

# 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.](#) 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.](#), 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.](#)

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](#), and [hundreds of new papers](#). 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](#) 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](#). 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](#), [WGAN-GP](#).

GANs are not the only way to train a generative model! For other approaches to generative modeling check out the [deep generative model chapter](#) of the Deep Learning [book](#). Another popular way of training neural networks as generative models is Variational Autoencoders (co-discovered [here](#) and [here](#)). 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

```
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'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
```

```
%autoreload 2

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

    fig = plt.figure(figsize=(sqrt_n, sqrt_n))
    gs = gridspec.GridSpec(sqrt_n, sqrt_n)
    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([sqrt_m, sqrt_m]))
    return
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```
# Colab users only
%cd drive/My\ Drive/$FOLDERNAME/
%cp -r gan-checks-tf.npz /content/
%cd /content/
```

```
[Errno 2] No such file or directory: 'drive/My Drive/$FOLDERNAME/'
/home/daizj/Homeworks/dl/hw4
cp: cannot create regular file '/content/': Not a directory
[Errno 2] No such file or directory: '/content/'
/home/daizj/Homeworks/dl/hw4
```

```
from cs231n.gan_pytorch import preprocess_img, deprocess_img, rel_error,
count_params, ChunkSampler

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](#) 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`.

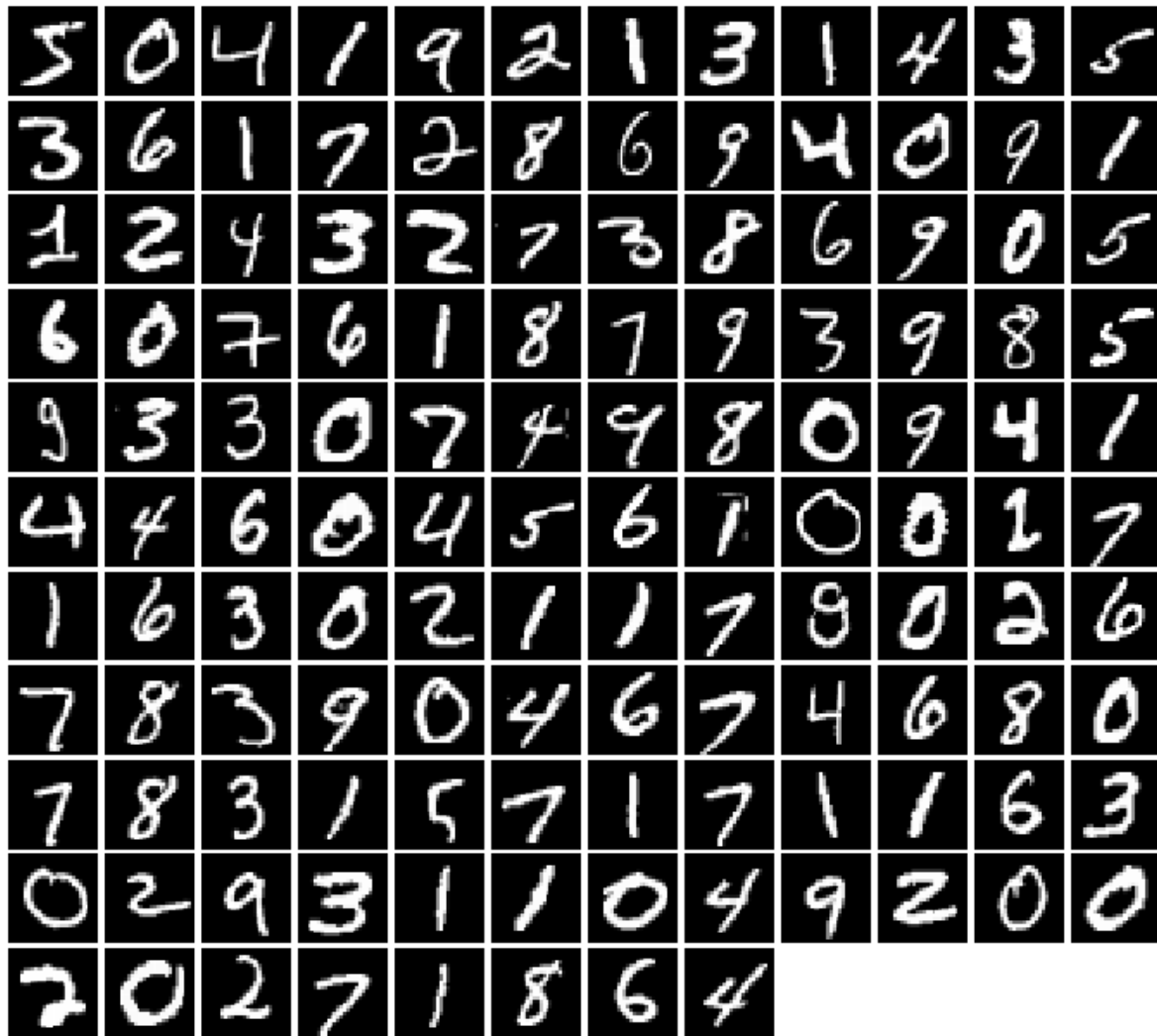
```
NUM_TRAIN = 50000
NUM_VAL = 5000

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)
```



## Random Noise

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

Implement `sample_noise` in `cs231n/gan_pytorch.py`.

Hint: use `torch.rand`.

Make sure noise is the correct shape and type:

```
from cs231n.gan_pytorch import sample_noise

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.

```
from cs231n.gan_pytorch import Flatten, Unflatten, initialize_weights
```

## 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.

**If you are a Colab user, it is recommended to change colab runtime to GPU.**

```
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  $\alpha$ ; 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.

Implement `discriminator` in `cs231n/gan_pytorch.py`

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

```

from cs231n.gan_pytorch import discriminator

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])

Implement `generator` in `cs231n/gan_pytorch.py`

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

```

from cs231n.gan_pytorch import generator

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.

Implement `bce_loss`, `discriminator_loss`, `generator_loss` in `cs231n/gan_pytorch.py`

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

```
from cs231n.gan_pytorch import bce_loss, discriminator_loss, generator_loss

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: 2.83811e-08

```
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,  $\beta_1=0.5$ ,  $\beta_2=0.999$ . You'll use this to construct optimizers for the generators and discriminators for the rest of the notebook.

Implement `get_optimizer` in `cs231n/gan_pytorch.py`

## Training a GAN!

We provide you the main training loop... you won't need to change `run_a_gan` in `cs231n/gan_pytorch.py`, but we encourage you to read through and understand it.

```
from cs231n.gan_pytorch import get_optimizer, run_a_gan
```

```
# Make the discriminator
D = discriminator().type(dtype)

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

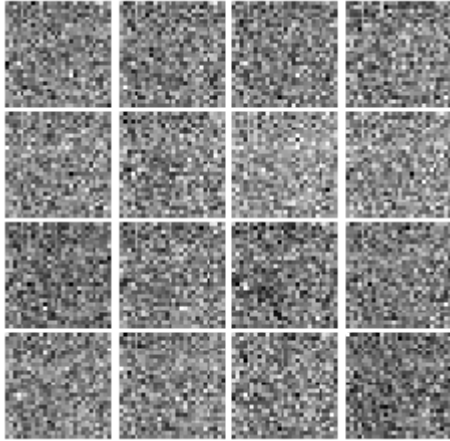
# Use the function you wrote earlier to get optimizers for the Discriminator and
the Generator
D_solver = get_optimizer(D)
G_solver = get_optimizer(G)
# Run it!
images = run_a_gan(D, G, D_solver, G_solver, discriminator_loss, generator_loss,
loader_train)
```

```
Iter: 0, D: 1.328, G:0.7202
Iter: 250, D: 1.31, G:1.058
Iter: 500, D: 1.273, G:1.105
Iter: 750, D: 1.173, G:1.106
Iter: 1000, D: 0.9179, G:1.434
Iter: 1250, D: 1.205, G:1.002
Iter: 1500, D: 1.304, G:0.9009
Iter: 1750, D: 1.327, G:0.7838
Iter: 2000, D: 1.287, G:0.9228
Iter: 2250, D: 1.383, G:0.8177
Iter: 2500, D: 1.374, G:0.7767
Iter: 2750, D: 1.329, G:0.7857
Iter: 3000, D: 1.344, G:0.8106
Iter: 3250, D: 1.289, G:0.8252
Iter: 3500, D: 1.305, G:0.8144
Iter: 3750, D: 1.33, G:0.8106
```

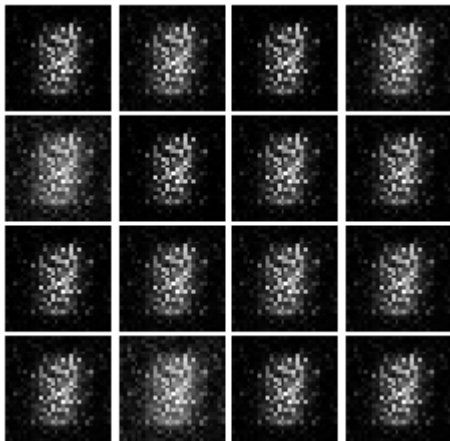
Run the cell below to show the generated images.

```
numIter = 0
for img in images:
    print("Iter: {}".format(numIter))
    show_images(img)
    plt.show()
    numIter += 250
    print()
```

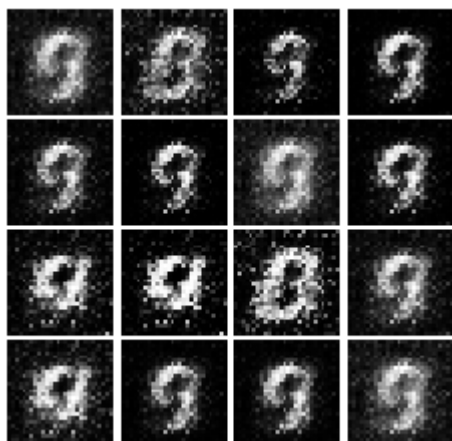
Iter: 0



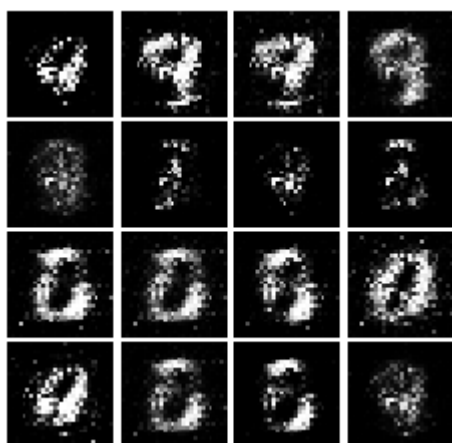
Iter: 250



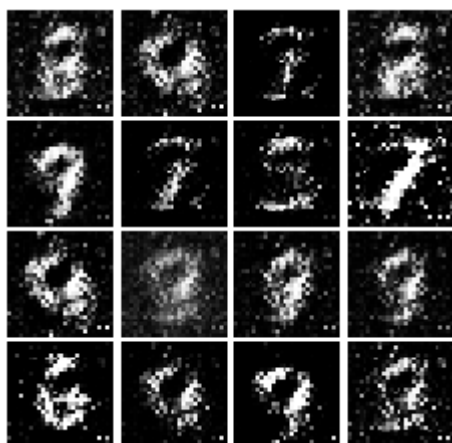
Iter: 500



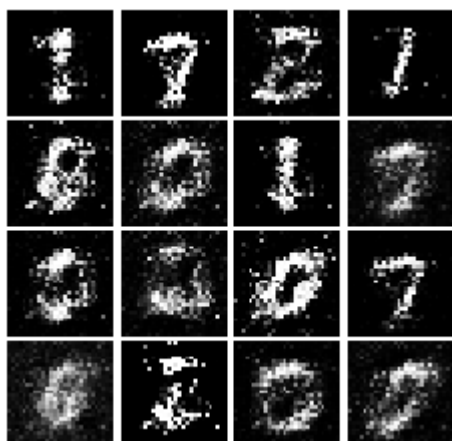
Iter: 750



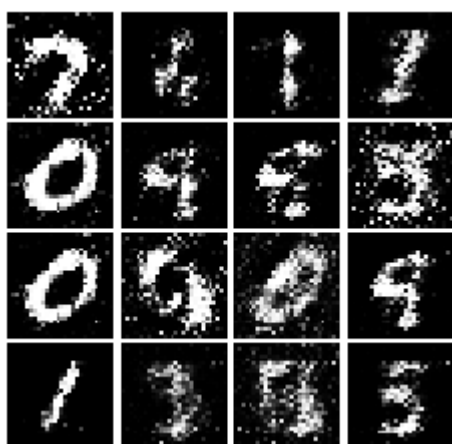
Iter: 1000



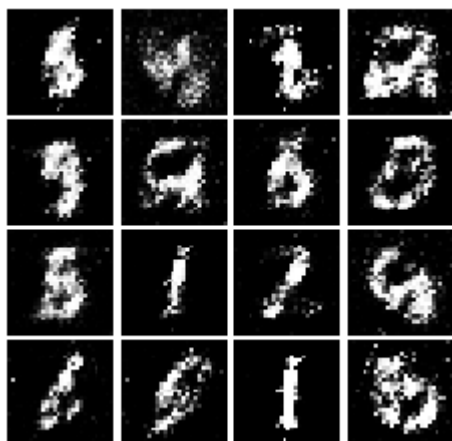
Iter: 1250



Iter: 1500



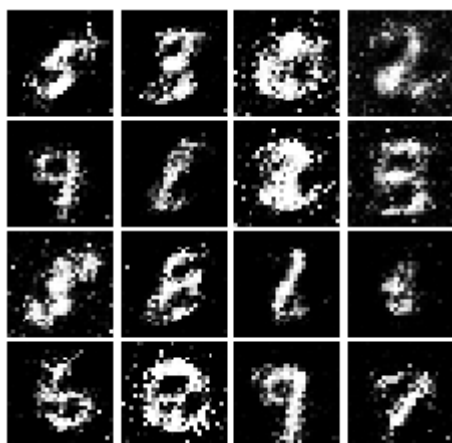
Iter: 1750



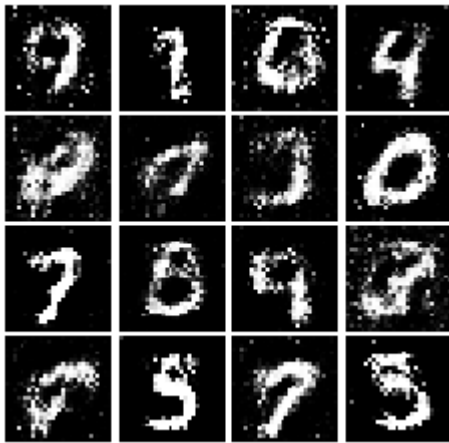
Iter: 2000



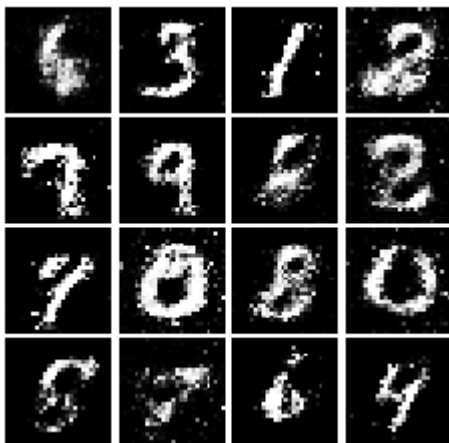
Iter: 2250



Iter: 2500



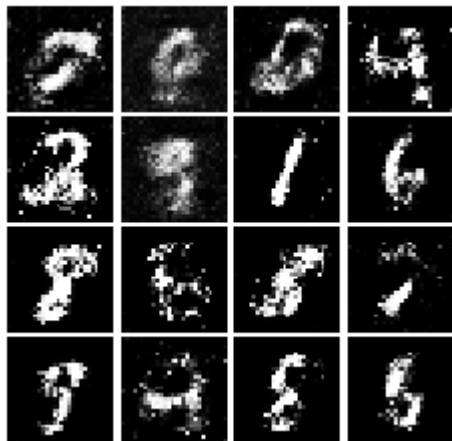
Iter: 2750



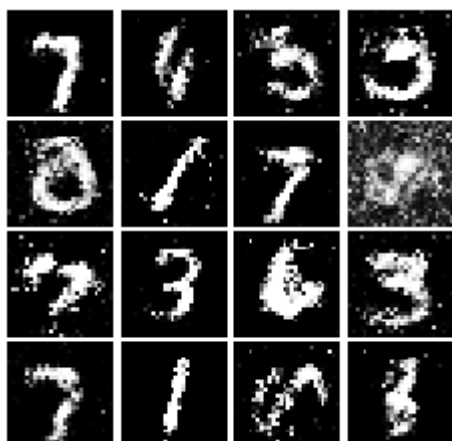
Iter: 3000



Iter: 3250



Iter: 3500



Iter: 3750



Please tag the cell below on Gradescope while submitting.

```
print("Vanilla GAN Final image:")  
show_images(images[-1])  
plt.show()
```

Vanilla GAN Final image:



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](#), 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:

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

and the discriminator loss:

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

**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` ).

Implement `ls_discriminator_loss`, `ls_generator_loss` in `cs231n/gan_pytorch.py`

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

```
from cs231n.gan_pytorch import ls_discriminator_loss, ls_generator_loss

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: 3.36961e-08
```

Run the following cell to train your model!

```
D_LS = discriminator().type(dtype)
G_LS = generator().type(dtype)

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

images = run_a_gan(D_LS, G_LS, D_LS_solver, G_LS_solver, ls_discriminator_loss,
                  ls_generator_loss, loader_train)
```

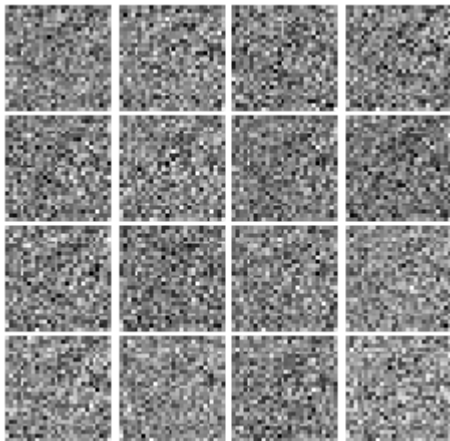
```
Iter: 0, D: 0.5689, G:0.51
Iter: 250, D: 0.09411, G:0.3328
Iter: 500, D: 0.1757, G:0.6882
Iter: 750, D: 0.1784, G:0.3486
Iter: 1000, D: 0.1765, G:0.2795
Iter: 1250, D: 0.219, G:0.3982
Iter: 1500, D: 0.1863, G:0.2372
Iter: 1750, D: 0.2372, G:0.2035
Iter: 2000, D: 0.2166, G:0.1922
```

```
Iter: 2250, D: 0.2273, G:0.154  
Iter: 2500, D: 0.2108, G:0.1511  
Iter: 2750, D: 0.2353, G:0.1767  
Iter: 3000, D: 0.2265, G:0.1452  
Iter: 3250, D: 0.2024, G:0.182  
Iter: 3500, D: 0.2348, G:0.161  
Iter: 3750, D: 0.2276, G:0.1784
```

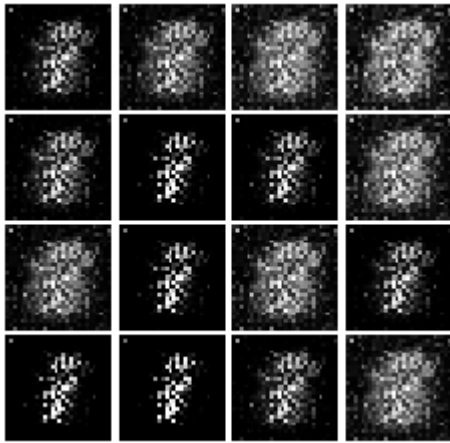
Run the cell below to show generated images.

```
numIter = 0  
for img in images:  
    print("Iter: {}".format(numIter))  
    show_images(img)  
    plt.show()  
    numIter += 250  
    print()
```

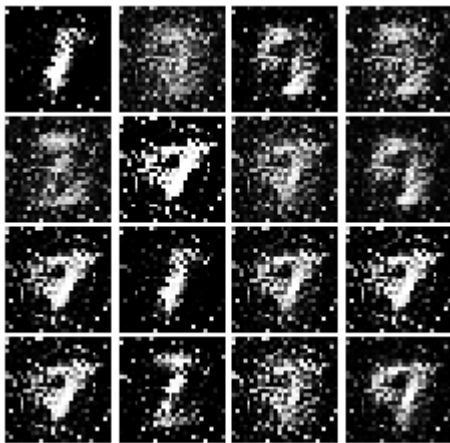
Iter: 0



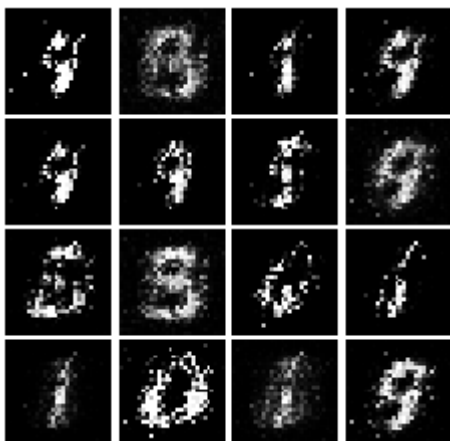
Iter: 250



Iter: 500



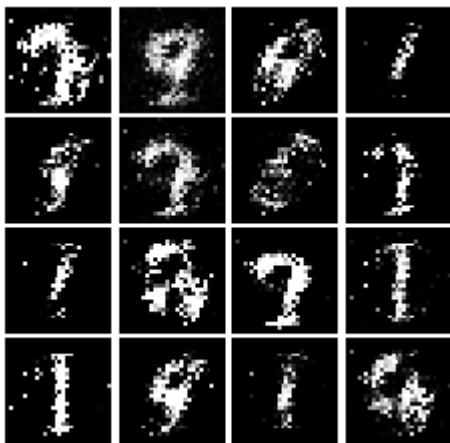
Iter: 750



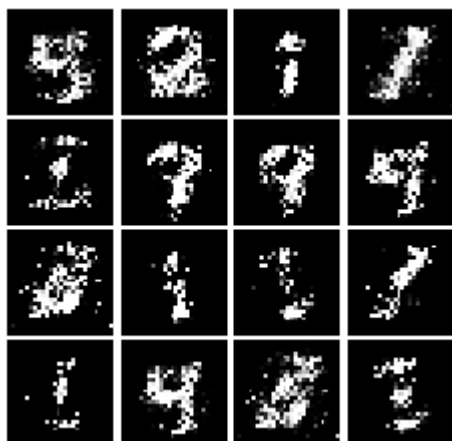
Iter: 1000



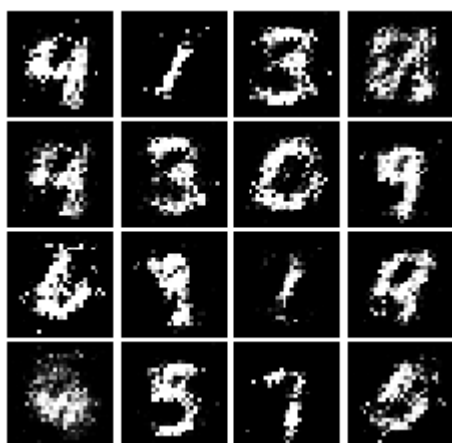
Iter: 1250



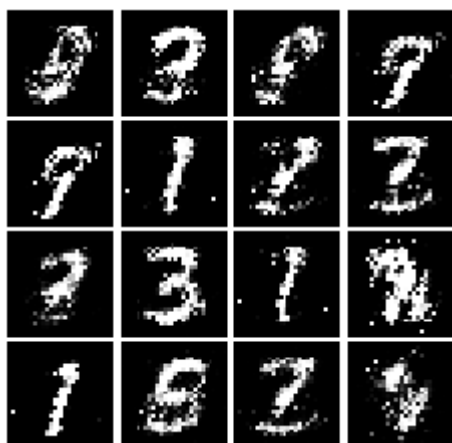
Iter: 1500



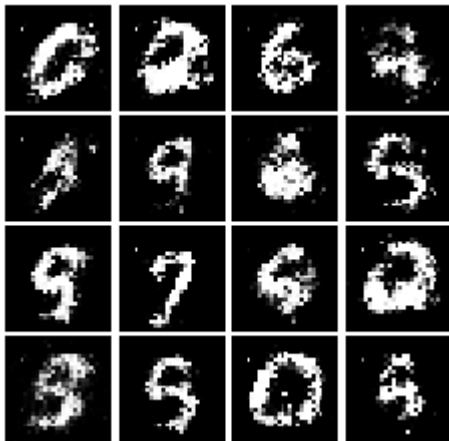
Iter: 1750



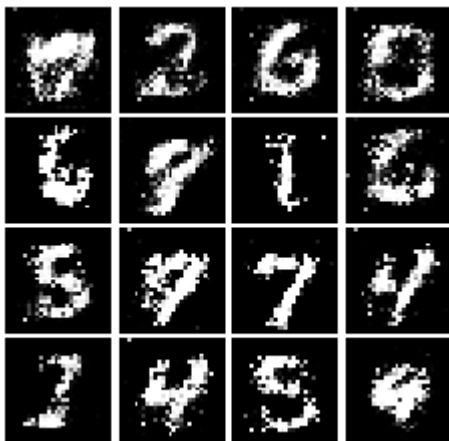
Iter: 2000



Iter: 2250



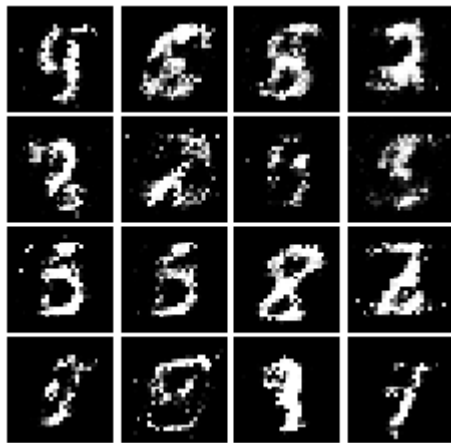
Iter: 2500



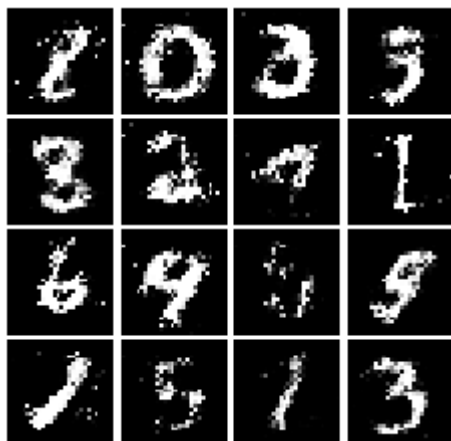
Iter: 2750



Iter: 3000



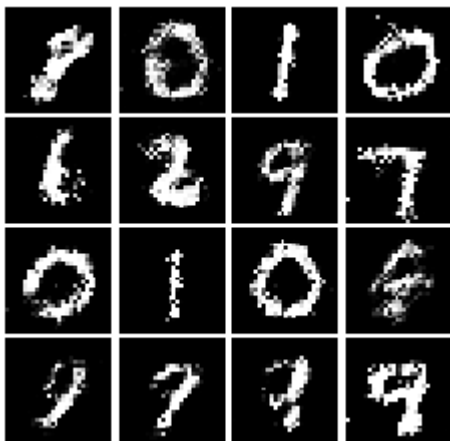
Iter: 3250



Iter: 3500



Iter: 3750

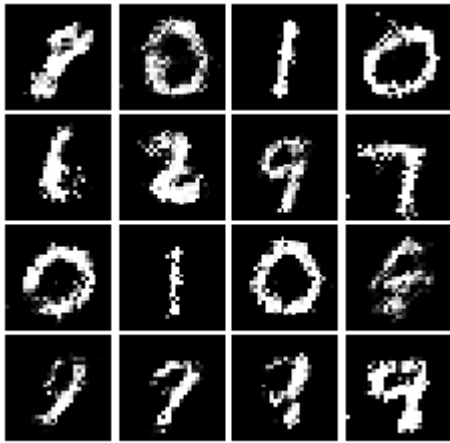


Please tag the cell below on Gradescope while submitting.

```
print("LSGAN Final image:")  
show_images(images[-1])  
plt.show()
```

LSGAN Final image:





## 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](#), 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

Implement `build_dc_classifier` in `cs231n/gan_pytorch.py`

```
from cs231n.gan_pytorch import build_dc_classifier

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

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

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

```
def test_dc_classifier(true_count=1102721):
    model = build_dc_classifier(batch_size)
    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_classifier()
```

Correct number of parameters in generator.

## Generator

For the generator, we will copy the architecture exactly from the [InfoGAN paper](#). See Appendix C.1 MNIST. See the documentation for [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 (use padding=1)
- ReLU
- BatchNorm
- Conv2D^T (Transpose): 1 filter of 4x4, stride 2, 'same' padding (use padding=1)
- TanH
- Should have a 28x28x1 image, reshape back into 784 vector

Implement `build_dc_generator` in `cs231n/gan_pytorch.py`

```
from cs231n.gan_pytorch import build_dc_generator

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()
```

torch.Size([128, 784])

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

```
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.

```
D_DC = build_dc_classifier(batch_size).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)

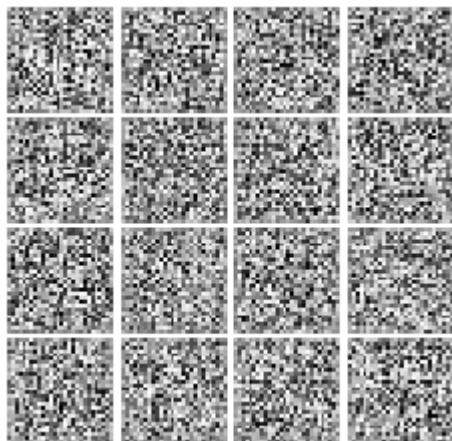
images = run_a_gan(D_DC, G_DC, D_DC_solver, G_DC_solver, discriminator_loss,
generator_loss, loader_train, num_epochs=5)
```

```
Iter: 0, D: 1.4, G:2.369
Iter: 250, D: 1.262, G:0.7922
Iter: 500, D: 1.178, G:0.9531
Iter: 750, D: 1.19, G:1.415
Iter: 1000, D: 1.276, G:0.874
Iter: 1250, D: 1.174, G:1.181
Iter: 1500, D: 1.126, G:1.056
Iter: 1750, D: 1.102, G:1.066
```

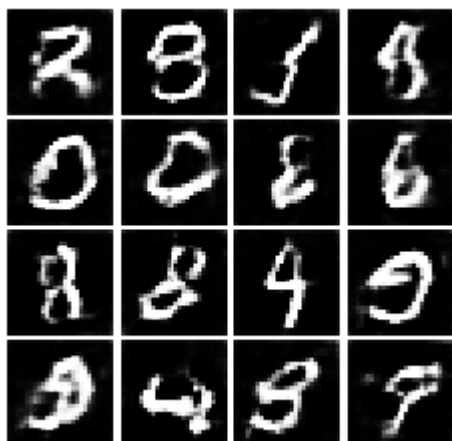
Run the cell below to show generated images.

```
numIter = 0
for img in images:
    print("Iter: {}".format(numIter))
    show_images(img)
    plt.show()
    numIter += 250
    print()
```

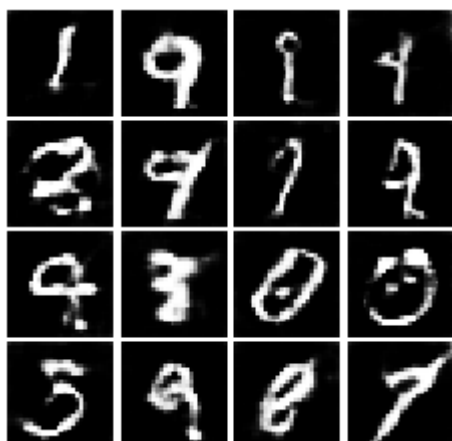
Iter: 0



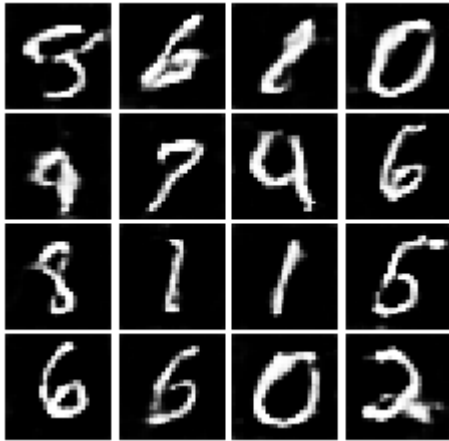
Iter: 250



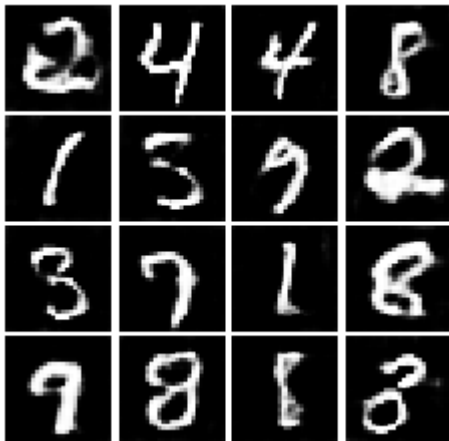
Iter: 500



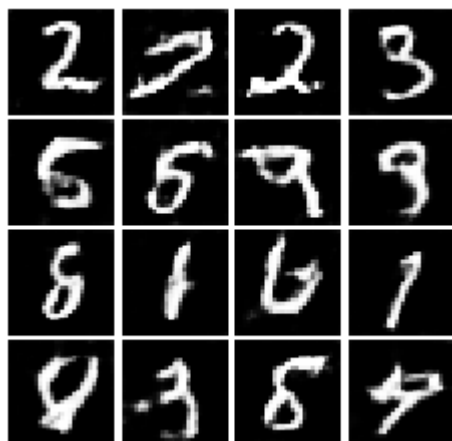
Iter: 750



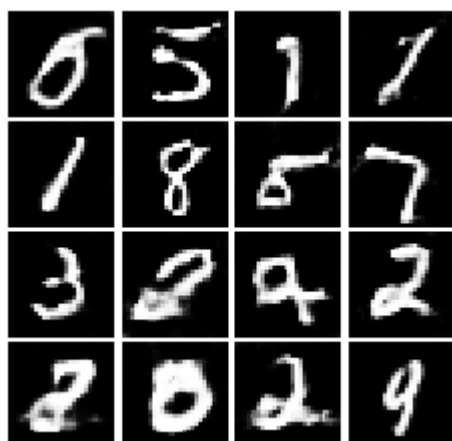
Iter: 1000



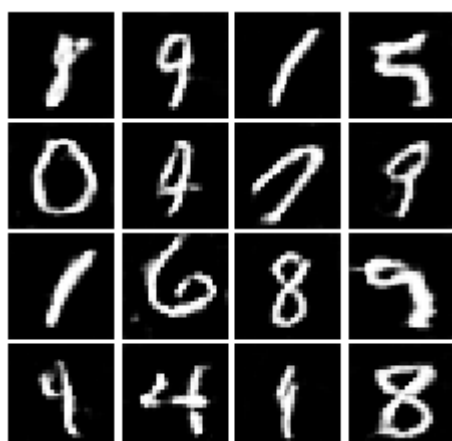
Iter: 1250



Iter: 1500



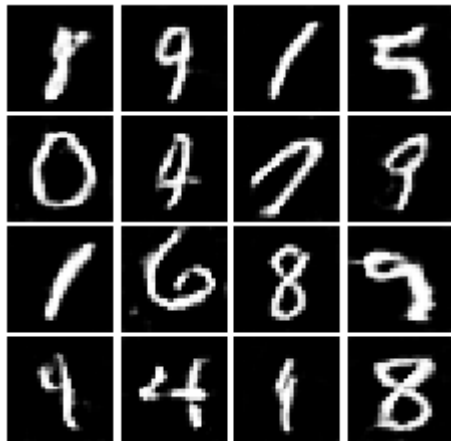
Iter: 1750



Please tag the cell below on Gradescope while submitting.

```
print("DCGAN Final image:")
show_images(images[-1])
plt.show()
```

DCGAN Final image:



## 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  $x$  using that updated  $y$ ) with step size **1**. **Here step size is the learning\_rate, and steps will be learning\_rate \* gradient.**

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.

Briefly explain what  $\min_x \max_y f(x, y)$  evaluates to and 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:**

No, because it will reach a local optimal value instead of a global one.

## 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:**

No, because the generator is easy to cheat the discriminator, which means the discriminator is weak.



# Style Transfer

In this notebook we will implement the style transfer technique from ["Image Style Transfer Using Convolutional Neural Networks" \(Gatys et al., CVPR 2015\)](#).

The general idea is to take two images, and produce a new image that reflects the content of one but the artistic "style" of the other. We will do this by first formulating a loss function that matches the content and style of each respective image in the feature space of a deep network, and then performing gradient descent on the pixels of the image itself.

The deep network we use as a feature extractor is [SqueezeNet](#), a small model that has been trained on ImageNet. You could use any network, but we chose SqueezeNet here for its small size and efficiency.

Here's an example of the images you'll be able to produce by the end of this notebook:

## Part 0: Setup

```
import torch
import torch.nn as nn
import torchvision
import torchvision.transforms as T
import PIL

import numpy as np

import matplotlib.pyplot as plt

from cs231n.image_utils import SQUEEZENET_MEAN, SQUEEZENET_STD
%matplotlib inline
# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-
# ipython
%load_ext autoreload
%autoreload 2
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

We provide you with some helper functions to deal with images, since for this part of the assignment we're dealing with real JPEGs, not CIFAR-10 data.

```
from cs231n.style_transfer_pytorch import preprocess, deprocess, rescale,
rel_error, features_from_img

CHECKS_PATH = None

# Local
CHECKS_PATH = 'style-transfer-checks.npz'
```

```
# Colab
#CHECKS_PATH = '/content/drive/My Drive/{}/{}'.format(FOLDERNAME, 'style-
transfer-checks.npz')

assert CHECKS_PATH is not None, "[!] Choose path to style-transfer-checks.npz"

STYLES_FOLDER = CHECKS_PATH.replace('style-transfer-checks.npz', 'styles')

answers = dict(np.load(CHECKS_PATH))
```

Pytorch has two separate types for Tensors that contain floating-point numbers: one for operations on the CPU (`torch.FloatTensor`), and one using CUDA for operations on the GPU (`torch.cuda.FloatTensor`).

We'll be using this type variable more later, so we need to set the tensor type to one of them.

```
dtype = torch.FloatTensor
# Uncomment out the following line if you're on a machine with a GPU set up for
PyTorch!
#dtype = torch.cuda.FloatTensor
```

```
# Load the pre-trained SqueezeNet model.
cnn = torchvision.models.squeezenet1_1(pretrained=True).features
cnn.type(dtype)

# We don't want to train the model any further, so we don't want PyTorch to waste
computation
# computing gradients on parameters we're never going to update.
for param in cnn.parameters():
    param.requires_grad = False
```

```
/home/daizj/anaconda3/envs/cs231n/lib/python3.7/site-
packages/torchvision/models/_utils.py:209: UserWarning: The parameter
'pretrained' is deprecated since 0.13 and will be removed in 0.15, please use
'weights' instead.
  f"The parameter '{pretrained_param}' is deprecated since 0.13 and will be
removed in 0.15, "
/home/daizj/anaconda3/envs/cs231n/lib/python3.7/site-
packages/torchvision/models/_utils.py:223: UserWarning: Arguments other than a
weight enum or `None` for 'weights' are deprecated since 0.13 and will be removed
in 0.15. The current behavior is equivalent to passing
`weights=SqueezeNet1_1_weights.IMAGENET1K_V1`. You can also use
`weights=SqueezeNet1_1_weights.DEFAULT` to get the most up-to-date weights.
  warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/squeezenet1_1-b8a52dc0.pth" to
/home/daizj/.cache/torch/hub/checkpoints/squeezenet1_1-b8a52dc0.pth
100.0%
```

## Part 1: Computing Loss

We're going to compute the three components of our loss function now. The loss function is a weighted sum of three terms: content loss + style loss + total variation loss. You'll fill in the functions that compute these weighted terms below.

## Part 1A: Content loss

We can generate an image that reflects the content of one image and the style of another by incorporating both in our loss function. We want to penalize deviations from the content of the content image and deviations from the style of the style image. We can then use this hybrid loss function to perform gradient descent **not on the parameters** of the model, but instead **on the pixel values** of our original image.

Let's first write the content loss function. Content loss measures how much the feature map of the generated image differs from the feature map of the source image. We only care about the content representation of one layer of the network (say, layer  $\ell$ ), that has feature maps  $A^\ell \in \mathbb{R}^{1 \times C_\ell \times H_\ell \times W_\ell}$ .  $C_\ell$  is the number of filters/channels in layer  $\ell$ ,  $H_\ell$  and  $W_\ell$  are the height and width. We will work with reshaped versions of these feature maps that combine all spatial positions into one dimension. Let  $F^\ell \in \mathbb{R}^{C_\ell \times M_\ell}$  be the feature map for the current image and  $P^\ell \in \mathbb{R}^{C_\ell \times M_\ell}$  be the feature map for the content source image where  $M_\ell = H_\ell \times W_\ell$  is the number of elements in each feature map. Each row of  $F^\ell$  or  $P^\ell$  represents the vectorized activations of a particular filter, convolved over all positions of the image. Finally, let  $w_c$  be the weight of the content loss term in the loss function.

Then the content loss is given by:

$$L_c = w_c \times \sum_{i,j} (F_{ij}^\ell - P_{ij}^\ell)^2$$

Implement `content_loss` in `cs231n/style_transfer_pytorch.py`

Test your content loss. You should see errors less than 0.001.

```
from cs231n.style_transfer_pytorch import content_loss, extract_features,
features_from_img

def content_loss_test(correct):
    content_image = '%s/tubingen.jpg' % (STYLES_FOLDER)
    image_size = 192
    content_layer = 3
    content_weight = 6e-2

    c_feats, content_img_var = features_from_img(content_image, image_size, cnn)

    bad_img = torch.zeros(*content_img_var.data.size()).type(dtype)
    feats = extract_features(bad_img, cnn)

    student_output = content_loss(content_weight, c_feats[content_layer],
    feats[content_layer]).cpu().data.numpy()
    error = rel_error(correct, student_output)
    print('Maximum error is {:.3f}'.format(error))

content_loss_test(answers['c1_out'])
```

Maximum error is 0.000

## Part 1B: Style loss

Now we can tackle the style loss. For a given layer  $\ell$ , the style loss is defined as follows:

First, compute the Gram matrix  $G$  which represents the correlations between the values in each channel of the feature map (i.e. the "responses" of the filter responsible for that channel), where  $F$  is as above. The Gram matrix is an approximation of the covariance matrix -- it tells us how every channel's values (i.e. that filter's activations) correlate with every other channel's values. If we have  $C$  channels, matrix  $G$  will be of shape  $(C, C)$  to capture these correlations.

We want the activation statistics of our generated image to match the activation statistics of our style image, and matching the (approximate) covariance is one way to do that. There are a variety of ways you could do this, but the Gram matrix is nice because it's easy to compute and in practice shows good results.

Given a feature map  $F^\ell$  of shape  $(C_\ell, H_\ell, W_\ell)$ , we can flatten the height and width dimensions so they're just 1 dimension  $M_\ell = H_\ell \times W_\ell$ : the new shape of  $F^\ell$  is  $(C_\ell, M_\ell)$ . Then, the Gram matrix has shape  $(C_\ell, C_\ell)$  where each element is given by the equation:

$$G_{ij}^\ell = \sum_k F_{ik}^\ell F_{jk}^\ell$$

Assuming  $G^\ell$  is the Gram matrix from the feature map of the current image,  $A^\ell$  is the Gram Matrix from the feature map of the source style image, and  $w_\ell$  a scalar weight term, then the style loss for the layer  $\ell$  is simply the weighted Euclidean distance between the two Gram matrices:

$$L_s^\ell = w_\ell \sum_{i,j} \left( G_{ij}^\ell - A_{ij}^\ell \right)^2$$

In practice we usually compute the style loss at a set of layers  $\mathcal{L}$  rather than just a single layer  $\ell$ ; then the total style loss is the sum of style losses at each layer:

$$L_s = \sum_{\ell \in \mathcal{L}} L_s^\ell$$

Begin by implementing the Gram matrix computation function `gram_matrix` inside `cs231n/style_transfer_pytorch.py`:

Test your Gram matrix code. You should see errors less than 0.001.

```
from cs231n.style_transfer_pytorch import gram_matrix
def gram_matrix_test(correct):
    style_image = '%s/starry_night.jpg' % (STYLES_FOLDER)
    style_size = 192
    feats, _ = features_from_img(style_image, style_size, cnn)
    student_output = gram_matrix(feats[5].clone()).cpu().data.numpy()
    error = rel_error(correct, student_output)
    print('Maximum error is {:.3f}'.format(error))

gram_matrix_test(answers['gm_out'])
```

Maximum error is 0.000

Next, put it together and implement the style loss function `style_loss` in `cs231n/style_transfer_pytorch.py`

Test your style loss implementation. The error should be less than 0.001.

```

from cs231n.style_transfer_pytorch import style_loss
def style_loss_test(correct):
    content_image = '%s/tubingen.jpg' % (STYLES_FOLDER)
    style_image = '%s/starry_night.jpg' % (STYLES_FOLDER)
    image_size = 192
    style_size = 192
    style_layers = [1, 4, 6, 7]
    style_weights = [300000, 1000, 15, 3]

    c_feats, _ = features_from_img(content_image, image_size, cnn)
    feats, _ = features_from_img(style_image, style_size, cnn)
    style_targets = []
    for idx in style_layers:
        style_targets.append(gram_matrix(feats[idx].clone()))

    student_output = style_loss(c_feats, style_layers, style_targets,
style_weights).cpu().data.numpy()
    error = rel_error(correct, student_output)
    print('Error is {:.3f}'.format(error))

style_loss_test(answers['s1_out'])

```

Error is 0.000

## Part 1C: Total-variation regularization

It turns out that it's helpful to also encourage smoothness in the image. We can do this by adding another term to our loss that penalizes wiggles or "total variation" in the pixel values.

You can compute the "total variation" as the sum of the squares of differences in the pixel values for all pairs of pixels that are next to each other (horizontally or vertically). Here we sum the total-variation regularization for each of the 3 input channels (RGB), and weight the total summed loss by the total variation weight,  $w_t$ :

$$L_{tv} = w_t \times \left( \sum_{c=1}^3 \sum_{i=1}^{H-1} \sum_{j=1}^W (x_{i+1,j,c} - x_{i,j,c})^2 + \sum_{c=1}^3 \sum_{i=1}^H \sum_{j=1}^{W-1} (x_{i,j+1,c} - x_{i,j,c})^2 \right)$$

In `cs231n/style_transfer_pytorch.py`, fill in the definition for the TV loss term in `tv_loss`. To receive full credit, your implementation should not have any loops.

Test your TV loss implementation. Error should be less than 0.0001.

```

from cs231n.style_transfer_pytorch import tv_loss
from inspect import getsourcelines
import re

def tv_loss_test(correct):
    content_image = '%s/tubingen.jpg' % (STYLES_FOLDER)
    image_size = 192
    tv_weight = 2e-2

    content_img = preprocess(PIL.Image.open(content_image),
size=image_size).type(dtype)

```

```

student_output = tv_loss(content_img, tv_weight).cpu().data.numpy()
error = rel_error(correct, student_output)
print('Error is {:.4f}'.format(error))
lines, _ = getsourcelines(tv_loss)
used_loop = any(bool(re.search(r"for\s*\s*in", line)) for line in lines)
if used_loop:
    print("WARNING!!!! - Your implementation of tv_loss contains a loop! To
receive full credit, your implementation should not have any loops")

tv_loss_test(answers['tv_out'])

```

Error is 0.0000

## Part 2: Style Transfer

Now we're ready to string it all together (you shouldn't have to modify this function):

```

def style_transfer(content_image, style_image, image_size, style_size,
content_layer, content_weight,
                    style_layers, style_weights, tv_weight, init_random = False):
    """
    Run style transfer!

    Inputs:
    - content_image: filename of content image
    - style_image: filename of style image
    - image_size: size of smallest image dimension (used for content loss and
generated image)
    - style_size: size of smallest style image dimension
    - content_layer: layer to use for content loss
    - content_weight: weighting on content loss
    - style_layers: list of layers to use for style loss
    - style_weights: list of weights to use for each layer in style_layers
    - tv_weight: weight of total variation regularization term
    - init_random: initialize the starting image to uniform random noise
    """

    # Extract features for the content image
    content_img = preprocess(PIL.Image.open(content_image),
size=image_size).type(dtype)
    feats = extract_features(content_img, cnn)
    content_target = feats[content_layer].clone()

    # Extract features for the style image
    style_img = preprocess(PIL.Image.open(style_image),
size=style_size).type(dtype)
    feats = extract_features(style_img, cnn)
    style_targets = []
    for idx in style_layers:
        style_targets.append(gram_matrix(feats[idx].clone()))

    # Initialize output image to content image or noise
    if init_random:
        img = torch.Tensor(content_img.size()).uniform_(0, 1).type(dtype)

```

```

else:
    img = content_img.clone().type(dtype)

# We do want the gradient computed on our image!
img.requires_grad_()

# Set up optimization hyperparameters
initial_lr = 3.0
decayed_lr = 0.1
decay_lr_at = 180

# Note that we are optimizing the pixel values of the image by passing
# in the img Torch tensor, whose requires_grad flag is set to True
optimizer = torch.optim.Adam([img], lr=initial_lr)

f, axarr = plt.subplots(1,2)
axarr[0].axis('off')
axarr[1].axis('off')
axarr[0].set_title('Content Source Img.')
axarr[1].set_title('Style Source Img.')
axarr[0].imshow(deprocess(content_img.cpu()))
axarr[1].imshow(deprocess(style_img.cpu()))
plt.show()
plt.figure()

for t in range(200):
    if t < 190:
        img.data.clamp_(-1.5, 1.5)
        optimizer.zero_grad()

        feats = extract_features(img, cnn)

        # Compute loss
        c_loss = content_loss(content_weight, feats[content_layer],
content_target)
        s_loss = style_loss(feats, style_layers, style_targets, style_weights)
        t_loss = tv_loss(img, tv_weight)
        loss = c_loss + s_loss + t_loss

        loss.backward()

        # Perform gradient descents on our image values
        if t == decay_lr_at:
            optimizer = torch.optim.Adam([img], lr=decayed_lr)
            optimizer.step()

        if t % 100 == 0:
            print('Iteration {}'.format(t))
            plt.axis('off')
            plt.imshow(deprocess(img.data.cpu()))
            plt.show()
print('Iteration {}'.format(t))
plt.axis('off')
plt.imshow(deprocess(img.data.cpu()))
plt.show()

```

# Generate some pretty pictures!

Try out `style_transfer` on the three different parameter sets below. Make sure to run all three cells. Feel free to add your own, but make sure to include the results of style transfer on the third parameter set (starry night) in your submitted notebook.

- The `content_image` is the filename of content image.
- The `style_image` is the filename of style image.
- The `image_size` is the size of smallest image dimension of the content image (used for content loss and generated image).
- The `style_size` is the size of smallest style image dimension.
- The `content_layer` specifies which layer to use for content loss.
- The `content_weight` gives weighting on content loss in the overall loss function. Increasing the value of this parameter will make the final image look more realistic (closer to the original content).
- `style_layers` specifies a list of which layers to use for style loss.
- `style_weights` specifies a list of weights to use for each layer in `style_layers` (each of which will contribute a term to the overall style loss). We generally use higher weights for the earlier style layers because they describe more local/smaller scale features, which are more important to texture than features over larger receptive fields. In general, increasing these weights will make the resulting image look less like the original content and more distorted towards the appearance of the style image.
- `tv_weight` specifies the weighting of total variation regularization in the overall loss function. Increasing this value makes the resulting image look smoother and less jagged, at the cost of lower fidelity to style and content.

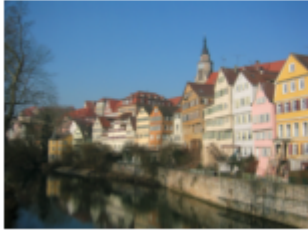
Below the next three cells of code (in which you shouldn't change the hyperparameters), feel free to copy and paste the parameters to play around them and see how the resulting image changes.

```
# Composition VII + Tubingen
params1 = {
    'content_image' : '%s/tubingen.jpg' % (STYLES_FOLDER),
    'style_image' : '%s/composition_vii.jpg' % (STYLES_FOLDER),
    'image_size' : 192,
    'style_size' : 512,
    'content_layer' : 3,
    'content_weight' : 5e-2,
    'style_layers' : (1, 4, 6, 7),
    'style_weights' : (20000, 500, 12, 1),
    'tv_weight' : 5e-2
}

style_transfer(**params1)
```



Content Source Img.



Style Source Img.



Iteration 0



Iteration 100



Iteration 199



```
# Scream + Tübingen
params2 = {
    'content_image': '%s/tubingen.jpg' % (STYLES_FOLDER),
    'style_image': '%s/the_scream.jpg' % (STYLES_FOLDER),
    'image_size': 192,
    'style_size': 224,
    'content_layer': 3,
    'content_weight': 3e-2,
    'style_layers': [1, 4, 6, 7],
    'style_weights': [200000, 800, 12, 1],
    'tv_weight': 2e-2
}

style_transfer(**params2)
```



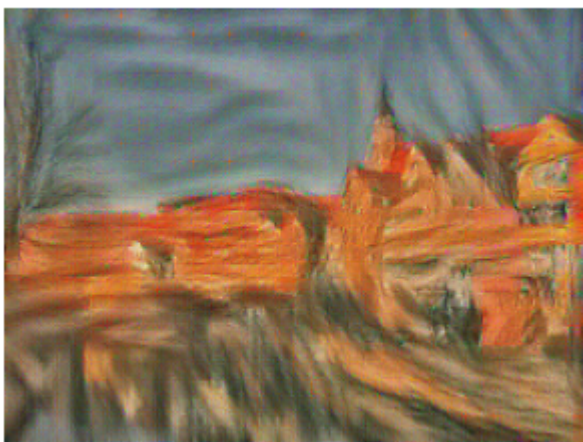
Iteration 0



Iteration 100



Iteration 199



```
# Starry Night + Tübingen  
params3 = {
```



```
'content_image' : '%s/tubingen.jpg' % (STYLES_FOLDER),  
'style_image' : '%s/starry_night.jpg' % (STYLES_FOLDER),  
'image_size' : 192,  
'style_size' : 192,  
'content_layer' : 3,  
'content_weight' : 6e-2,  
'style_layers' : [1, 4, 6, 7],  
'style_weights' : [300000, 1000, 15, 3],  
'tv_weight' : 2e-2  
}  
  
style_transfer(**params3)
```

Content Source Img.



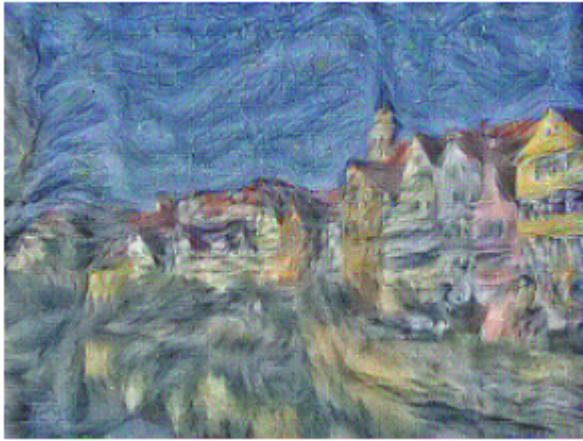
Style Source Img.



Iteration 0



Iteration 100



Iteration 199



## Part 3: Feature Inversion

---

The code you've written can do another cool thing. In an attempt to understand the types of features that convolutional networks learn to recognize, a recent paper "[Understanding Deep Image Representations by Inverting Them](#)" attempts to reconstruct an image from its feature representation. We can easily implement this idea using image gradients from the pretrained network, which is exactly what we did above (but with two different feature representations).

Now, if you set the style weights to all be 0 and initialize the starting image to random noise instead of the content source image, you'll reconstruct an image from the feature representation of the content source image. You're starting with total noise, but you should end up with something that looks quite a bit like your original image.

(Similarly, you could do "texture synthesis" from scratch if you set the content weight to 0 and initialize the starting image to random noise, but we won't ask you to do that here.)

Run the following cell to try out feature inversion.

[1] Aravindh Mahendran, Andrea Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR 2015

```
# Feature Inversion -- Starry Night + Tübingen
params_inv = {
    'content_image' : '%s/tubingen.jpg' % (STYLES_FOLDER),
    'style_image' : '%s/starry_night.jpg' % (STYLES_FOLDER),
    'image_size' : 192,
    'style_size' : 192,
    'content_layer' : 3,
    'content_weight' : 6e-2,
    'style_layers' : [1, 4, 6, 7],
    'style_weights' : [0, 0, 0, 0], # we discard any contributions from style to
the loss
    'tv_weight' : 2e-2,
    'init_random': True # we want to initialize our image to be random
}

style_transfer(**params_inv)
```

Content Source Img.



Style Source Img.



Iteration 0



Iteration 100



Iteration 199

