       wrong
128
129                generator.zero_grad()
130                loss_generator.backward()
131                opt_generator.step()
132
133                if batch_idx % logStep == 0: # tensorboard logging
134                    with torch.no_grad():
135                        fake_images = generator(fixed_noise).reshape(-1, 1, 28, 28) # reshape
       the fake images to 1 channel and 28x28
```

```python
136                    real_images = real.reshape(-1, 1, 28, 28) # reshape the real images to
       1 channel and 28x28

137
138                    imgGridFake = torchvision.utils.make_grid(fake_images, normalize=True)
       # make a grid of fake images
139                    imgGridReal = torchvision.utils.make_grid(real_images, normalize=True)
       # make a grid of real images

140
141                    # Denormalize: [-1, 1] → [0, 1]
142                    fake_images = denormalize(fake_images)
143                    real_images = denormalize(fake_images)

144
145                    writer.add_image("MNIST Fake Images", imgGridFake, global_step=step)
146                    writer.add_image("MNIST Real Images", imgGridReal, global_step=step)
147                    writer.add_scalar("Loss Discriminator", loss_discriminator.item(),
       step)
148                    writer.add_scalar("Loss Generator", loss_generator.item(), step)

149
150                    step += 1
151          for param_group in opt_generator.param_groups:
152              lr_g = param_group['lr']
153          for param_group in opt_discriminator.param_groups:
154              lr_d = param_group['lr']

155
156          # At the end of each epoch
157          gen_losses.append(loss_generator.item())
158          disc_losses.append(loss_discriminator.item())
159          elapsed = time.time() - start_time
160          print(f"Epoch [{epoch+1}/{epochs}] | Loss D: {loss_discriminator:.4f} | lr D:
       {lr_d} | Loss G: {loss_generator:.4f} | lr G: {lr_g} | {elapsed:.2f}s")

161
162          get_fake_images(generator, fixed_noise, latent_dimension, num_examples, epoch) #
       get some fake images for the epoch

163
164      return gen_losses, disc_losses

165
166
167  def denormalize(images):
168      return images * 0.5 + 0.5

169
170
171
172  def get_fake_images(generator, fixed_noise, latent_dimension, num_examples, epoch):
173      """
174      Returns a batch of fake images from the generator.
175      """

176
177      generator.eval()

178
179      # Generate images from noise
180      with torch.no_grad():
181          generated_images = generator(fixed_noise).reshape(-1, 1, 28, 28)
182          generated_images = denormalize(generated_images) # denormalize the images

183
```

```python
184        # Create a grid for visualization
185        grid = torchvision.utils.make_grid(generated_images.cpu(), nrow=4, normalize=True)
186
187        # Plot the generated images
188        plt.figure(figsize=(8, 8))
189        plt.title(f"Fake Images Generated by GAN at epoch {epoch+1}")
190        plt.axis("off")
191        plt.imshow(grid.permute(1, 2, 0).squeeze())
192        plt.savefig(f"Figures/fake_images_{epoch+1}.png", bbox_inches='tight')
193        plt.close()
194
195
196
197  def plot_losses(gen_losses, disc_losses):
198        plt.figure(figsize=(10, 5))
199        plt.plot(gen_losses, label="Generator Loss")
200        plt.plot(disc_losses, label="Discriminator Loss")
201        plt.xlabel("Batch")
202        plt.ylabel("Loss")
203        plt.title("Generator and Discriminator Loss During Training")
204        plt.legend()
205        plt.grid(True)
206        plt.tight_layout()
207        plt.savefig("Figures/loss_plot.png", dpi=300)
208
209
210
211  # ========================================= Run Code
     =========================================================
212  discriminator = Discriminator().to(device)
213  generator = Generator().to(device)
214
215  lr_G = 1e-4
216  lr_D = 5e-5
217
218  opt_generator = optim.Adam(generator.parameters(), lr=lr_G, betas=(0.5, 0.999))
219  opt_discriminator = optim.Adam(discriminator.parameters(), lr=lr_D, betas=(0.5, 0.999))
220  criterion = nn.BCELoss() # Binary Cross Entropy Loss
221
222  epochs = 30
223  num_examples = 3
224
225  fixed_noise = torch.randn(num_examples, latent_dimension).to(device) # fixed noise for
     reproducibility after every epoch
226  writer = SummaryWriter("runs/GAN_MNIST")
227
228  gen_losses, disc_losses = train_models(discriminator, generator, fixed_noise,
     num_examples, opt_discriminator, opt_generator, criterion, epochs, writer)
229  plot_losses(gen_losses, disc_losses)
230
231  writer.close()
232  print("Complete!")
233
234
```

```
235
236
237
238
239
240
```

# References

[1]  Alec Radford, Luke Metz, and Soumith Chintala. "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks". In: *arXiv preprint arXiv:1511.06434* (2016). URL: https://arxiv.org/abs/1511.06434.