

Generative Adversarial Networks (GANs)

So far in CS231N, all the applications of neural networks that we have explored have been **discriminative models** that take an input and are trained to produce a labeled output. This has ranged from straightforward classification of image categories to sentence generation (which was still phrased as a classification problem, our labels were in vocabulary space and we'd learned a recurrence to capture multi-word labels). In this notebook, we will expand our repertoire, and build **generative models** using neural networks. Specifically, we will learn how to build models which generate novel images that resemble a set of training images.

What is a GAN?

In 2014, [Goodfellow et al. \(https://arxiv.org/abs/1406.2661\)](https://arxiv.org/abs/1406.2661) presented a method for training generative models called Generative Adversarial Networks (GANs for short). In a GAN, we build two different neural networks. Our first network is a traditional classification network, called the **discriminator**. We will train the discriminator to take images, and classify them as being real (belonging to the training set) or fake (not present in the training set). Our other network, called the **generator**, will take random noise as input and transform it using a neural network to produce images. The goal of the generator is to fool the discriminator into thinking the images it produced are real.

We can think of this back and forth process of the generator (G) trying to fool the discriminator (D), and the discriminator trying to correctly classify real vs. fake as a minimax game:

$$\underset{G}{\text{minimize}} \underset{D}{\text{maximize}} \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

where $z \sim p(z)$ are the random noise samples, $G(z)$ are the generated images using the neural network generator G , and D is the output of the discriminator, specifying the probability of an input being real. In [Goodfellow et al. \(https://arxiv.org/abs/1406.2661\)](https://arxiv.org/abs/1406.2661), they analyze this minimax game and show how it relates to minimizing the Jensen-Shannon divergence between the training data distribution and the generated samples from G .

To optimize this minimax game, we will alternate between taking gradient *descent* steps on the objective for G , and gradient *ascent* steps on the objective for D :

1. update the **generator** (G) to minimize the probability of the **discriminator making the correct choice**.
2. update the **discriminator** (D) to maximize the probability of the **discriminator making the correct choice**.

While these updates are useful for analysis, they do not perform well in practice. Instead, we will use a different objective when we update the generator: maximize the probability of the **discriminator making the incorrect choice**. This small change helps to alleviate problems with the generator gradient vanishing when the discriminator is confident. This is the standard update used in most GAN papers, and was used in the original paper from [Goodfellow et al. \(https://arxiv.org/abs/1406.2661\)](https://arxiv.org/abs/1406.2661).

In this assignment, we will alternate the following updates:

1. Update the generator (G) to maximize the probability of the discriminator making the incorrect choice on generated data:

$$\underset{G}{\text{maximize}} \mathbb{E}_{z \sim p(z)} [\log D(G(z))]$$

2. Update the discriminator (D), to maximize the probability of the discriminator making the correct choice on real and generated data:

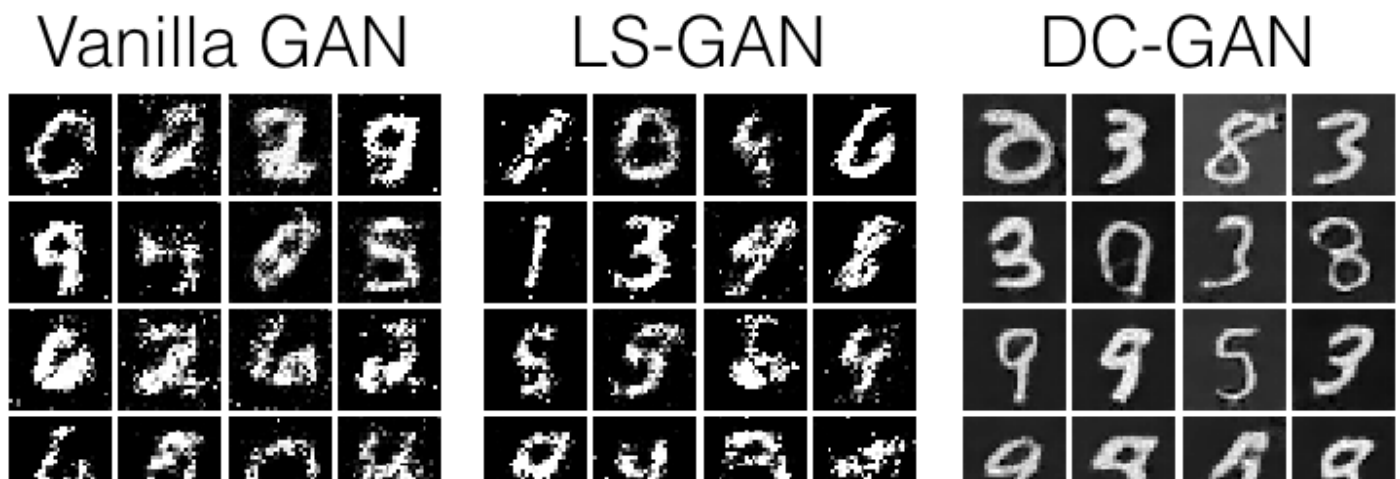
$$\underset{D}{\text{maximize}} \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

What else is there?

Since 2014, GANs have exploded into a huge research area, with massive [workshops](https://sites.google.com/site/nips2016adversarial/) (<https://sites.google.com/site/nips2016adversarial/>), and [hundreds of new papers](https://github.com/hindupuravinash/the-gan-zoo) (<https://github.com/hindupuravinash/the-gan-zoo>). Compared to other approaches for generative models, they often produce the highest quality samples but are some of the most difficult and finicky models to train (see [this github repo](https://github.com/soumith/ganhacks) (<https://github.com/soumith/ganhacks>) that contains a set of 17 hacks that are useful for getting models working). Improving the stability and robustness of GAN training is an open research question, with new papers coming out every day! For a more recent tutorial on GANs, see [here](https://arxiv.org/abs/1701.00160) (<https://arxiv.org/abs/1701.00160>). There is also some even more recent exciting work that changes the objective function to Wasserstein distance and yields much more stable results across model architectures: [WGAN](https://arxiv.org/abs/1701.07875) (<https://arxiv.org/abs/1701.07875>), [WGAN-GP](https://arxiv.org/abs/1704.00028) (<https://arxiv.org/abs/1704.00028>).

GANs are not the only way to train a generative model! For other approaches to generative modeling check out the [deep generative model chapter](http://www.deeplearningbook.org/contents/generative_models.html) (http://www.deeplearningbook.org/contents/generative_models.html) of the Deep Learning [book](http://www.deeplearningbook.org) (<http://www.deeplearningbook.org>). Another popular way of training neural networks as generative models is Variational Autoencoders (co-discovered [here](https://arxiv.org/abs/1312.6114) (<https://arxiv.org/abs/1312.6114>) and [here](https://arxiv.org/abs/1401.4082) (<https://arxiv.org/abs/1401.4082>)). Variational autoencoders combine neural networks with variational inference to train deep generative models. These models tend to be far more stable and easier to train but currently don't produce samples that are as pretty as GANs.

Here's an example of what your outputs from the 3 different models you're going to train should look like... note that GANs are sometimes finicky, so your outputs might not look exactly like this... this is just meant to be a *rough* guideline of the kind of quality you can expect:



Setup

```

In [1]: import torch
import torch.nn as nn
from torch.nn import init
import torchvision
import torchvision.transforms as T
import torch.optim as optim
from torch.utils.data import DataLoader
from torch.utils.data import sampler
import torchvision.datasets as dset

import numpy as np

import matplotlib.pyplot as plt
import matplotlib.gridspec as gridspec

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

def show_images(images):
    images = np.reshape(images, [images.shape[0], -1]) # images reshape
    to (batch_size, D)
    sqrtn = int(np.ceil(np.sqrt(images.shape[0])))
    sqrtimg = int(np.ceil(np.sqrt(images.shape[1])))

    fig = plt.figure(figsize=(sqrtn, sqrtn))
    gs = gridspec.GridSpec(sqrtn, sqrtn)
    gs.update(wspace=0.05, hspace=0.05)

    for i, img in enumerate(images):
        ax = plt.subplot(gs[i])
        plt.axis('off')
        ax.set_xticklabels([])
        ax.set_yticklabels([])
        ax.set_aspect('equal')
        plt.imshow(img.reshape([sqrtimg, sqrtimg]))
    return

def preprocess_img(x):
    return 2 * x - 1.0

def deprocess_img(x):
    return (x + 1.0) / 2.0

def rel_error(x, y):
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y)
    )))

def count_params(model):
    """Count the number of parameters in the current TensorFlow graph
    """
    param_count = np.sum([np.prod(p.size()) for p in model.parameters
    ()])
    return param_count

```

```
answers = dict(np.load('gan-checks-tf.npz'))
```

Dataset

GANs are notoriously finicky with hyperparameters, and also require many training epochs. In order to make this assignment approachable without a GPU, we will be working on the MNIST dataset, which is 60,000 training and 10,000 test images. Each picture contains a centered image of white digit on black background (0 through 9). This was one of the first datasets used to train convolutional neural networks and it is fairly easy -- a standard CNN model can easily exceed 99% accuracy.

To simplify our code here, we will use the PyTorch MNIST wrapper, which downloads and loads the MNIST dataset. See the [documentation \(https://github.com/pytorch/vision/blob/master/torchvision/datasets/mnist.py\)](https://github.com/pytorch/vision/blob/master/torchvision/datasets/mnist.py) for more information about the interface. The default parameters will take 5,000 of the training examples and place them into a validation dataset. The data will be saved into a folder called `MNIST_data`.

```

In [2]: class ChunkSampler(sampler.Sampler):
        """Samples elements sequentially from some offset.
        Arguments:
            num_samples: # of desired datapoints
            start: offset where we should start selecting from
        """
        def __init__(self, num_samples, start=0):
            self.num_samples = num_samples
            self.start = start

        def __iter__(self):
            return iter(range(self.start, self.start + self.num_samples))

        def __len__(self):
            return self.num_samples

NUM_TRAIN = 50000
NUM_VAL = 5000

NOISE_DIM = 96
batch_size = 128

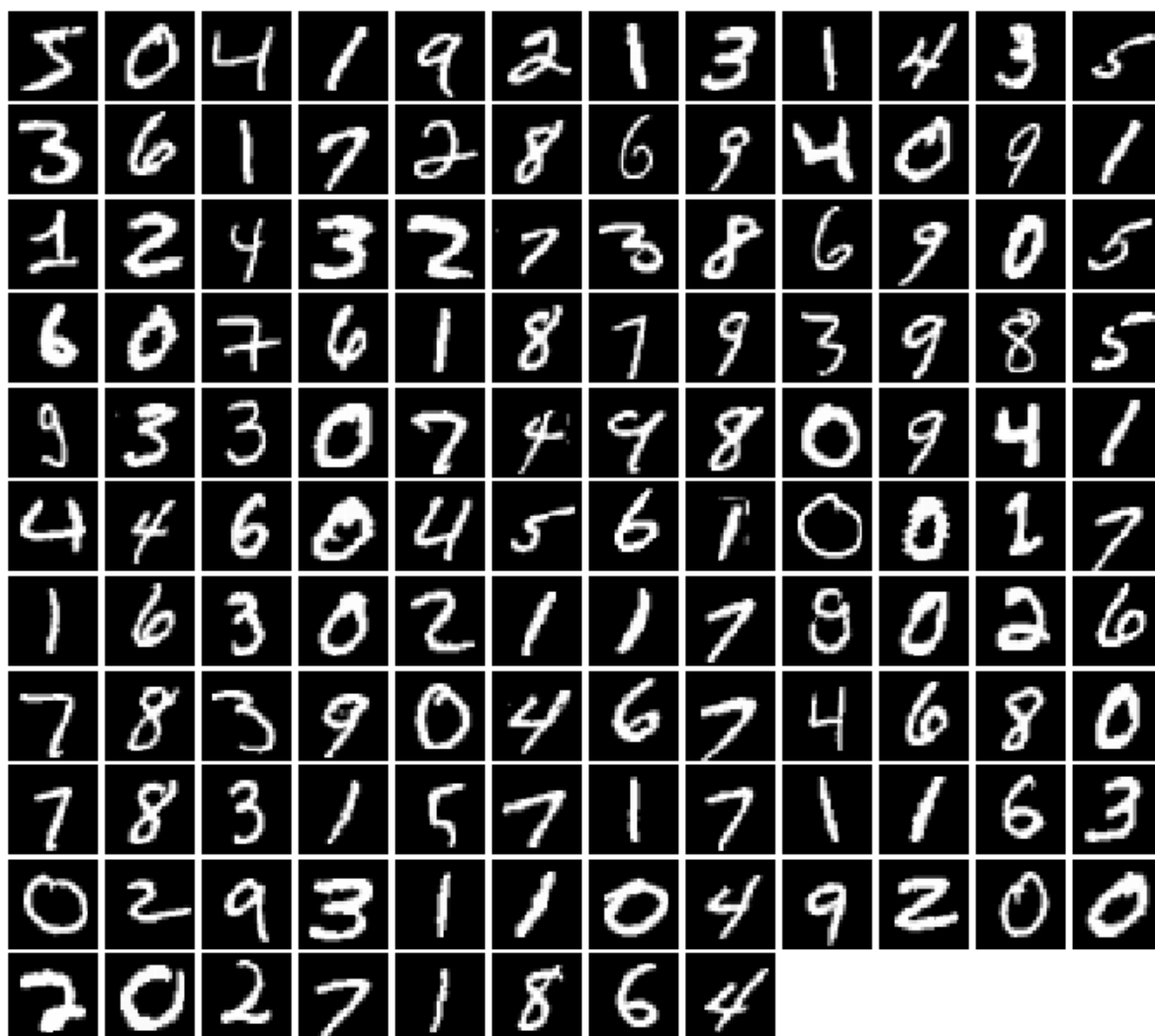
mnist_train = dset.MNIST('./cs231n/datasets/MNIST_data', train=True, download=True,
                        transform=T.ToTensor())
loader_train = DataLoader(mnist_train, batch_size=batch_size,
                        sampler=ChunkSampler(NUM_TRAIN, 0))

mnist_val = dset.MNIST('./cs231n/datasets/MNIST_data', train=True, download=True,
                        transform=T.ToTensor())
loader_val = DataLoader(mnist_val, batch_size=batch_size,
                        sampler=ChunkSampler(NUM_VAL, NUM_TRAIN))

imgs = loader_train.__iter__().next()[0].view(batch_size, 784).numpy().squeeze()
show_images(imgs)

```

Downloading <http://yann.lecun.com/exdb/mnist/train-images-idx3-ubyte.gz>
 Downloading <http://yann.lecun.com/exdb/mnist/train-labels-idx1-ubyte.gz>
 Downloading <http://yann.lecun.com/exdb/mnist/t10k-images-idx3-ubyte.gz>
 Downloading <http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-ubyte.gz>
 Processing...
 Done!



Random Noise

Generate uniform noise from -1 to 1 with shape `[batch_size, dim]`.

Hint: use `torch.rand`.

```
In [3]: def sample_noise(batch_size, dim):
        """
        Generate a PyTorch Tensor of uniform random noise.

        Input:
        - batch_size: Integer giving the batch size of noise to generate.
        - dim: Integer giving the dimension of noise to generate.

        Output:
        - A PyTorch Tensor of shape (batch_size, dim) containing uniform
          random noise in the range (-1, 1).
        """
        return 2*torch.rand(batch_size, dim) - 1
```

Make sure noise is the correct shape and type:

```
In [4]: def test_sample_noise():
        batch_size = 3
        dim = 4
        torch.manual_seed(231)
        z = sample_noise(batch_size, dim)
        np_z = z.cpu().numpy()
        assert np_z.shape == (batch_size, dim)
        assert torch.is_tensor(z)
        assert np.all(np_z >= -1.0) and np.all(np_z <= 1.0)
        assert np.any(np_z < 0.0) and np.any(np_z > 0.0)
        print('All tests passed!')

        test_sample_noise()

All tests passed!
```

Flatten

Recall our Flatten operation from previous notebooks... this time we also provide an Unflatten, which you might want to use when implementing the convolutional generator. We also provide a weight initializer (and call it for you) that uses Xavier initialization instead of PyTorch's uniform default.


```
In [5]: class Flatten(nn.Module):
        def forward(self, x):
            N, C, H, W = x.size() # read in N, C, H, W
            return x.view(N, -1) # "flatten" the C * H * W values into a single vector per image

        class Unflatten(nn.Module):
            """
            An Unflatten module receives an input of shape (N, C*H*W) and reshapes it
            to produce an output of shape (N, C, H, W).
            """
            def __init__(self, N=-1, C=128, H=7, W=7):
                super(Unflatten, self).__init__()
                self.N = N
                self.C = C
                self.H = H
                self.W = W
            def forward(self, x):
                return x.view(self.N, self.C, self.H, self.W)

        def initialize_weights(m):
            if isinstance(m, nn.Linear) or isinstance(m, nn.ConvTranspose2d):
                init.xavier_uniform_(m.weight.data)
```

CPU / GPU

By default all code will run on CPU. GPUs are not needed for this assignment, but will help you to train your models faster. If you do want to run the code on a GPU, then change the `dtype` variable in the following cell.

```
In [6]: dtype = torch.FloatTensor
        #dtype = torch.cuda.FloatTensor ## UNCOMMENT THIS LINE IF YOU'RE ON A GPU!
```

Discriminator

Our first step is to build a discriminator. Fill in the architecture as part of the `nn.Sequential` constructor in the function below. All fully connected layers should include bias terms. The architecture is:

- Fully connected layer with input size 784 and output size 256
- LeakyReLU with alpha 0.01
- Fully connected layer with input_size 256 and output size 256
- LeakyReLU with alpha 0.01
- Fully connected layer with input size 256 and output size 1

Recall that the Leaky ReLU nonlinearity computes $f(x) = \max(\alpha x, x)$ for some fixed constant α ; for the LeakyReLU nonlinearities in the architecture above we set $\alpha = 0.01$.

The output of the discriminator should have shape `[batch_size, 1]`, and contain real numbers corresponding to the scores that each of the `batch_size` inputs is a real image.

```
In [62]: def discriminator():
    """
    Build and return a PyTorch model implementing the architecture above.
    """
    model = nn.Sequential(
        Flatten(),
        nn.Linear(784, 256),
        nn.LeakyReLU(negative_slope=0.01),
        nn.Linear(256, 256),
        nn.LeakyReLU(negative_slope=0.01),
        nn.Linear(256, 1)
    )
    return model
```

Test to make sure the number of parameters in the discriminator is correct:

```
In [63]: def test_discriminator(true_count=267009):
    model = discriminator()
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in discriminator. Check your architecture.')
    else:
        print('Correct number of parameters in discriminator.')

test_discriminator()

Correct number of parameters in discriminator.
```

Generator

Now to build the generator network:

- Fully connected layer from noise_dim to 1024
- ReLU
- Fully connected layer with size 1024
- ReLU
- Fully connected layer with size 784
- TanH (to clip the image to be in the range of [-1,1])

```
In [64]: def generator(noise_dim=NOISE_DIM):
        """
        Build and return a PyTorch model implementing the architecture above.
        """
        model = nn.Sequential(
            nn.Linear(noise_dim, 1024),
            nn.ReLU(),
            nn.Linear(1024, 1024),
            nn.ReLU(),
            nn.Linear(1024, 784),
            nn.Tanh()
        )
        return model
```

Test to make sure the number of parameters in the generator is correct:

```
In [65]: def test_generator(true_count=1858320):
        model = generator(4)
        cur_count = count_params(model)
        if cur_count != true_count:
            print('Incorrect number of parameters in generator. Check your architecture.')
        else:
            print('Correct number of parameters in generator.')

        test_generator()
```

Correct number of parameters in generator.

GAN Loss

Compute the generator and discriminator loss. The generator loss is:

$$\ell_G = -\mathbb{E}_{z \sim p(z)} [\log D(G(z))]$$

and the discriminator loss is:

$$\ell_D = -\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] - \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

Note that these are negated from the equations presented earlier as we will be *minimizing* these losses.

HINTS: You should use the `bce_loss` function defined below to compute the binary cross entropy loss which is needed to compute the log probability of the true label given the logits output from the discriminator. Given a score $s \in \mathbb{R}$ and a label $y \in \{0, 1\}$, the binary cross entropy loss is

$$bce(s, y) = -y * \log(s) - (1 - y) * \log(1 - s)$$

A naive implementation of this formula can be numerically unstable, so we have provided a numerically stable implementation for you below.

You will also need to compute labels corresponding to real or fake and use the `logit` arguments to determine their size. Make sure you cast these labels to the correct data type using the global `dtype` variable, for example:

```
true_labels = torch.ones(size).type(dtype)
```

Instead of computing the expectation of $\log D(G(z))$, $\log D(x)$ and $\log(1 - D(G(z)))$, we will be averaging over elements of the minibatch. so make sure to combine the loss by averaging instead of summing.

```
In [66]: def bce_loss(input, target):
    """
    Numerically stable version of the binary cross-entropy loss function.

    As per https://github.com/pytorch/pytorch/issues/751
    See the TensorFlow docs for a derivation of this formula:
    https://www.tensorflow.org/api_docs/python/tf/nn/sigmoid_cross_entropy_with_logits

    Inputs:
    - input: PyTorch Tensor of shape (N, ) giving scores.
    - target: PyTorch Tensor of shape (N,) containing 0 and 1 giving targets.

    Returns:
    - A PyTorch Tensor containing the mean BCE loss over the minibatch of input data.
    """
    neg_abs = - input.abs()
    loss = input.clamp(min=0) - input * target + (1 + neg_abs.exp()).log()
    return loss.mean()
```

```
In [67]: def discriminator_loss(logits_real, logits_fake):
        """
        Computes the discriminator loss described above.

        Inputs:
        - logits_real: PyTorch Tensor of shape (N,) giving scores for the real data.
        - logits_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

        Returns:
        - loss: PyTorch Tensor containing (scalar) the loss for the discriminator.
        """
        loss = bce_loss(logits_real, torch.ones(logits_real.shape).type(dtype))
        loss += bce_loss(logits_fake, torch.zeros(logits_fake.shape).type(dtype))
        return loss

def generator_loss(logits_fake):
    """
    Computes the generator loss described above.

    Inputs:
    - logits_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

    Returns:
    - loss: PyTorch Tensor containing the (scalar) loss for the generator.
    """
    loss = bce_loss(logits_fake, torch.ones(logits_fake.shape).type(dtype))
    return loss
```

Test your generator and discriminator loss. You should see errors $< 1e-7$.

```
In [68]: def test_discriminator_loss(logits_real, logits_fake, d_loss_true):
        d_loss = discriminator_loss(torch.Tensor(logits_real).type(dtype),
                                     torch.Tensor(logits_fake).type(dtype)).cpu().numpy()
        print("Maximum error in d_loss: %g"%rel_error(d_loss_true, d_loss))

test_discriminator_loss(answers['logits_real'], answers['logits_fake'],
                        answers['d_loss_true'])

Maximum error in d_loss: 3.97058e-09
```

```
In [69]: def test_generator_loss(logits_fake, g_loss_true):
          g_loss = generator_loss(torch.Tensor(logits_fake).type(dtype)).cpu().numpy()
          print("Maximum error in g_loss: %g"%rel_error(g_loss_true, g_loss))

          test_generator_loss(answers['logits_fake'], answers['g_loss_true'])
```

Maximum error in g_loss: 4.4518e-09

Optimizing our loss

Make a function that returns an `optim.Adam` optimizer for the given model with a `1e-3` learning rate, `beta1=0.5`, `beta2=0.999`. You'll use this to construct optimizers for the generators and discriminators for the rest of the notebook.

```
In [70]: def get_optimizer(model):
          """
          Construct and return an Adam optimizer for the model with learning rate
          1e-3,
          beta1=0.5, and beta2=0.999.

          Input:
          - model: A PyTorch model that we want to optimize.

          Returns:
          - An Adam optimizer for the model with the desired hyperparameters.
          """
          optimizer = optim.Adam(model.parameters(), lr=1e-3, betas=(0.5, 0.999))
          return optimizer
```

Training a GAN!

We provide you the main training loop... you won't need to change this function, but we encourage you to read through and understand it.

```

In [71]: def run_a_gan(D, G, D_solver, G_solver, discriminator_loss, generator_loss,
                show_every=250,
                batch_size=128, noise_size=96, num_epochs=10):
    """
    Train a GAN!

    Inputs:
    - D, G: PyTorch models for the discriminator and generator
    - D_solver, G_solver: torch.optim Optimizers to use for training the
      discriminator and generator.
    - discriminator_loss, generator_loss: Functions to use for computing
      the generator and
      discriminator loss, respectively.
    - show_every: Show samples after every show_every iterations.
    - batch_size: Batch size to use for training.
    - noise_size: Dimension of the noise to use as input to the generator.
    - num_epochs: Number of epochs over the training dataset to use for
      training.
    """
    iter_count = 0
    for epoch in range(num_epochs):
        for x, _ in loader_train:
            if len(x) != batch_size:
                continue
            D_solver.zero_grad()
            real_data = x.type(dtype)
            logits_real = D(2* (real_data - 0.5)).type(dtype)

            g_fake_seed = sample_noise(batch_size, noise_size).type(dtype)
e)
            fake_images = G(g_fake_seed).detach()
            logits_fake = D(fake_images.view(batch_size, 1, 28, 28))

            d_total_error = discriminator_loss(logits_real, logits_fake)
            d_total_error.backward()
            D_solver.step()

            G_solver.zero_grad()
            g_fake_seed = sample_noise(batch_size, noise_size).type(dtype)
e)
            fake_images = G(g_fake_seed)

            gen_logits_fake = D(fake_images.view(batch_size, 1, 28, 28))
            g_error = generator_loss(gen_logits_fake)
            g_error.backward()
            G_solver.step()

            if (iter_count % show_every == 0):
                print('Iter: {}, D: {:.4}, G:{:.4}'.format(iter_count, d_total_error.item(), g_error.item()))
                imgs_numpy = fake_images.data.cpu().numpy()
                show_images(imgs_numpy[0:16])
                plt.show()
                print()
                iter_count += 1

```

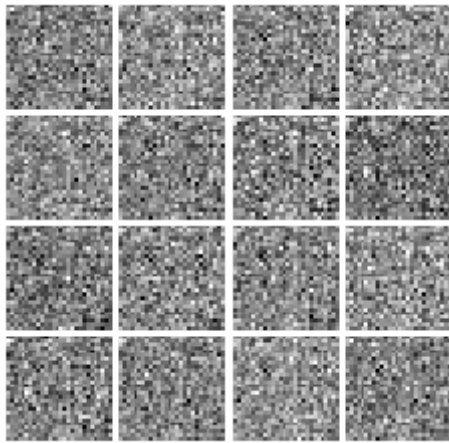


```
In [72]: # Make the discriminator
D = discriminator().type(dtype)

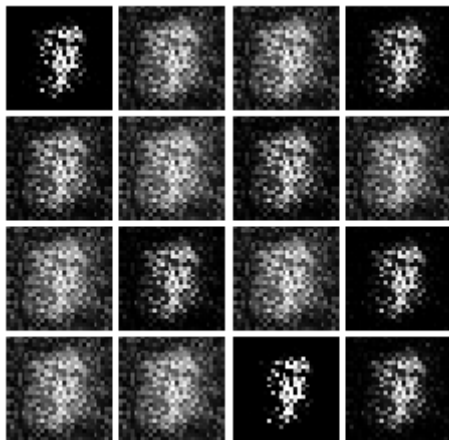
# Make the generator
G = generator().type(dtype)

# Use the function you wrote earlier to get optimizers for the Discriminator and the Generator
D_solver = get_optimizer(D)
G_solver = get_optimizer(G)
# Run it!
run_a_gan(D, G, D_solver, G_solver, discriminator_loss, generator_loss)
```

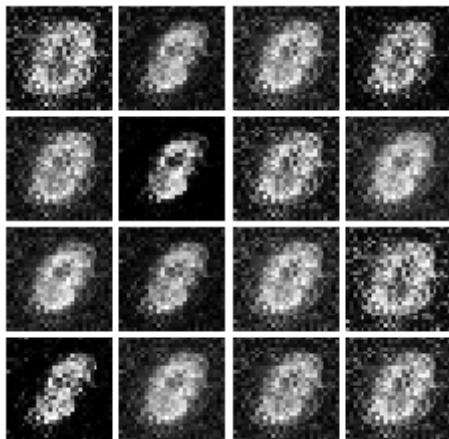
Iter: 0, D: 1.36, G:0.6958



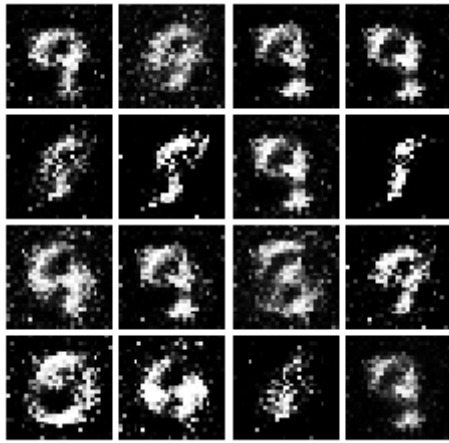
Iter: 250, D: 1.382, G:0.8853



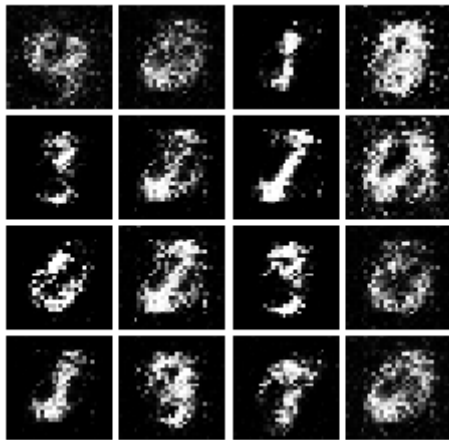
Iter: 500, D: 1.678, G:1.451



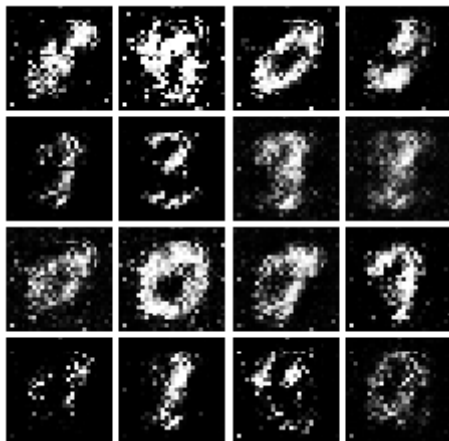
Iter: 750, D: 1.195, G:1.42



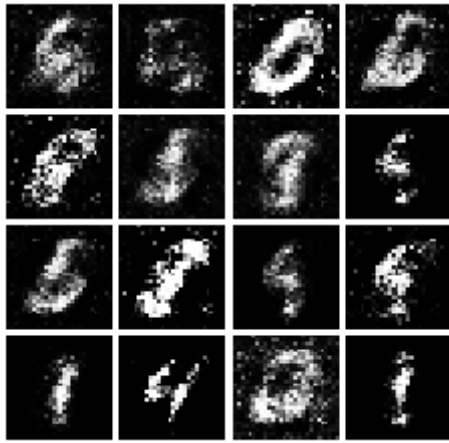
Iter: 1000, D: 1.205, G:1.162



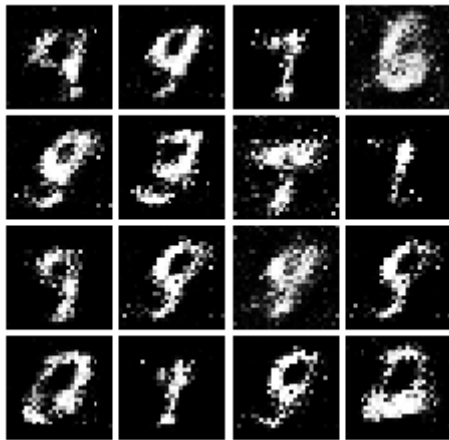
Iter: 1250, D: 1.402, G:1.158



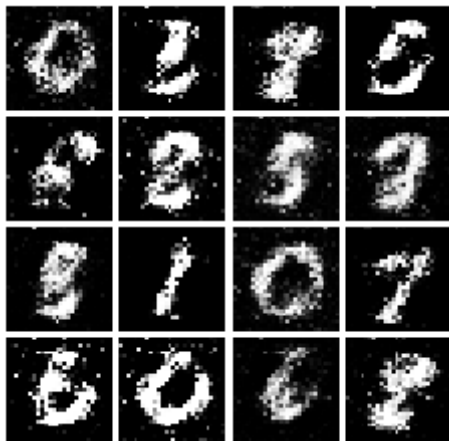
Iter: 1500, D: 1.185, G:1.01



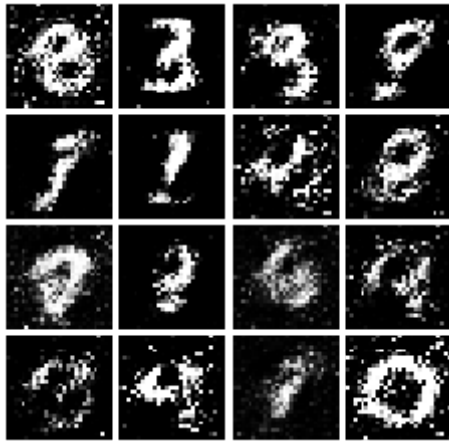
Iter: 1750, D: 1.23, G:0.7494



Iter: 2000, D: 1.42, G:0.7609



Iter: 2250, D: 1.307, G:0.7953



Iter: 2500, D: 1.422, G:0.7597



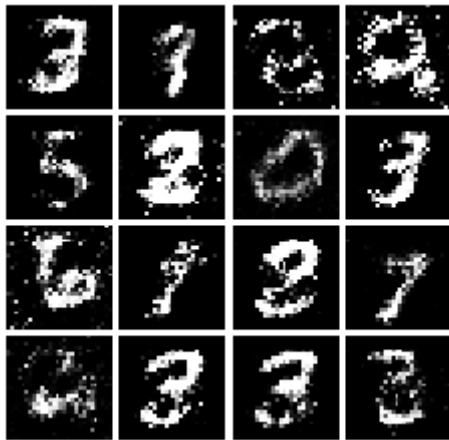
Iter: 2750, D: 1.379, G:0.8201



Iter: 3000, D: 1.354, G:0.774



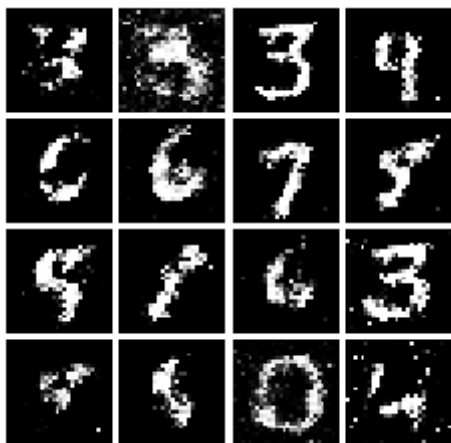
Iter: 3250, D: 1.285, G:0.929



Iter: 3500, D: 1.292, G:0.9089



Iter: 3750, D: 1.361, G:0.7867



Well that wasn't so hard, was it? In the iterations in the low 100s you should see black backgrounds, fuzzy shapes as you approach iteration 1000, and decent shapes, about half of which will be sharp and clearly recognizable as we pass 3000.

Least Squares GAN

We'll now look at Least Squares GAN (<https://arxiv.org/abs/1611.04076>), a newer, more stable alternative to the original GAN loss function. For this part, all we have to do is change the loss function and retrain the model. We'll implement equation (9) in the paper, with the generator loss:

$$\ell_G = \frac{1}{2} \mathbb{E}_{z \sim p(z)} [(D(G(z)) - 1)^2]$$

and the discriminator loss:

$$\ell_D = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} [(D(x) - 1)^2] + \frac{1}{2} \mathbb{E}_{z \sim p(z)} [(D(G(z)))^2]$$

HINTS: Instead of computing the expectation, we will be averaging over elements of the minibatch, so make sure to combine the loss by averaging instead of summing. When plugging in for $D(x)$ and $D(G(z))$ use the direct output from the discriminator (`scores_real` and `scores_fake`).

```
In [73]: def ls_discriminator_loss(scores_real, scores_fake):
        """
        Compute the Least-Squares GAN loss for the discriminator.

        Inputs:
        - scores_real: PyTorch Tensor of shape (N,) giving scores for the real data.
        - scores_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

        Outputs:
        - loss: A PyTorch Tensor containing the loss.
        """
        loss = torch.mean((scores_real - 1)**2)/2 + torch.mean((scores_fake)**2)/2
        return loss

def ls_generator_loss(scores_fake):
    """
    Computes the Least-Squares GAN loss for the generator.

    Inputs:
    - scores_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

    Outputs:
    - loss: A PyTorch Tensor containing the loss.
    """
    loss = torch.mean((scores_fake - 1)**2)/2
    return loss
```

Before running a GAN with our new loss function, let's check it:

```
In [74]: def test_lsgan_loss(score_real, score_fake, d_loss_true, g_loss_true):
        score_real = torch.Tensor(score_real).type(dtype)
        score_fake = torch.Tensor(score_fake).type(dtype)
        d_loss = ls_discriminator_loss(score_real, score_fake).cpu().numpy()
        g_loss = ls_generator_loss(score_fake).cpu().numpy()
        print("Maximum error in d_loss: %g"%rel_error(d_loss_true, d_loss))
        print("Maximum error in g_loss: %g"%rel_error(g_loss_true, g_loss))

test_lsgan_loss(answers['logits_real'], answers['logits_fake'],
                answers['d_loss_lsgan_true'], answers['g_loss_lsgan_true'])

Maximum error in d_loss: 1.53171e-08
Maximum error in g_loss: 2.7837e-09
```

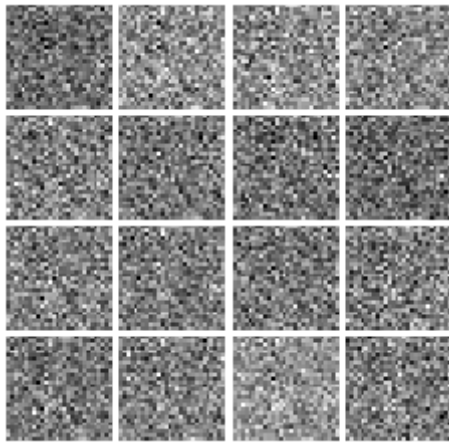
Run the following cell to train your model!


```
In [75]: D_LS = discriminator().type(dtype)
          G_LS = generator().type(dtype)

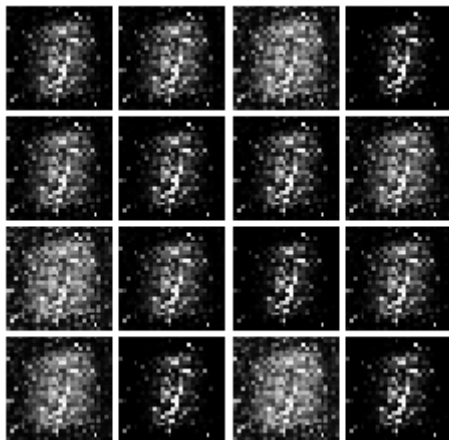
          D_LS_solver = get_optimizer(D_LS)
          G_LS_solver = get_optimizer(G_LS)

          run_a_gan(D_LS, G_LS, D_LS_solver, G_LS_solver, ls_discriminator_loss, l
                    s_generator_loss)
```

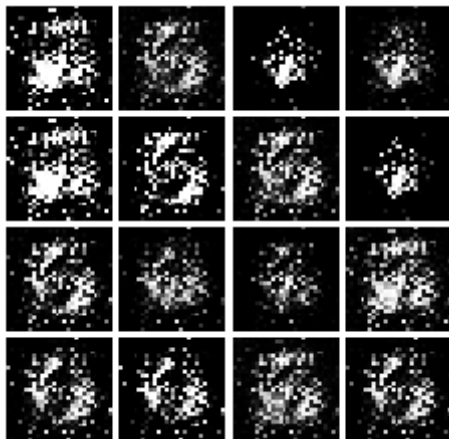
Iter: 0, D: 0.6202, G:0.5



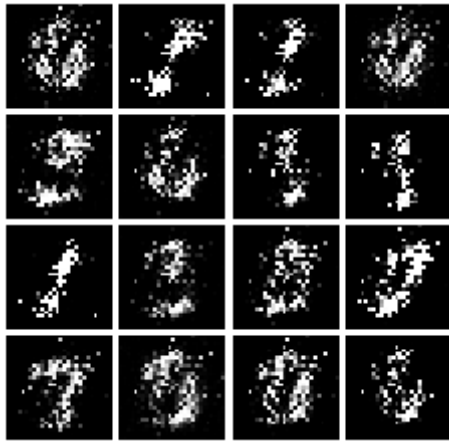
Iter: 250, D: 0.1284, G:0.7519



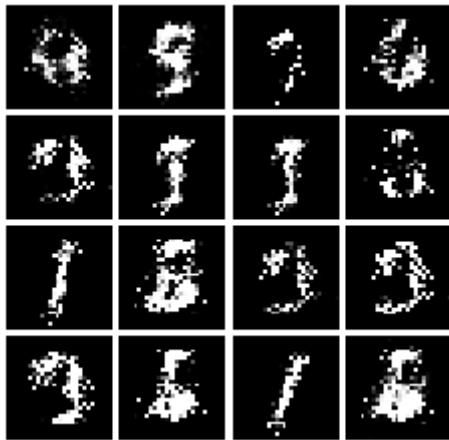
Iter: 500, D: 0.5919, G:0.06593



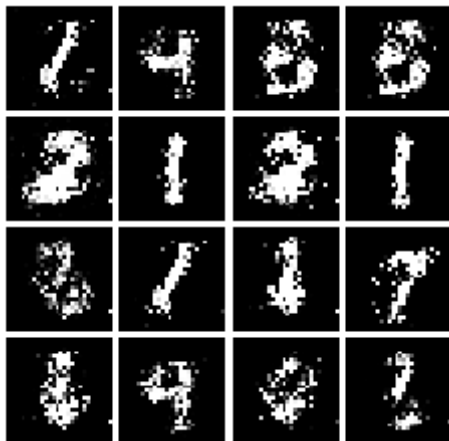
Iter: 750, D: 0.1399, G:0.2599



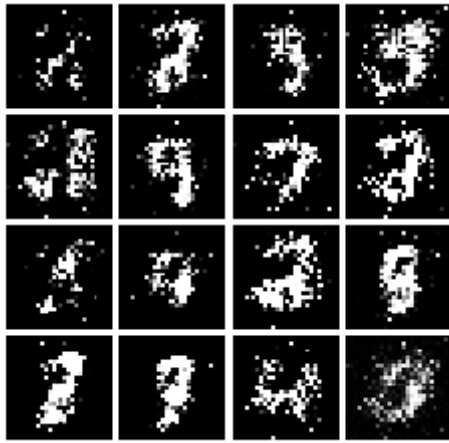
Iter: 1000, D: 0.17, G:0.2505



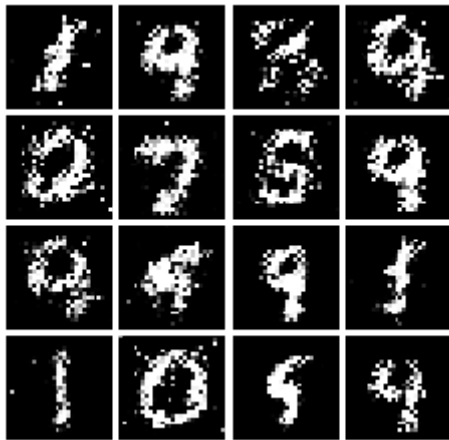
Iter: 1250, D: 0.1303, G:0.3523



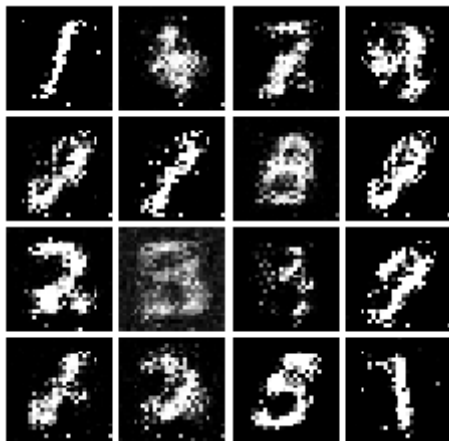
Iter: 1500, D: 0.1983, G:0.2713



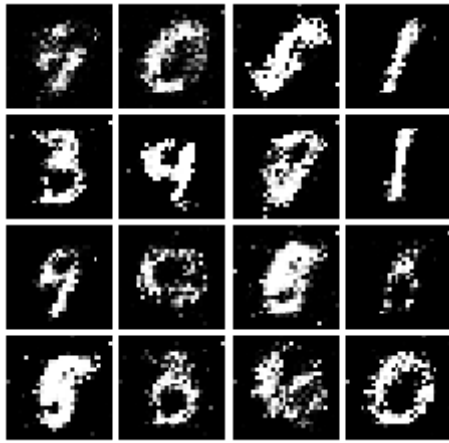
Iter: 1750, D: 0.2297, G:0.1875



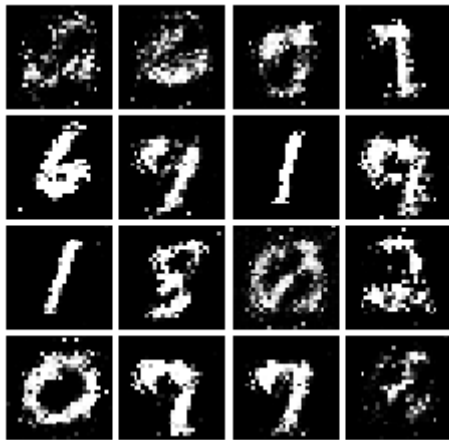
Iter: 2000, D: 0.2038, G:0.2138



Iter: 2250, D: 0.2298, G:0.1659



Iter: 2500, D: 0.2322, G:0.1518



Iter: 2750, D: 0.2196, G:0.1523



Iter: 3000, D: 0.2202, G:0.1195



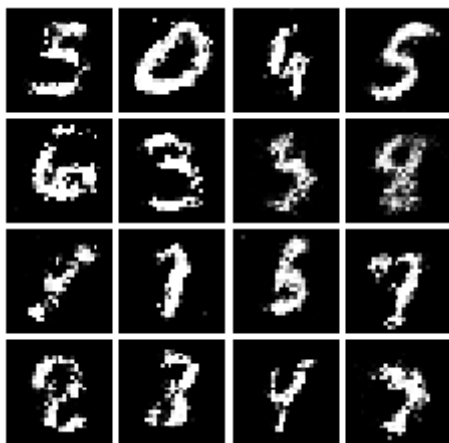
Iter: 3250, D: 0.2282, G:0.1776



Iter: 3500, D: 0.2032, G:0.1171



Iter: 3750, D: 0.2389, G:0.1679



Deeply Convolutional GANs

In the first part of the notebook, we implemented an almost direct copy of the original GAN network from Ian Goodfellow. However, this network architecture allows no real spatial reasoning. It is unable to reason about things like "sharp edges" in general because it lacks any convolutional layers. Thus, in this section, we will implement some of the ideas from DCGAN (<https://arxiv.org/abs/1511.06434>), where we use convolutional networks

Discriminator

We will use a discriminator inspired by the TensorFlow MNIST classification tutorial, which is able to get above 99% accuracy on the MNIST dataset fairly quickly.

- Reshape into image tensor (Use Unflatten!)
- Conv2D: 32 Filters, 5x5, Stride 1
- Leaky ReLU(alpha=0.01)
- Max Pool 2x2, Stride 2
- Conv2D: 64 Filters, 5x5, Stride 1
- Leaky ReLU(alpha=0.01)
- Max Pool 2x2, Stride 2
- Flatten
- Fully Connected with output size 4 x 4 x 64
- Leaky ReLU(alpha=0.01)
- Fully Connected with output size 1

```
In [76]: def build_dc_classifier():
        """
        Build and return a PyTorch model for the DCGAN discriminator implementing
        the architecture above.
        """
        return nn.Sequential(
            #####
            ##### TO DO #####
            #####
            Unflatten(batch_size, 1, 28, 28),
            nn.Conv2d(1, 32, 5),
            nn.LeakyReLU(0.01),
            nn.MaxPool2d(2),
            nn.Conv2d(32, 64, 5),
            nn.LeakyReLU(0.01),
            nn.MaxPool2d(2),
            Flatten(),
            nn.Linear(4*4*64, 4*4*64),
            nn.LeakyReLU(0.01),
            nn.Linear(4*4*64, 1)
        )

data = next(enumerate(loader_train))[-1][0].type(dtype)
b = build_dc_classifier().type(dtype)
out = b(data)
print(out.size())

torch.Size([128, 1])
```

Check the number of parameters in your classifier as a sanity check:

```
In [79]: def test_dc_classifer(true_count=1102721):
        model = build_dc_classifier()
        cur_count = count_params(model)
        if cur_count != true_count:
            print('Incorrect number of parameters in generator. Check your architecture.')
        else:
            print('Correct number of parameters in generator.')

test_dc_classifer()

Correct number of parameters in generator.
```


Generator

For the generator, we will copy the architecture exactly from the [InfoGAN paper](https://arxiv.org/pdf/1606.03657.pdf) (<https://arxiv.org/pdf/1606.03657.pdf>). See Appendix C.1 MNIST. See the documentation for `tf.nn.conv2d transpose` (https://www.tensorflow.org/api_docs/python/tf/nn/conv2d_transpose). We are always "training" in GAN mode.

- Fully connected with output size 1024
- ReLU
- BatchNorm
- Fully connected with output size 7 x 7 x 128
- ReLU
- BatchNorm
- Reshape into Image Tensor of shape 7, 7, 128
- Conv2D^T (Transpose): 64 filters of 4x4, stride 2, 'same' padding
- ReLU
- BatchNorm
- Conv2D^T (Transpose): 1 filter of 4x4, stride 2, 'same' padding
- TanH
- Should have a 28x28x1 image, reshape back into 784 vector

```
In [82]: def build_dc_generator(noise_dim=NOISE_DIM):
    """
    Build and return a PyTorch model implementing the DCGAN generator using
    the architecture described above.
    """
    return nn.Sequential(
        #####
        ##### TO DO #####
        #####
        nn.Linear(noise_dim, 1024),
        nn.ReLU(),
        nn.BatchNorm1d(1024),
        nn.Linear(1024, 128*7*7),
        nn.ReLU(),
        nn.BatchNorm1d(128*7*7),
        Unflatten(-1, 128, 7, 7),
        nn.ConvTranspose2d(128, 64, 4, stride=2, padding=1),
        nn.ReLU(),
        nn.BatchNorm2d(64),
        nn.ConvTranspose2d(64, 1, 4, stride=2, padding=1),
        nn.Tanh(),
        Flatten()
    )

test_g_gan = build_dc_generator().type(dtype)
test_g_gan.apply(initialize_weights)

fake_seed = torch.randn(batch_size, NOISE_DIM).type(dtype)
fake_images = test_g_gan.forward(fake_seed)
fake_images.size()
```

```
Out[82]: torch.Size([128, 784])
```

Check the number of parameters in your generator as a sanity check:

```
In [83]: def test_dc_generator(true_count=6580801):
    model = build_dc_generator(4)
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in generator. Check your architecture.')
    else:
        print('Correct number of parameters in generator.')

test_dc_generator()

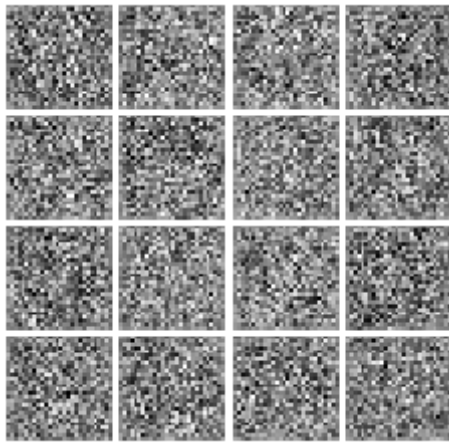
Correct number of parameters in generator.
```

```
In [84]: D_DC = build_dc_classifier().type(dtype)
D_DC.apply(initialize_weights)
G_DC = build_dc_generator().type(dtype)
G_DC.apply(initialize_weights)

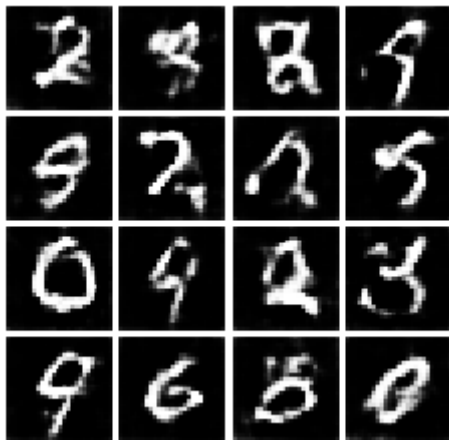
D_DC_solver = get_optimizer(D_DC)
G_DC_solver = get_optimizer(G_DC)

run_a_gan(D_DC, G_DC, D_DC_solver, G_DC_solver, discriminator_loss, generator_loss, num_epochs=5)
```

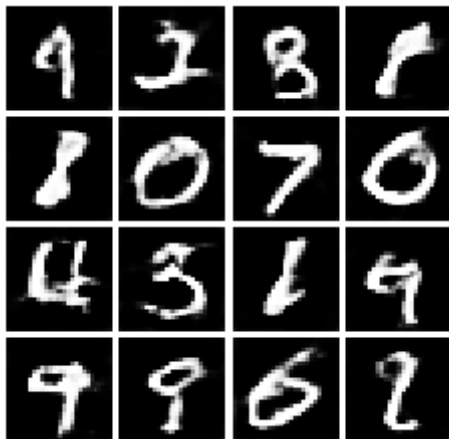
Iter: 0, D: 1.402, G:0.4124



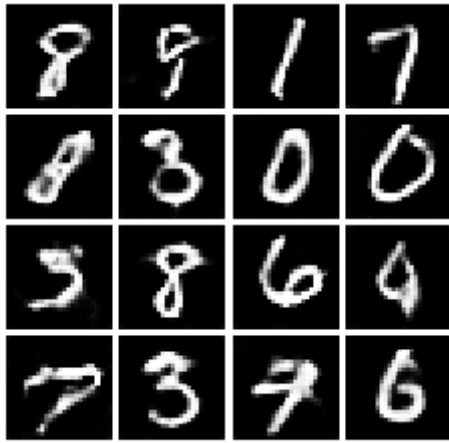
Iter: 250, D: 1.298, G:0.6493



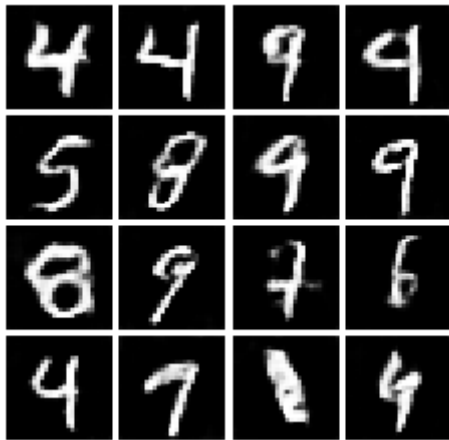
Iter: 500, D: 1.088, G:1.069



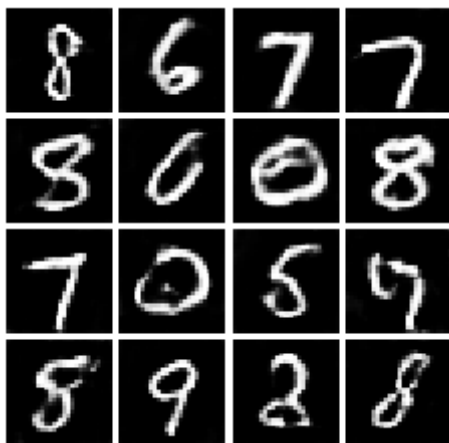
Iter: 750, D: 1.208, G:0.951



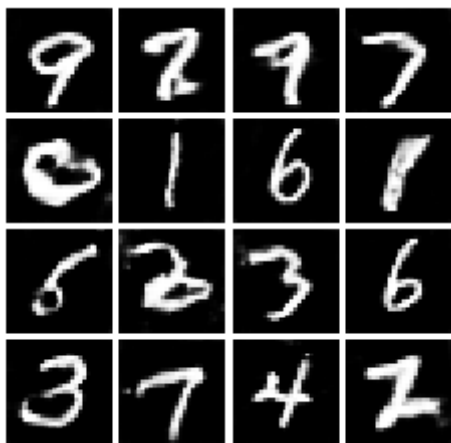
Iter: 1000, D: 1.232, G:0.9834



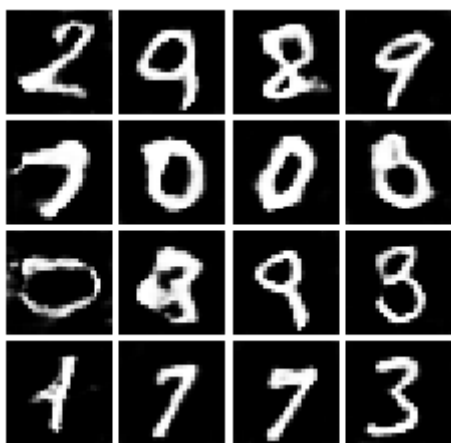
Iter: 1250, D: 1.287, G:0.9408



Iter: 1500, D: 1.192, G:1.001



Iter: 1750, D: 1.388, G:0.7739



INLINE QUESTION 1

We will look at an example to see why alternating minimization of the same objective (like in a GAN) can be tricky business.

Consider $f(x, y) = xy$. What does $\min_x \max_y f(x, y)$ evaluate to? (Hint: minmax tries to minimize the maximum value achievable.)

Now try to evaluate this function numerically for 6 steps, starting at the point $(1, 1)$, by using alternating gradient (first updating y , then updating x) with step size 1. You'll find that writing out the update step in terms of $x_t, y_t, x_{t+1}, y_{t+1}$ will be useful.

Record the six pairs of explicit values for (x_t, y_t) in the table below.

Your answer:

y_0	y_1	y_2	y_3	y_4	y_5	y_6
1	2	1	-1	-2	-1	1
x_0	x_1	x_2	x_3	x_4	x_5	x_6
1	-1	-2	-1	1	2	1

INLINE QUESTION 2

Using this method, will we ever reach the optimal value? Why or why not?

Your answer:

(y_0, x_0) start at (1,1) and loop to (1,1) again during gradient update which indicate all future (x, y) would stuck at this path and never reach the optimal value.

INLINE QUESTION 3

If the generator loss decreases during training while the discriminator loss stays at a constant high value from the start, is this a good sign? Why or why not? A qualitative answer is sufficient

Your answer:

It is not a good sign, discriminator loss stays at a constant high value means the discriminator can't recognize fake image as fake and real image as real, which even the generator loss decreases, a bad discriminator can not indicate generator is good enough