

# Project Report: Image-to-Image Translation

Gurdit Singh Siyan  
[gurdit@cmi.ac.in](mailto:gurdit@cmi.ac.in)  
[gurdit.siyang@gmail.com](mailto:gurdit.siyang@gmail.com)

*Chennai Mathematical Institute*

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Background</b>	<b>6</b>
2.1	Neural Networks . . . . .	6
2.1.1	Computing with Neural Networks . . . . .	6
2.1.2	Learning Parameters for Neural Networks . . . . .	7
2.2	Convolutional Neural Networks . . . . .	8
2.2.1	Convolution Layers . . . . .	8
2.2.2	Pooling Layers . . . . .	9
2.2.3	Architecture of a CNN . . . . .	9
2.3	Generative Adversarial Networks . . . . .	10
2.3.1	Training the Network . . . . .	10
2.3.2	Theoretical Results . . . . .	11
<b>3</b>	<b>Supervised Image-to-Image Translation</b>	<b>14</b>
3.1	Pix2Pix . . . . .	15
3.2	Objective . . . . .	15
3.3	Discriminator Architecture . . . . .	16
3.4	Generator Architecture . . . . .	17
3.5	Experiments . . . . .	18
3.5.1	Satellite2Map . . . . .	18
3.5.2	Edges2Shoes . . . . .	22
3.5.3	Anime Sketch Colorization . . . . .	23
<b>4</b>	<b>Unsupervised Image-to-Image Translation</b>	<b>26</b>
4.1	CycleGAN . . . . .	27
4.2	Objective . . . . .	27
4.3	Discriminator Architecture . . . . .	28
4.4	Generator Architecture . . . . .	28
4.5	Experiments . . . . .	29
4.5.1	Horse2Zebra . . . . .	29
<b>A</b>	<b>PyTorch Implementation</b>	<b>31</b>

---

A.1	constants.py . . . . .	31
A.2	datasets.py . . . . .	31
A.3	networks.py . . . . .	33
A.4	models.py . . . . .	36
A.5	pix2pix.py . . . . .	39
A.6	cyclegan.py . . . . .	42
A.7	models.py . . . . .	46
	<b>Bibliography</b>	<b>49</b>

# Chapter 1

## Introduction

Many problems in image processing, graphics, and vision involve translating an input image into a corresponding output image. There exist a wide variety of domain specific solutions for many problems, even though the goal is always the same: mapping pixels to pixels. The objective is always to find a map from an input domain  $X$  to an output domain  $Y$ . That is, for image domains  $X$  and  $Y$ , we want to find a function  $G : X \rightarrow Y$ . One of the questions we can ask is: does there exist a general solution to this problem?

Solving such a general problem using traditional Computer Vision techniques and algorithms can indeed feel impossible. However, there has been a relatively recent breakthrough in the field of deep learning. The introduction of a class of generative deep learning models known as Generative Adversarial Networks [3], can be used to approach generative problems which were considered very hard(or downright impossible) for a long time.

The goal of this project is to implement two of the most influential models in two settings of the Image-to-Image translation problem:

- **Supervised setting:** In this setting we have paired input-output images between domain  $X$  and domain  $Y$ .
- **Unsupervised setting:** In this harder setting we just have two image domains(sets)  $X$  and  $Y$  and paired examples are unavailable.



Figure 1.1: Colorization: Black & White → Color



Figure 1.2: Semantic Segmentation: Cityscape → Segmented Image



Figure 1.3: Style Transfer: Portrait → Stylized Image



Figure 1.4: Background Removal: Image → Image without Background

# Chapter 2

## Background

### 2.1 Neural Networks

Neural networks [2], also called **feedforward neural networks** or **multi-layer perceptrons** (MLPs) are a network or a circuit composed of artificial neurons. The objective for a neural network is to approximate some function  $f^*$ . For example, for a classifier,  $y = f^*(\mathbf{x})$  maps an input  $\mathbf{x}$  to a category  $y$ . A neural network defines a mapping  $\hat{y} = f(\mathbf{x}; \boldsymbol{\theta})$  and learns the value of the weights(or parameters)  $\boldsymbol{\theta}$  that result in the best function approximation.

These models are called **networks** because they are represented by composing together many different functions called **layers**. And, these models are called **feedforward** because to evaluate the network at some value  $\mathbf{x}$  we pass it through a sequence of layers whose result depend upon the values of previous layers, i.e., under evaluation the computation flows in one direction only.

#### 2.1.1 Computing with Neural Networks

The model is associated with a directed acyclic graph describing how the functions are composed together. For example, in the figure 2.1 we have 3 layers(functions)  $f^{(1)}$ ,  $f^{(2)}$  and  $f^{(3)}$  connected in a chain to form the network  $f(\mathbf{x}) = f^{(1)}(f^{(2)}(f^{(3)}(\mathbf{x})))$ , the final function  $f^{(3)}$  is called the output layer.

These layer functions  $f^{(i)}$  in general look like:

$$f^{(i)}(\mathbf{a}^{(i-1)}) = g^{(i)}((\mathbf{W}^{(i)})^T \mathbf{a}^{(i-1)} + \mathbf{b}^{(i)}) \quad (2.1)$$

Here  $\mathbf{a}^{(i-1)}$  is the result vector of computation of the previous layer, note that  $\mathbf{a}^{(0)}$  is defined to be  $\mathbf{x}$ .  $\mathbf{W}^{(i)}$  is defined to be the weight matrix for the layer  $i$  where each column correspond to weighted edges connected to a neuron, and  $\mathbf{b}^{(i)}$  is the bias vector for the layer  $(i)$ . To elaborate if the number of elements

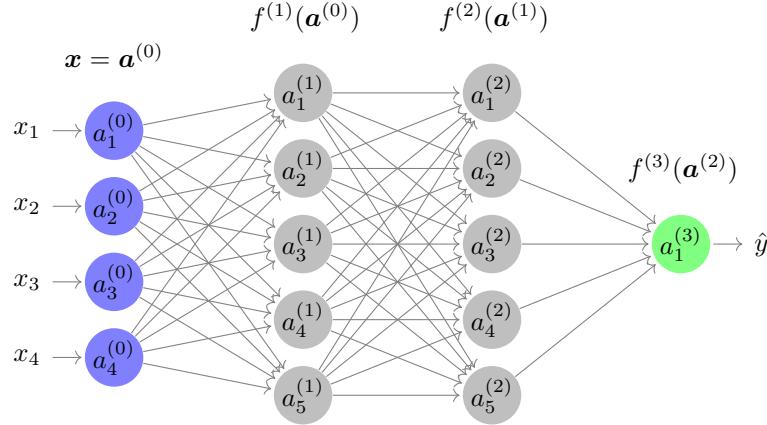


Figure 2.1: A simple neural network.

in vector  $\mathbf{a}^{(i-1)}$  is  $n^{(i-1)}$  and in  $\mathbf{a}^{(i)}$  is  $n^{(i)}$ , then the shape of matrix  $\mathbf{W}^{(i)}$  is  $(n^{(i-1)} \times n^{(i)})$  and length for the vector  $\mathbf{b}$  is  $n^{(i)}$ . Also,  $g^{(i)} : \mathbb{R} \rightarrow \mathbb{R}$  is just a real valued function which is applied(broadcasted) to every element of the input vector. Some commonly used activation functions are:

- **Linear** is just the identity function.

$$g(x) = x$$

- **ReLU** (Rectified Linear Unit) activation function is an activation function defined as the positive part of its argument:

$$g(x) = \max(0, x)$$

- **Sigmoid** function is a mathematical function having a characteristic “S”-shaped curve or sigmoid curve. Some examples are:

<b>Logistic:</b> <b>Hyperbolic Tangent:</b>	$g(x) = \frac{1}{1 + e^{-x}}$ $g(x) = \tanh x$
--	--

### 2.1.2 Learning Parameters for Neural Networks

Neural networks are trained similarly to other machine learning models and that is via Gradient Descent.

Let the training set for the model be defined as:

$$\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(i)}, \mathbf{y}^{(i)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}$$

---

Where  $\mathbf{x}^{(i)}$  is the input and  $\mathbf{y}^{(i)}$  is the corresponding output. We also need to define a **cost function**  $J$  which for parameters  $\boldsymbol{\theta}$  measures the distance between the predicted vector  $f(\mathbf{x}^{(i)}) = \hat{\mathbf{y}}^{(i)}$  and the actual vector  $\mathbf{y}^{(i)}$ . One common loss function is the **Mean Squared Error** (MSE) loss, defined as:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)})^2$$

Then we use any form of gradient based learning to minimize this loss function. For example, for **Stochastic Gradient Descent**(SGD) we estimate the weights as follows:

---

**Algorithm 1** Stochastic Gradient Descent training for Neural Networks

---

```

• Initialize random weights  $\boldsymbol{\theta}$  and chose learning rate  $\eta$ .
for epoch  $\leftarrow$  number of training iterations do
    for  $i \leftarrow$  number of training samples do
        • Update  $\boldsymbol{\theta}$  by descending its gradient evaluated at sample  $i$ .
        
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \cdot \nabla_{\boldsymbol{\theta}} J^{(i)}(\boldsymbol{\theta})$$

    end for
end for

```

---

## 2.2 Convolutional Neural Networks

**Convolutional Neural Networks** [2] or CNNs, are a specialized form of Neural Network for processing input which has a sequential grid-like format. For example, in one dimensional data we can use CNNs for time-series data, natural language represented as a sequence of word vectors, etc., and in two dimensions we most commonly use CNNs for images. We will only discuss 2D CNNs in this section.

The input to a CNN is a tensor of shape  $m \times h \times w \times c$ , where  $m$  is the number of samples,  $h$  and  $w$  are the height and width of the input image and  $c$  is the number of channels in the image.

For CNNs, in addition to **fully connected** layers as discussed in the previous section, we also have extra kind of layers such as **convolution** and **pooling** layers.

### 2.2.1 Convolution Layers

Convolution layers “convolve” the input and pass its result to the next layer. The convolution operation is a mathematical operation defined as follows, if we have an input image  $I$ , and we apply a convolution filter  $K$ (which is also a

---


$$\begin{bmatrix}
 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\
 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 \\
 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\
 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\
 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0
 \end{bmatrix}_I * \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}_K = \begin{bmatrix} 1 & 4 & 3 & 4 & 1 \\ 1 & 2 & 4 & 3 & 3 \\ 1 & 2 & 3 & 4 & 1 \\ 1 & 3 & 3 & 1 & 1 \\ 3 & 3 & 1 & 1 & 0 \end{bmatrix}_{I * K}$$

Figure 2.2: Cross-correlation in 2D.

matrix), the resulting computation would be a matrix  $S = I * k$ . Note: here  $*$  is the convolution operator.

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i - m, j - n)K(m, n)$$

However, in practice we implement a related function called the **cross-correlation** 2.2, which is the same as convolution except we don't flip the kernel.

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n)K(m, n)$$

For CNNs the entire matrix  $K$  is learned during training.

### 2.2.2 Pooling Layers

After passing the input to some convolution layers and activation functions, we also use a new kind of layer called the pooling layer. A pooling layer, like an activation function, contains no learnable parameters but rather is just a normal function.

Similar to a convolution, pooling layer is applied to a rectangular region and returns a single number. Some common examples of pooling layers include max pooling, which takes the maximum value in a region, average pooling which takes the average value, weighted average,  $L^2$  norm of the region etc.

### 2.2.3 Architecture of a CNN

A CNN usually consist of many alternating convolution-pooling layer pairs and ultimately connected to a fully connected layer. Learning the weights for CNN work similarly to a normal neural network.

## 2.3 Generative Adversarial Networks

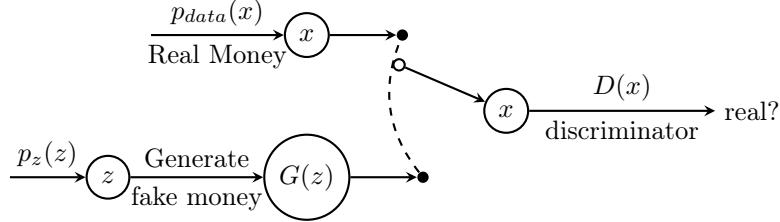


Figure 2.3: A visual representation of Generative Adversarial Network

**Generative adversarial networks** [3] are a framework for estimating generative models via an adversarial process, in which two models: a generative model  $\mathbf{G}$  that captures the data distribution, and a discriminative model  $\mathbf{D}$  that estimates the probability that the sample came from the training data rather than  $\mathbf{G}$ , are trained simultaneously.

The canonical example to understand GANs is as follows: imagine  $G$  as a fraud trying to create fake currency and use it without detection and imagine  $D$  is working for the police who is trying to detect counterfeit money.

The goal of the generator is to generate samples that are indistinguishable from the real data, i.e.,  $p_G(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$ .

To build the optimization function for the discriminator  $D$ , we want to make it recognize the samples from the actual distribution. That is, we try to maximize:

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \log D(\mathbf{x})$$

But we also want to make it so that the discriminator  $D$  rejects samples generated by  $G$ :

$$\mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} \log(1 - D(G(\mathbf{z})))$$

Maximizing these terms for  $D$  mean  $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$  implies  $D(\mathbf{x}) \approx 1$ , and  $D(G(\mathbf{z})) \approx 0$ .

The generator  $G$  will try to minimize these terms, which makes this a two-player minimax game with the value function  $V(G, D)$ :

$$\min_G \max_D V(G, D) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \log D(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} \log(1 - D(G(\mathbf{z}))) \quad (2.2)$$

### 2.3.1 Training the Network

The Generative Adversarial Network can be trained by training both the discriminator  $D$  and the generator  $G$  on a turn by turn basis.

Note, in practice rather than minimizing(maximizing)  $\log(1 - D(G(\mathbf{z})))$  we maximize(minimize)  $\log(D(G(\mathbf{z})))$  which provides much stronger gradients early in the training when  $G$  is weak, but results in the same fixed point dynamics.

---

**Algorithm 2** Minibatch SGD training of Generative Adversarial Networks

---

```
for number of training iterations do
    for k steps do
        • Sample minibatch m noise examples { $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)}$ } from  $p_z(\mathbf{z})$ .
        • Sample minibatch m examples { $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}$ } from the actual distribution  $p_{\text{data}}(\mathbf{x})$ .
        • Update the Discriminator  $D$  by ascending its stochastic gradient:
            
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(\mathbf{x}^{(i)}) + \log \left( 1 - D(G(\mathbf{z}^{(i)})) \right) \right]$$

    end for
    • Sample minibatch m noise examples { $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)}$ } from  $p_z(\mathbf{z})$ .
    • Update the Generator  $G$  by descending its stochastic gradient:
        
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log \left( 1 - D(G(\mathbf{z}^{(i)})) \right) \right]$$

end for
```

---

### 2.3.2 Theoretical Results

In the space of arbitrary functions  $G$  and  $D$ , a unique solution exists, with  $G$  recovering the training data and  $D$  being equal to  $\frac{1}{2}$  everywhere.

**Proposition 1.** *For a fixed  $G$ , the optimal discriminator  $D$  is*

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \quad (2.3)$$

*Proof.* The training criterion for  $D$  is to maximize the following function:

$$\begin{aligned} V(G, D) &= \int_x p_{\text{data}}(\mathbf{x}) \log D(\mathbf{x}) dx + \int_z p_z(\mathbf{z}) \log(1 - D(G(\mathbf{z}))) dz \\ &= \int_x \left( p_{\text{data}}(\mathbf{x}) \log D(\mathbf{x}) + p_G(\mathbf{x}) \log(1 - D(G(\mathbf{x}))) \right) dx \end{aligned}$$

The equation reduces to finding the maximum for the equation

$$f(y) = a \log y + b \log(1 - y)$$

To find the maxima, we calculate the derivatives:

$$\begin{aligned}f'(y) &= \frac{a}{y} - \frac{b}{1-y} \\f''(y) &= -\frac{a}{y^2} - \frac{b}{(1-y)^2}\end{aligned}$$

We have the extreme point when,

$$f'(y) = 0 \Rightarrow y = \frac{a}{a+b}$$

But also note:

$$f''\left(\frac{a}{a+b}\right) = -\frac{a}{\left(\frac{a}{a+b}\right)^2} - \frac{b}{\left(1-\left(\frac{a}{a+b}\right)\right)^2} < 0$$

So, the function  $f$  is maximum at  $\frac{a}{a+b}$ , which means  $V(G, D)$  is maximum when you let

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})}$$

□

By the previous proposition we can see that for an optimal generator  $G$  i.e., when  $p_G(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$ , there exists a unique discriminator  $D$ , which is:

$$\begin{aligned}D_G^*(\mathbf{x}) &= \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \\&= \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{data}}(\mathbf{x})} \\&= \frac{1}{2}\end{aligned}$$

**Theorem 1.** *The global minimum of the virtual training criterion  $C(G) = \max_D V(G, D)$  is reached if and only if  $p_G(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$ , and the minimum value is  $-\log(4)$*

*Proof.* ( $\Leftarrow$ ) Assuming  $p_G(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$ ,

$$\begin{aligned}V(G, D_G^*) &= \int_x \left( p_{\text{data}}(\mathbf{x}) \log \frac{1}{2} + p_{\text{data}}(\mathbf{x}) \log(1 - \frac{1}{2}) \right) dx \\&= -2 \log 2 \int_x p_{\text{data}}(\mathbf{x}) dx \\&= -2 \log 2 \\&= -\log 4\end{aligned}$$

---

We will later prove that this is the absolute minimum.

( $\Rightarrow$ ) Plugging in value for  $D_G^*$  we get,

$$\begin{aligned} C(G) &= \int_x \left( p_{\text{data}}(\mathbf{x}) \log \left( \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right. \\ &\quad \left. + p_G(\mathbf{x}) \log \left( 1 - \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right) dx \\ &= \int_x \left( p_{\text{data}}(\mathbf{x}) \log \left( \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right. \\ &\quad \left. + p_G(\mathbf{x}) \log \left( \frac{p_G(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right) dx \end{aligned}$$

Adding and subtracting  $\log 2$

$$\begin{aligned} &= \int_x \left( p_{\text{data}}(\log 2 - \log 2) + p_{\text{data}}(\mathbf{x}) \log \left( \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right. \\ &\quad \left. + p_G(\log 2 - \log 2) + p_G(\mathbf{x}) \log \left( \frac{p_G(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right) dx \\ &= -\log 2 \int_x p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x}) \\ &\quad + \int_x p_{\text{data}}(\mathbf{x}) \log \left( \frac{p_{\text{data}}(\mathbf{x})}{(p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x}))/2} \right) dx \\ &\quad + \int_x p_G(\mathbf{x}) \log \left( \frac{p_G(\mathbf{x})}{(p_{\text{data}}(\mathbf{x}) + p_G(\mathbf{x}))/2} \right) dx \\ &= -\log 4 + KL \left( p_{\text{data}} \middle| \frac{p_{\text{data}} + p_G}{2} \right) + KL \left( p_G \middle| \frac{p_{\text{data}} + p_G}{2} \right) \end{aligned}$$

Here  $KL$  is the Kullback-Leibler divergence, which is always non-negative, and thus validates our assumption that  $-\log 4$  is the global minimum.

Also note, we can rewrite the last equation as

$$= -\log 4 + 2 \cdot JSD(p_{\text{data}} \mid p_G)$$

where  $JSD$  is the Jenson-Shannon divergence, which is 0 only when  $p_{\text{data}} = p_G$ , which means the global minimum occurs when  $p_{\text{data}} = p_G$ , i.e., the generative model perfectly replicates the training data.  $\square$

## Chapter 3

# Supervised Image-to-Image Translation

In supervised Image-to-Image translation there exist “paired” examples of input-output images for the problem we are solving.

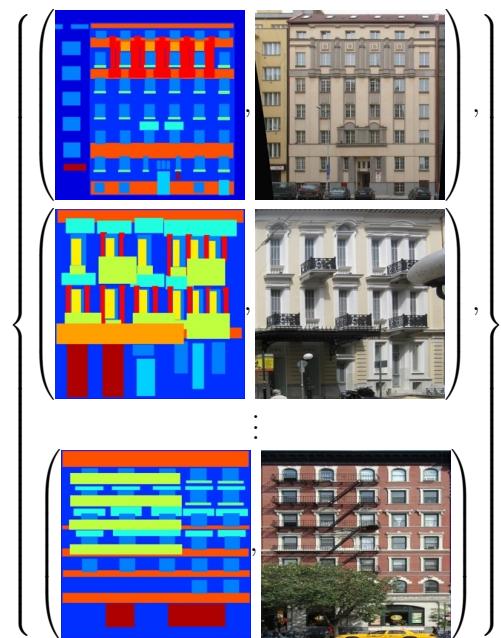


Figure 3.1: Supervised Image-to-Image Translation



Figure 3.2: Example output on Satellite2Map dataset

### 3.1 Pix2Pix

**Pix2Pix** was one of the first general-purpose Image-to-Image translation network which was introduced in the seminal paper “Image-to-Image Translation with Conditional Adversarial Networks” [4]. They applied the same architecture on a wide variety of problems and got very promising results.

The first question we should ask before using GANs for this problem is: do we need GANs at all? or can we just use a simple CNN? There exist several generative CNN models — for example Encoder-Decoder, and U-Net — so why don’t we use them? The problem is that CNNs learn to minimize a loss function and although the learning process is automatic, we still need to carefully design the loss function. That is, we need to tell the CNN what we want to minimize.

But if we take a naive approach and tell the CNN to minimize the Euclidean distance between the input image and output image, then we will get blurry results. This happens because it is much “safer” for the L2 loss to predict the mean of the distribution, because this minimizes the mean pixel-wise error, but results in a blurry averaged image. [12]

Rather than specifying a carefully crafted model and loss for each Image-to-Image problem. It would be better if we could specify a goal like: *“make the output indistinguishable from reality”* While this may seem like a vague goal, this is exactly what Generative Adversarial Networks claim to do [3].

The main idea in the **Pix2Pix** paper was to use a Generative Adversarial Network to learn the loss function, rather than hand-coding a loss function such as Euclidean distance between target and input images.

### 3.2 Objective

The goal is to generate a target image based on some source image and we don’t want the GAN to generate random images, so we use a variant of GANs called **Conditional GANs** [8]

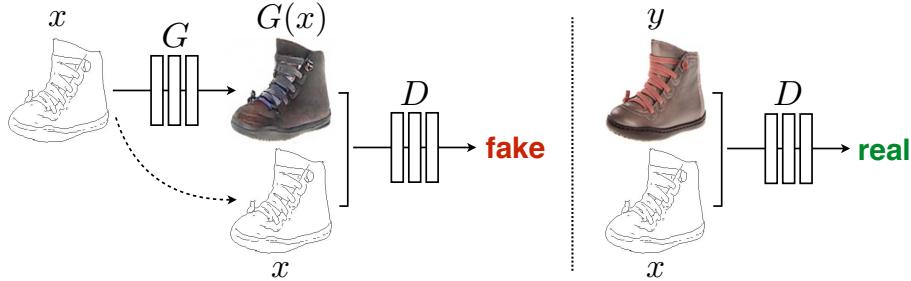


Figure 3.3: Conditional GAN used in Pix2Pix

As the name suggest Conditional GANs (figure 3.3) take an input to both the generator  $G$  and the discriminator  $D$ , for example, if we have a training pair  $(\mathbf{c}, \mathbf{y})$ , where we want to generate the item  $\mathbf{y}$  based on the value of  $\mathbf{c}$ , we pass  $\mathbf{c}$  to both  $G$  and  $D$ , then the loss function becomes:

$$\mathcal{L}_{\text{cGAN}}(G, D) = \mathbb{E}_{\mathbf{x}, \mathbf{c}}[\log D(\mathbf{c}, \mathbf{x})] + \mathbb{E}_{\mathbf{z}, \mathbf{c}}[\log(1 - D(\mathbf{c}, G(\mathbf{c}, \mathbf{z})))] \quad (3.1)$$

In their experiments they found mixing both the adversarial loss and a more traditional loss such as L1 was helpful for training the GAN. The loss function for the discriminator remains unchanged but the generator now not only has to fool the discriminator but also minimize the L1 loss between the predicted and target image.

$$\mathcal{L}_{\text{L1}}(G, D) = \mathbb{E}_{\mathbf{y}, \mathbf{c}, \mathbf{z}}[\|\mathbf{y} - G(\mathbf{c}, \mathbf{z})\|] \quad (3.2)$$

The final objective becomes for the generator becomes:

$$G^* = \arg \min_G \max_D \mathcal{L}_{\text{cGAN}}(G, D) + \lambda \mathcal{L}_{\text{L1}}(G, D) \quad (3.3)$$

For training the generator, instead of providing random noise as a direct input, they provide noise in the form of dropout.

### 3.3 Discriminator Architecture

Since we are already adding a traditional loss such as L1(or L2) which — although produce blurry results when used alone — are very useful for enforcing correctness at the low frequency. This is the main idea behind the **PatchGAN** [4] discriminator.

**PatchGAN** architecture restricts the discriminator to model only high frequency structure in the generated image. To model this, it is sufficient to restrict the attention to local image patches (figure 3.4). Therefore, rather than generating a single probability in range  $[0, 1]$ , the discriminator generates a  $k \times k$  matrix of values in range  $[0, 1]$  where each element in the matrix corresponds to a local patch(or receptive field) of the image.

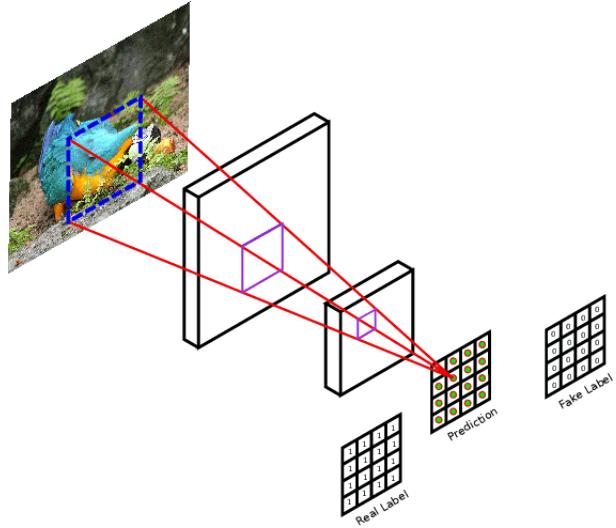


Figure 3.4: Visualization of **PatchGAN**

One advantage of **PatchGAN** is that the same architecture can extend to images of higher resolution with same number of parameters, and it would generate a similar, but larger,  $k' \times k'$  matrix.

Inspired by **DCGAN** [9], the discriminator is implemented as a Deep CNN mostly build up of repetitions of modules of form: Convolution - BatchNorm - ReLU, where we down-sample until we reach the required  $k \times k$  matrix, which we then pass through the logistic function to get the elements in range  $[0, 1]$ .

### 3.4 Generator Architecture

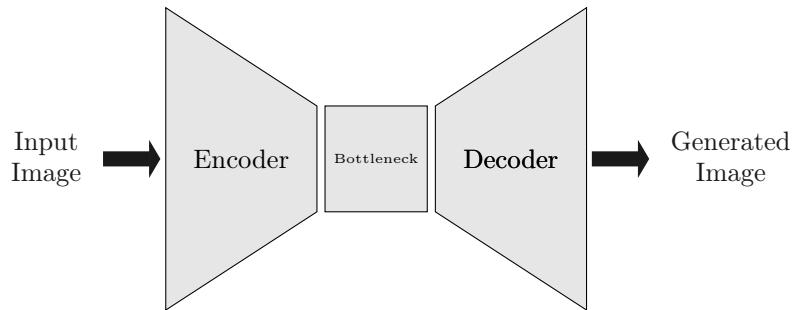


Figure 3.5: Encoder-Decoder Architecture

---

One of the obvious choices for a generator architecture for Image-to-Image translation, is the encoder-decoder network (figure 3.5). In this kind of network, the input image is passed through a series of layers which down-sample the resolution but increases the number of channels until we reach an image of size  $1 \times 1$ , after that the process is reversed until we reach the resolution of the original image. Here also we use modules of the form: Convolution - BatchNorm - LeakyReLU for encoder, for decoder we use: TransposeConvolution - BatchNorm - ReLU.

For many interesting image-translation problems, there is a lot of low-level information shared between the input and output image. For example, in the case of colorization, the input and output share the location of edges. To preserve this information, we add skip connections to create a “U-net” [10] like architecture.

This network is similar to Encoder-Decoder network , but we add some extra skip-connections between layer  $i$  in the encoder and layer  $n - i$  in the decoder.

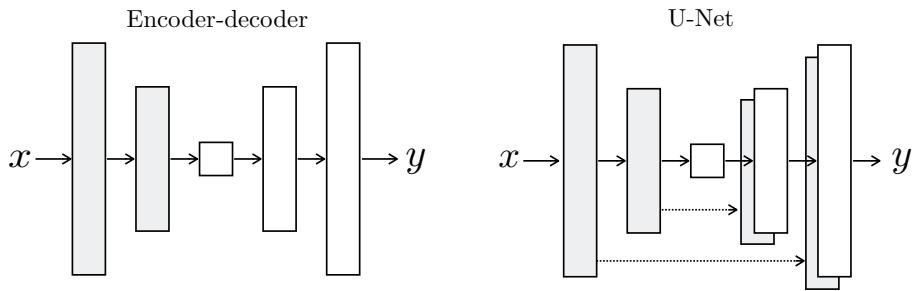


Figure 3.6: Difference between Encoder-Decoder and U-Net architectures.

## 3.5 Experiments

To demonstrate the versatility of the **Pix2Pix** approach, I applied the model on three different datasets with seemingly very different problem setting.

The datasets used are: **Satellite2Map** [4], **Edges2Shoes** [4], **Anime Sketch Colorization** (link). All datasets were trained on a single GPU and took  $\approx 1$  day to train.

I used the **Adam** Optimizer [6] with hyperparameters  $\beta_1 = 0.5$  and  $\beta_2 = 0.999$ . Learning rate was set to 0.0002, and  $\lambda$  used was 100.

### 3.5.1 Satellite2Map

This dataset [4] consists of paired images of aerial photography and their corresponding map view scraped using google maps. There are 1096 training images of size  $256 \times 256$ . The model was trained for 150 epochs with batch size 1.



Figure 3.7: Example output on Satellite2Map dataset

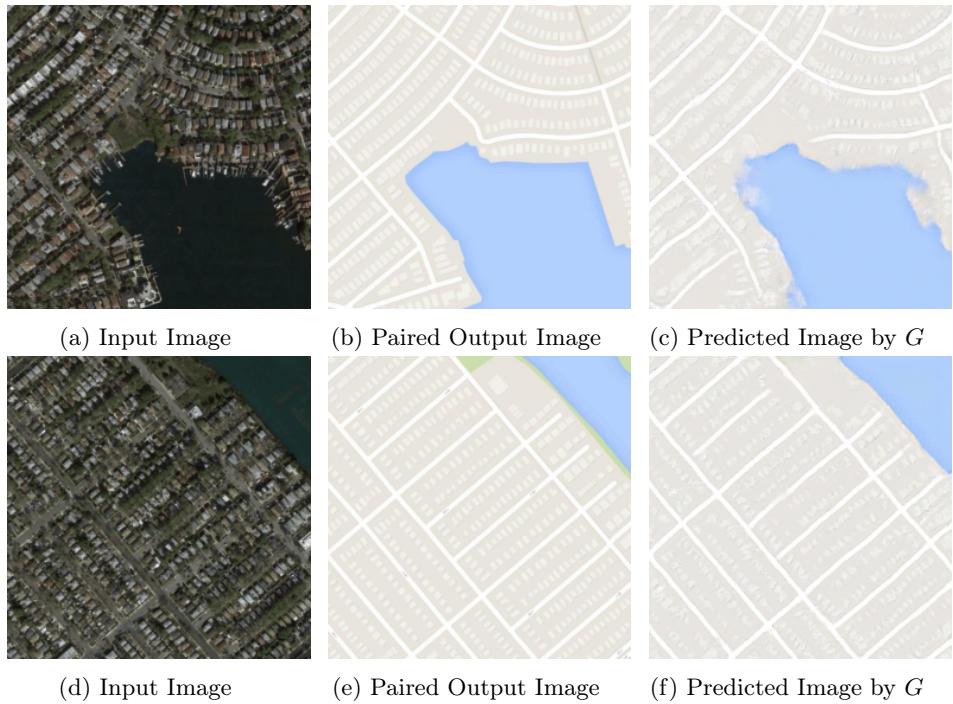


Figure 3.8: Example output on Satellite2Map dataset

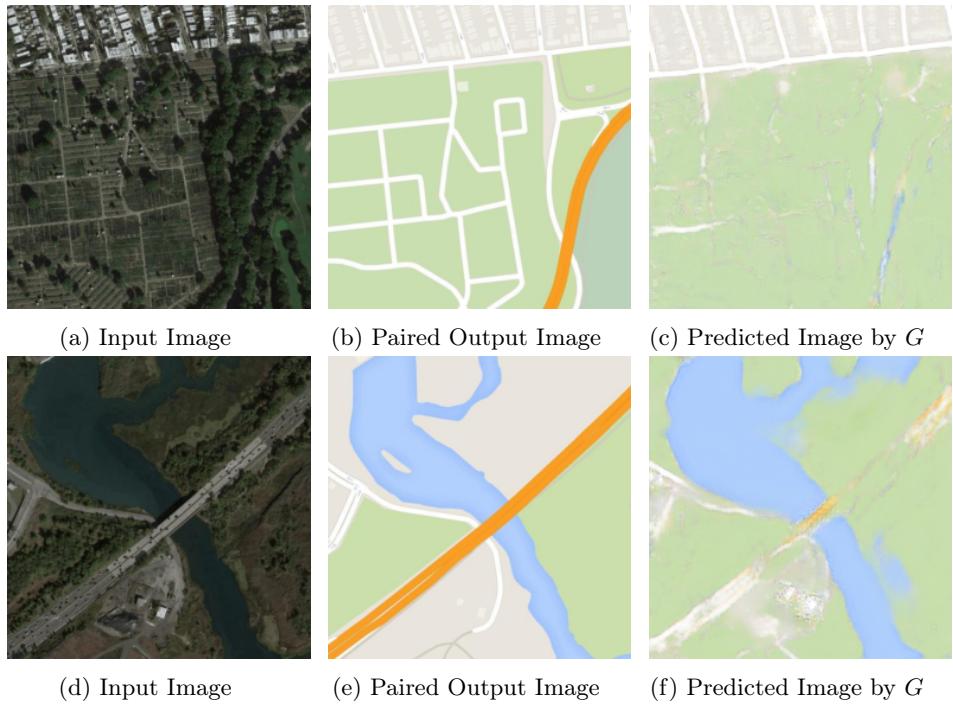


Figure 3.9: Example bad cases on Satellite2Map dataset

---

### 3.5.2 Edges2Shoes

This dataset [4] consists of paired images of shoe sketch image and the corresponding actual image (created using edge detection). There are 50k training images of size  $256 \times 256$ . The model was trained for 30 epochs with batch size 1.



Figure 3.10: Example output on Edges2Shoes dataset



Figure 3.11: Example output on Edges2Shoes dataset

### 3.5.3 Anime Sketch Colorization

This dataset ([link](#)) consists of paired images of anime line sketch images and the corresponding colored image. There are 14k training images of size  $256 \times 256$ . The model was trained for 75 epochs with batch size 1.



Figure 3.12: Example output on Anime Sketch Colorization dataset

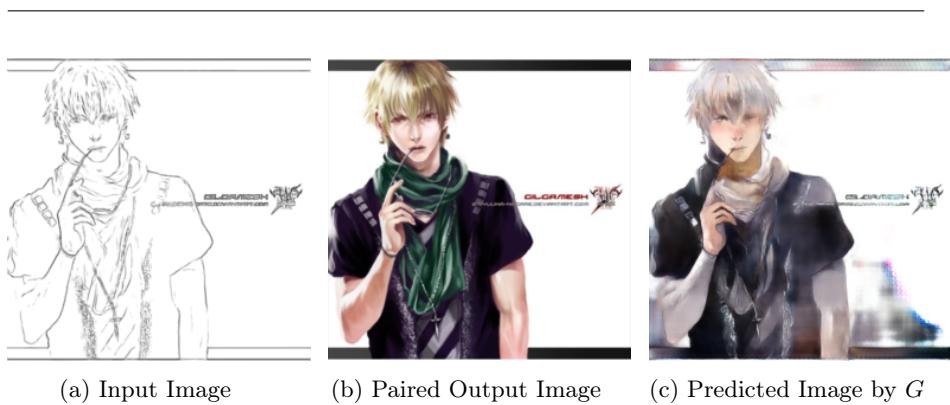


Figure 3.13: Example output on Anime Sketch Colorization dataset

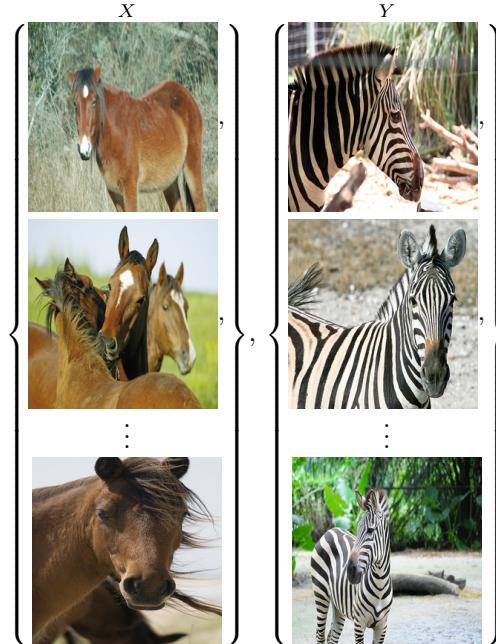


Figure 3.14: Example bad case on Anime Sketch Colorization dataset

## Chapter 4

# Unsupervised Image-to-Image Translation

In unsupervised Image-to-Image translation there does not exist “paired” examples of input-output images, but only a “domain” (set) of images  $X$  and  $Y$  and thus, usual models cannot be applied.



## 4.1 CycleGAN

In the previous chapter we saw how we can utilize Conditional Generative Adversarial Networks [8] to tackle the supervised problem. Surprisingly, there is a way to extend the previous approach to work in an unsupervised setting!

**CycleGAN** [13] is a new model which allows for approaching the Image-to-Image translation problem from an unsupervised learning point of view. It is built using two GANs, capturing both transformations between domains as  $G : X \rightarrow Y$  and also  $F : Y \rightarrow X$ .

The main idea behind this model is to add a cycle-consistency loss in addition to an adversarial loss. That is, if  $K \in X$  and  $K' \in Y$ , then if we apply  $G$  and  $F$  sequentially on  $K$  and vice-versa on  $K'$ , then we should get the same image back:

$$K \approx F(G(K)) \text{ and } K' \approx G(F(K'))$$

## 4.2 Objective

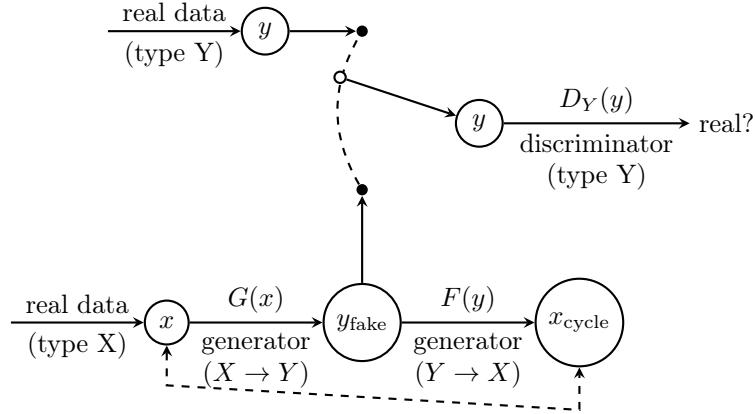


Figure 4.1: CycleGAN Visual Representation

The objective for the full CycleGAN model is the combination of adversarial loss for the two GANs:

$$\mathcal{L}_{\text{GAN}}(G, D_Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)}[\log D_Y(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log(1 - D_Y(G(x)))]$$

$$\mathcal{L}_{\text{GAN}}(F, D_X) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D_X(x)] + \mathbb{E}_{y \sim p_{\text{data}}(y)}[\log(1 - D_X(F(y)))]$$

And an additional cycle-consistency loss:

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) &= \mathbb{E}_{x \sim p_{\text{data}}(x)} \|F(G(x)) - x\|_1 \\ &\quad + \mathbb{E}_{y \sim p_{\text{data}}(y)} \|G(F(y)) - y\|_1 \end{aligned}$$

---

The final objective becomes:

$$\begin{aligned} G^*, F^* &= \arg \min_{G, F} \max_{D_X, D_Y} \mathcal{L}(G, F, D_X, D_Y) \\ &= \arg \min_{G, F} \max_{D_X, D_Y} \mathcal{L}_{\text{GAN}}(G, D_Y) + \mathcal{L}_{\text{GAN}}(F, D_X) + \lambda \mathcal{L}_{\text{cyc}}(G, F) \end{aligned}$$

### 4.3 Discriminator Architecture

The discriminator architecture is mostly the same as in the **Pix2Pix** model, but there are a few changes:

- We replace BatchNorm with InstanceNorm [11]
- We replace the negative log-likelihood objective by a least-squares loss, as this loss was found to be more stable during training and generate higher quality images [7].

That is, we now train  $G$  to minimize:

$$\mathbb{E}_{x \sim p_{\text{data}}(x)}[(D_Y(G(x)) - 1)^2]$$

And we train  $D_Y$  to minimize:

$$\mathbb{E}_{y \sim p_{\text{data}}(y)}[(D_Y(y) - 1)^2] + \mathbb{E}_{x \sim p_{\text{data}}(x)}[D_Y(G(x))^2]$$

The objective is similarly replaced for  $F$  and  $D_X$ .

### 4.4 Generator Architecture

The generator has a completely new architecture inspired networks presented in Perceptual Loss paper [5].

Similarly to the previous generator, we build the model in modules of Convolution - InstanceNorm - ReLU, but instead of the usual Encoder-Decoder network we have a network consisting of:

- 3 Convolution - InstanceNorm - ReLU Blocks.
- Several residual blocks (9 in the case of  $256 \times 256$  images).
- 2 fractionally-strided convolution blocks to sample up.
- And finally a convolution block to map to RGB (or any variable number of channels).

---

## 4.5 Experiments

**CycleGAN** is a considerably heavier model than pix2pix and consists of at least twice as many parameters, but, there is also additional complexity since the two GANs are in some sense learning from each-other.

The only dataset used for this model is the **Horse2Zebra** [13] dataset. It was trained on a single GPU and took  $\approx 1$  day to train.

I used the **Adam** Optimizer [6] with hyperparameters  $\beta_1 = 0.5$  and  $\beta_2 = 0.999$ . Learning rate was set to 0.0002, and  $\lambda$  used was 10.

### 4.5.1 Horse2Zebra

This dataset [13] consists of two unpaired sets of images of horses and zebras, extracted from the ImageNet dataset [1]. There are 1067 images of horses and 1334 images of zebras, all of size  $256 \times 256$ . The model was trained for 10 epochs with batch size 1.

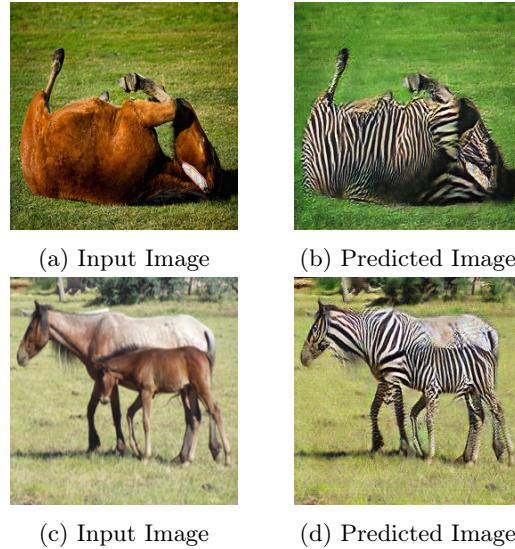
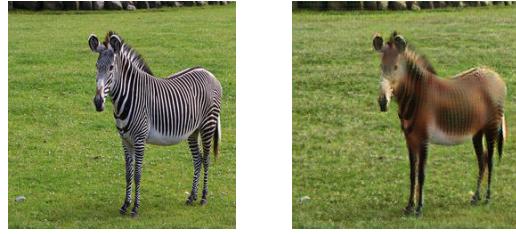


Figure 4.2: Horse  $\rightarrow$  Zebra Examples



(a) Input Image

(b) Predicted Image



(c) Input Image

(d) Predicted Image

Figure 4.3: Zebra → Horse Examples



(a) Input Image

(b) Predicted Image

Figure 4.4: Horse → Zebra Bad Case



(a) Input Image

(b) Predicted Image

Figure 4.5: Zebra → Horse Bad Case

## Appendix A

# PyTorch Implementation

### A.1 constants.py

```
1 import torch
2 import torchvision.transforms as transforms
3
4 torch.backends.cudnn.benchmark = True
5 DEVICE = torch.device("cuda" if torch.cuda.is_available() else "cpu")
6
7 NUM_WORKERS = 2
8 LOAD_CHECKPOINT = True
9 SAVE_CHECKPOINT = False
10
11 NUM_EPOCHS = 200
12 BATCH_SIZE = 1
13
14 LEARNING_RATE = 0.0002
15 L1_LAMBDA = 100
16 L1_LAMBDA_CYCLE = 10
17 L1_LAMBDA_IDENTITY = 0
18 BETAS = (0.5, 0.999)
19
20 BASE_TRANSFORMATIONS = transforms.Compose([
21     transforms.Resize((256, 256)),
22     transforms.ToTensor(),
23     transforms.Normalize(mean=[0.5, 0.5, 0.5], std=[0.5, 0.5, 0.5])
24     ,
25 ])
26
```

### A.2 datasets.py

```
1 import constants
2 import os
3 from PIL import Image
4 from torch.utils.data import Dataset
5
```

```

6  class SideBySideDataset(Dataset):
7      def __init__(self, root_dir, output = 'right'):
8          self.root_dir = root_dir
9          self.list_files = os.listdir(self.root_dir)
10         self.output = output
11
12     def __len__(self):
13         return len(self.list_files)
14
15     def __getitem__(self, index):
16         img_file = self.list_files[index]
17         img_path = os.path.join(self.root_dir, img_file)
18         image = Image.open(img_path).convert('RGB')
19
20         w,h = image.size
21
22         img1 = image.crop((0, 0, w//2, h))
23         img2 = image.crop((w//2, 0, w, h))
24
25         img1 = constants.BASE_TRANSFORMATIONS(img1)
26         img2 = constants.BASE_TRANSFORMATIONS(img2)
27
28         if self.output == 'right':
29             input_image, output_image = img1, img2
30         elif self.output == 'left':
31             input_image, output_image = img2, img1
32
33         return input_image, output_image
34
35 class UnpairedDataset(Dataset):
36     def __init__(self, root_dir_X, root_dir_Y):
37         self.root_dir_X = root_dir_X
38         self.list_files_X = os.listdir(self.root_dir_X)
39
40         self.root_dir_Y = root_dir_Y
41         self.list_files_Y = os.listdir(self.root_dir_Y)
42
43     def __len__(self):
44         return max(len(self.list_files_X), len(self.list_files_Y))
45
46     def __getitem__(self, index):
47         img_file_X = self.list_files_X[index % len(self.
48         list_files_X)]
49         img_path_X = os.path.join(self.root_dir_X, img_file_X)
50         image_X = Image.open(img_path_X).convert('RGB')
51
52         img_file_Y = self.list_files_Y[index % len(self.
53         list_files_Y)]
54         img_path_Y = os.path.join(self.root_dir_Y, img_file_Y)
55         image_Y = Image.open(img_path_Y).convert('RGB')
56
57         image_X = constants.BASE_TRANSFORMATIONS(image_X)
58         image_Y = constants.BASE_TRANSFORMATIONS(image_Y)
59
60         return image_X, image_Y

```

---

### A.3 networks.py

```
1 import torch
2 import torch.nn as nn
3
4 class ConvBlock(nn.Module):
5     def __init__(self, in_c, out_c, kind, act, drop=None, **kwargs):
6         :
7             super().__init__()
8             Conv = nn.Conv2d if kind=='down' else nn.ConvTranspose2d
9
10            self.block = [Conv(in_c, out_c, **kwargs)]
11            self.block.append(
12                nn.InstanceNorm2d(out_c))
13
14            if act == 'relu':
15                self.block.append(nn.ReLU(True))
16            elif act == 'leaky':
17                self.block.append(nn.LeakyReLU(0.2, True))
18            elif act != 'none':
19                raise Exception(f'Invalid Activation')
20
21            if drop is not None:
22                self.block.append(nn.Dropout(drop))
23
24            self.block = nn.Sequential(*self.block)
25
26    def forward(self, X):
27        return self.block(X)
28
29 class EncoderBlock(nn.Module):
30     def __init__(self, in_channels, out_channels, stride=2):
31         super().__init__()
32         self.block = ConvBlock(
33             in_channels, out_channels, kind='down', act='leaky',
34             kernel_size=4, stride=stride, padding=1, bias=False,
35             padding_mode='reflect')
36
37     def forward(self, X):
38         return self.block(X)
39
40 class DecoderBlock(nn.Module):
41     def __init__(self, in_channels, out_channels, stride=2,
42                  dropout_rate=None):
43         super().__init__()
44         self.block = ConvBlock(
45             in_channels, out_channels, kind='up', act='relu', drop
46             =0.5,
47             kernel_size=4, stride=stride, padding=1, bias=False)
48
49     def forward(self, X):
50         return self.block(X)
51
52 class PatchDiscriminator(nn.Module):
53     def __init__(self, net_in_channels, conditional=True, filters
54                  =[64, 128, 256, 512]):
55         super().__init__()
```

```

51     # Notes:
52     # - For conditional, We get 2 images stacked on last axis
53     #   so first in_channels = in_channels*2
54
55     self.conditional = conditional
56     net_in_channels = net_in_channels*2 if conditional else
57     net_in_channels
58     initial = nn.Sequential(
59         nn.Conv2d(net_in_channels, filters[0], 4, 2, padding=1,
60         padding_mode='reflect'),
61         nn.LeakyReLU(0.2, inplace=True)
62     )
63
64     layers = nn.ModuleList([initial])
65     in_channels = filters[0]
66
67     for out_channels in filters[1:]:
68         layers.append(
69             EncoderBlock(in_channels, out_channels, stride=(1
70             if out_channels == filters[-1] else 2))
71         )
72         in_channels = out_channels
73
74     layers.append(
75         nn.Conv2d(
76             in_channels, out_channels=1, kernel_size=4,
77             stride=1, padding=1, padding_mode='reflect'
78         )
79     )
80
81     layers.append(nn.Sigmoid())
82
83     self.model = nn.Sequential(*layers)
84
85     def forward(self, X, Y=None):
86         if self.conditional:
87             XY = torch.cat([X, Y], dim = 1)
88             return self.model(XY)
89         else:
90             return self.model(X)
91
92     class UNetGenerator(nn.Module):
93         def __init__(_
94             self,
95             net_in_channels,
96             net_out_channels,
97             encoder_filters = [64, 128, 256, 512, 512, 512, 512,
98             512],
99             decoder_filters = [512, 512, 512, 512, 256, 128, 64]):
100             super().__init__()
101
102             # Encoder
103             self.encoder_layers = nn.ModuleList([nn.Sequential(
104                 nn.Conv2d(net_in_channels, encoder_filters[0], 4, 2, 1,
105                 padding_mode="reflect"),
106                 nn.LeakyReLU(0.2),
107             )])

```

```

103
104     for i in range(len(encoder_filters)-2):
105         self.encoder_layers.append(
106             EncoderBlock(encoder_filters[i], encoder_filters[i
107             +1])
108         )
109
110         # Note: See Errata
111         self.encoder_layers.append(nn.Sequential(
112             nn.Conv2d(encoder_filters[-2], encoder_filters[-1], 4,
113             2, 1),
114             nn.ReLU()
115         ))
116
117         # Decoder
118         self.decoder_layers = nn.ModuleList([DecoderBlock(
119             encoder_filters[-1], decoder_filters[0])])
120
121         for i in range(len(decoder_filters)-1):
122             self.decoder_layers.append(DecoderBlock(2*
123             decoder_filters[i], decoder_filters[i+1]))
124
125             self.decoder_layers.append(nn.Sequential(
126                 nn.ConvTranspose2d(2*decoder_filters[-1],
127                 net_out_channels, 4, 2, 1),
128                 nn.Tanh()
129             ))
130
131     def forward(self, X):
132         enc_activations = [X]
133         for enc in self.encoder_layers:
134             enc_activations.append(enc(enc_activations[-1]))
135
136         activation = self.decoder_layers[0](enc_activations[-1])
137
138         for i in range(1, len(self.decoder_layers)):
139             activation = torch.cat([activation, enc_activations[-(i
140             +1)]], dim = 1)
141             activation = self.decoder_layers[i](activation)
142
143         return activation
144
145 class ResNetBlock(nn.Module):
146     def __init__(self, channels):
147         super().__init__()
148         self.block = nn.Sequential(
149             ConvBlock(channels, channels, kind='down', act='relu',
150             kernel_size=3, padding=1),
151             ConvBlock(channels, channels, kind='down', act='none',
152             kernel_size=3, padding=1),
153         )
154
155     def forward(self, x):
156         return x + self.block(x)
157
158 class ResNetGenerator(nn.Module):

```

---

```

151     def __init__(self, net_in_channels, net_out_channels, filters =
152         [64, 128, 256], num_residuals=9):
153         super().__init__()
154
154         self.encoder_layers = nn.ModuleList()
155         self.encoder_layers.append(nn.Sequential(
156             nn.Conv2d(net_in_channels, filters[0], kernel_size=7,
157             stride=1, padding=3, padding_mode="reflect"),
158             nn.InstanceNorm2d(filters[0]),
159             nn.ReLU(inplace=True),
160         ))
161
161         for i in range(len(filters)-1):
162             self.encoder_layers.append(ConvBlock(
163                 filters[i], filters[i+1], kind='down', act='relu',
164                 kernel_size=3, stride=2, padding=1))
165
165         self.residual_layers = nn.Sequential(*[ResNetBlock(filters
166 [-1]) for _ in range(num_residuals)])
167
167         self.decoder_layers = nn.ModuleList()
168         for i in range(len(filters)-1):
169             self.decoder_layers.append(ConvBlock(
170                 filters[-(i+1)], filters[-(i+2)], kind='up', act=
171                 'relu', kernel_size=3, stride=2, padding=1, output_padding=1))
172
172         self.decoder_layers.append(nn.Conv2d(
173             filters[0], net_out_channels, kernel_size=7, stride=1,
174             padding=3, padding_mode="reflect"))
175
175         self.decoder_layers.append(nn.Tanh())
176
177     def forward(self, X):
178         for layer in self.encoder_layers:
179             X = layer(X)
180         for layer in self.residual_layers:
181             X = layer(X)
182         for layer in self.decoder_layers:
183             X = layer(X)
184
185         return X

```

## A.4 models.py

```

1 import os
2 from networks import ResNetGenerator
3
4 from pix2pix import Pix2Pix
5 from cyclegan import CycleGAN
6 from datasets import SideBySideDataset, UnpairedDataset
7
8 def satellite_to_map_model():
9     folder = os.path.abspath('data/satellite_to_map')
10    train_dataset = SideBySideDataset(os.path.join(folder, 'train')
11        , output='right')

```

```

11     test_dataset = SideBySideDataset(os.path.join(folder, 'val'),
12                                     output='right')
13
14     model = Pix2Pix(
15         train_dataset,
16         test_dataset,
17         os.path.join(folder, 'model.pth.tar'),
18         os.path.join(folder, 'example_outputs'),
19         num_examples=3
20     )
21
22     # model.train()
23     model.save_examples('final', 200)
24
25 def segmentation_to_facade_model():
26     folder = os.path.abspath('data/segmentation_to_facade')
27     train_dataset = SideBySideDataset(os.path.join(folder, 'train'),
28                                     output='left')
29     test_dataset = SideBySideDataset(os.path.join(folder, 'test'),
30                                     output='left')
31
32     model = Pix2Pix(
33         train_dataset,
34         test_dataset,
35         os.path.join(folder, 'model.pth.tar'),
36         os.path.join(folder, 'example_outputs'),
37         num_examples=3,
38     )
39
40     # model.train()
41     model.save_examples('final', 200)
42
43 def horse_zebra():
44     folder = os.path.abspath('data/horse_zebra')
45     train_dataset = UnpairedDataset(
46         root_dir_X = os.path.join(folder, 'trainA'),
47         root_dir_Y = os.path.join(folder, 'trainB')
48     )
49     test_dataset = UnpairedDataset(
50         root_dir_X = os.path.join(folder, 'testA'),
51         root_dir_Y = os.path.join(folder, 'testB')
52     )
53
54     model = CycleGAN(
55         train_dataset,
56         test_dataset,
57         os.path.join(folder, 'model.pth.tar'),
58         os.path.join(folder, 'example_outputs'),
59         num_examples=3,
60     )
61
62     # model.train()
63
64     model.save_examples('final', 200)
65
66 def apple_orange():
67     folder = os.path.abspath('data/apple_orange')

```

```

65     train_dataset = UnpairedDataset(
66         root_dir_X = os.path.join(folder, 'trainA'),
67         root_dir_Y = os.path.join(folder, 'trainB')
68     )
69     test_dataset = UnpairedDataset(
70         root_dir_X = os.path.join(folder, 'testA'),
71         root_dir_Y = os.path.join(folder, 'testB')
72     )
73
74     model = CycleGAN(
75         train_dataset,
76         test_dataset,
77         os.path.join(folder, 'model.pth.tar'),
78         os.path.join(folder, 'example_outputs'),
79         num_examples=3,
80     )
81
82     model.train()
83
84 def ukiyoe():
85     folder = os.path.abspath('data/ukiyoe2photo')
86     train_dataset = UnpairedDataset(
87         root_dir_X = os.path.join(folder, 'trainA'),
88         root_dir_Y = os.path.join(folder, 'trainB')
89     )
90     test_dataset = UnpairedDataset(
91         root_dir_X = os.path.join(folder, 'testA'),
92         root_dir_Y = os.path.join(folder, 'testB')
93     )
94
95     model = CycleGAN(
96         train_dataset,
97         test_dataset,
98         os.path.join(folder, 'model.pth.tar'),
99         os.path.join(folder, 'example_outputs'),
100        num_examples=3,
101    )
102
103     model.train()
104
105
106 def tester(p, out):
107     folder = os.path.abspath(p)
108     train_dataset = SideBySideDataset(os.path.join(folder, 'test'),
109         output=out)
110     test_dataset = SideBySideDataset(os.path.join(folder, 'test'),
111         output=out)
112
113     model = Pix2Pix(
114         train_dataset,
115         test_dataset,
116         os.path.join(folder, 'model.pth.tar'),
117         os.path.join(folder, 'example_outputs'),
118         num_examples=3,
119     )
120
121     # model.train()

```

---

```

120     model.save_examples('final', 200)
121
122 if __name__ == '__main__':
123     # satellite_to_map_model()
124     # segmentation_to_facade_model()
125     # horse_zebra()
126     # ukiyoe()
127     tester('data/shoes', 'right')
128     # apple_orange()

```

## A.5 pix2pix.py

```

1 import constants
2 from networks import ResNetGenerator, UNetGenerator,
3         PatchDiscriminator
4
5 import os
6
7 import torch
8 import torch.nn as nn
9 import torch.optim as optim
10 from torch.utils.data import DataLoader
11 from tqdm import tqdm
12 from torchvision.utils import save_image
13
14 class Pix2Pix:
15     def __init__(self,
16                  train_dataset,
17                  test_dataset,
18                  checkpoint_filename,
19                  examples_folder,
20                  discriminator = PatchDiscriminator,
21                  generator = ResNetGenerator,
22                  loss_type = 'mse',
23                  num_examples = 1
24      ):
25         self.train_dataset = train_dataset
26         self.test_dataset = test_dataset
27         self.checkpoint_filename = checkpoint_filename
28         self.examples_folder = examples_folder
29         self.num_examples = num_examples
30
31         self.train_loader = DataLoader(
32             self.train_dataset,
33             batch_size = constants.BATCH_SIZE,
34             shuffle = True,
35             num_workers = constants.NUM_WORKERS,
36         )
37         self.test_loader = DataLoader(test_dataset, batch_size=1,
38                                     shuffle=False)
39
40         self.D = discriminator(net_in_channels=3).to(constants.
41                               DEVICE)
42         self.G = generator(net_in_channels=3, net_out_channels=3).
43                     to(constants.DEVICE)

```

```

41
42     self.opt_D = optim.Adam(
43         self.D.parameters(), lr=constants.LEARNING_RATE, betas=
44         constants.BETAS)
45     self.opt_G = optim.Adam(
46         self.G.parameters(), lr=constants.LEARNING_RATE, betas=
47         constants.BETAS)
48
49     self.D_scaler = torch.cuda.amp.GradScaler()
50     self.G_scaler = torch.cuda.amp.GradScaler()
51
52     self.loss = nn.MSELoss() if loss_type == 'mse' else nn.
53     BCELoss()
54     self.l1_loss = nn.L1Loss()
55
56     if constants.LOAD_CHECKPOINT:
57         self.load_checkpoint()
58
59     def train(self):
60         for epoch in range(constants.NUM_EPOCHS):
61             print(f'Epoch {epoch}')
62             self.train_one_epoch()
63
64             if constants.SAVE_CHECKPOINT and epoch % 5 == 0:
65                 self.save_checkpoint()
66
67             self.save_examples(epoch)
68
69     def train_one_epoch(self):
70         pbar = tqdm(self.train_loader, leave=True)
71
72         for idx, (X, Y) in enumerate(pbar):
73             X = X.to(constants.DEVICE)
74             Y = Y.to(constants.DEVICE)
75
76             with torch.cuda.amp.autocast():
77                 D_real = self.D(X, Y)
78                 D_real_loss = self.loss(D_real, torch.ones_like(
79                     D_real))
80
81                 Y_fake = self.G(X)
82                 D_fake = self.D(X, Y_fake.detach())
83                 D_fake_loss = self.loss(D_fake, torch.zeros_like(
84                     D_fake))
85
86                 D_loss = (D_real_loss + D_fake_loss) / 2
87
88                 self.opt_D.zero_grad()
89                 self.D_scaler.scale(D_loss).backward()
90                 self.D_scaler.step(self.opt_D)
91                 self.D_scaler.update()
92
93                 with torch.cuda.amp.autocast():
94                     D_fake = self.D(X, Y_fake)
95                     G_fake_loss = self.loss(D_fake, torch.ones_like(
96                         D_fake))
97
98                     L1 = self.l1_loss(Y_fake, Y) * constants.L1_LAMBDA

```

```

92         G_loss = G_fake_loss + L1
93
94         self.opt_G.zero_grad()
95         self.G_scaler.scale(G_loss).backward()
96         self.G_scaler.step(self.opt_G)
97         self.G_scaler.update()
98
99         if idx % 10 == 0:
100             pbar.set_postfix(
101                 D_real=torch.sigmoid(D_real).mean().item(),
102                 D_fake=torch.sigmoid(D_fake).mean().item(),
103             )
104
105     def save_examples(self, epoch, n=None):
106         folder = self.examples_folder
107         i = iter(self.test_loader)
108         for n in range(self.num_examples if n is None else n):
109             X, Y = next(i)
110             X, Y = X.to(constants.DEVICE), Y.to(constants.DEVICE)
111
112             self.G.eval()
113             with torch.no_grad():
114                 Z = self.G(X)
115                 save_image(Z*0.5 + 0.5, os.path.join(folder, f"{epoch}_{n}_Y_fake.png"))
116                 save_image(X*0.5 + 0.5, os.path.join(folder, f"{epoch}_{n}_X.png"))
117                 save_image(Y*0.5 + 0.5, os.path.join(folder, f"{epoch}_{n}_Y.png"))
118             self.G.train()
119
120
121     def save_checkpoint(self):
122         filename = self.checkpoint_filename
123         checkpoint = {
124             "G_state_dict": self.G.state_dict(),
125             "opt_G_state_dict": self.opt_G.state_dict(),
126
127             "D_state_dict": self.D.state_dict(),
128             "opt_D_state_dict": self.opt_D.state_dict(),
129         }
130         torch.save(checkpoint, filename)
131         print(">> Checkpoint Saved")
132
133
134     def load_checkpoint(self):
135         checkpoint = torch.load(self.checkpoint_filename,
136                                 map_location=constants.DEVICE)
137         self.G.load_state_dict(checkpoint['G_state_dict'])
138         self.opt_G.load_state_dict(checkpoint['opt_G_state_dict'])
139
140         self.D.load_state_dict(checkpoint['D_state_dict'])
141         self.opt_D.load_state_dict(checkpoint['opt_D_state_dict'])
142
143         for optimizer in (self.opt_G, self.opt_D):
144             for param_group in optimizer.param_groups:
145                 param_group["lr"] = constants.LEARNING_RATE

```

---

```
145     print('>> Checkpoint Loaded')
```

## A.6 cyclegan.py

```
1 import constants
2 from networks import ResNetGenerator, PatchDiscriminator
3
4 import os
5
6 import torch
7 import torch.nn as nn
8 import torch.optim as optim
9 from torch.utils.data import DataLoader
10 from tqdm import tqdm
11 from torchvision.utils import save_image
12
13 class CycleGAN:
14     def __init__(self,
15                  train_dataset,
16                  test_dataset,
17                  checkpoint_filename,
18                  examples_folder,
19                  discriminator = PatchDiscriminator,
20                  generator = ResNetGenerator,
21                  loss_type = 'mse',
22                  num_examples = 1):
23
24         self.train_dataset = train_dataset
25         self.test_dataset = test_dataset
26         self.checkpoint_filename = checkpoint_filename
27         self.examples_folder = examples_folder
28         self.num_examples = num_examples
29
30         self.train_loader = DataLoader(
31             self.train_dataset,
32             batch_size = constants.BATCH_SIZE,
33             shuffle = True,
34             num_workers = constants.NUM_WORKERS)
35
36         self.test_loader = DataLoader(test_dataset, batch_size=1,
37                                     shuffle=False)
38
39         # Discriminator of the domain X
40         self.D_X = discriminator(net_in_channels=3, conditional=
41 False).to(constants.DEVICE)
42
43         # Discriminator of the domain Y
44         self.D_Y = discriminator(net_in_channels=3, conditional=
45 False).to(constants.DEVICE)
46
47         # Generator from Y -> X
48         self.G_YX = generator(net_in_channels=3, net_out_channels
=3).to(constants.DEVICE)
```

```

49     # Generator from X -> Y
50     self.G_XY = generator(net_in_channels=3, net_out_channels
51     =3).to(constants.DEVICE)
52
53     self.opt_D = optim.Adam(
54         list(self.D_X.parameters()) + list(self.D_Y.parameters
55         ()),
56         lr=constants.LEARNING_RATE, betas=constants.BETAS)
57
58     self.opt_G = optim.Adam(
59         list(self.G_YX.parameters()) + list(self.G_XY.
60         parameters()),
61         lr=constants.LEARNING_RATE, betas=constants.BETAS)
62
63     # Float-16 Training
64     self.D_scaler = torch.cuda.amp.GradScaler()
65     self.G_scaler = torch.cuda.amp.GradScaler()
66
67     self.loss = nn.MSELoss() if loss_type == 'mse' else nn.
68     BCELoss()
69     self.l1_loss = nn.L1Loss()
70
71     if constants.LOAD_CHECKPOINT:
72         self.load_checkpoint()
73
74     def train(self):
75         for epoch in range(constants.NUM_EPOCHS):
76             print(f'Epoch {epoch}')
77             self.train_one_epoch(epoch)
78
79             if constants.SAVE_CHECKPOINT and epoch % 5 == 0:
80                 self.save_checkpoint()
81
82             self.save_examples(epoch)
83
84     def train_one_epoch(self, epoch):
85         pbar = tqdm(self.train_loader, leave=True)
86
87         for idx, (image_X, image_Y) in enumerate(pbar):
88             image_X = image_X.to(constants.DEVICE)
89             image_Y = image_Y.to(constants.DEVICE)
90
91             # Train Discriminators
92             with torch.cuda.amp.autocast():
93                 # D_X
94                 fake_X = self.G_YX(image_Y)
95
96                 D_X_real = self.D_X(image_X)
97                 D_X_fake = self.D_X(fake_X.detach())
98
99                 D_X_real_loss = self.loss(D_X_real, torch.ones_like
(D_X_real))
D_X_fake_loss = self.loss(D_X_fake, torch.
zeros_like(D_X_fake))
100
D_X_loss = D_X_real_loss + D_X_fake_loss

```

```

100         # D_Y
101         fake_Y = self.G_XY(image_X)
102
103         D_Y_real = self.D_Y(image_Y)
104         D_Y_fake = self.D_Y(fake_Y.detach())
105
106         D_Y_real_loss = self.loss(D_Y_real, torch.ones_like
107             (D_Y_real))
108         D_Y_fake_loss = self.loss(D_Y_fake, torch.
109             zeros_like(D_Y_fake))
110
111         D_Y_loss = D_Y_real_loss + D_Y_fake_loss
112
113         # Total Loss
114         D_loss = (D_X_loss + D_Y_loss)/2
115
116         self.opt_D.zero_grad()
117         self.D_scaler.scale(D_loss).backward()
118         self.D_scaler.step(self.opt_D)
119         self.D_scaler.update()
120
121         # Train Generators
122         with torch.cuda.amp.autocast():
123             # Adversarial Loss
124             D_X_fake = self.D_X(fake_X)
125             D_Y_fake = self.D_Y(fake_Y)
126
127             G_YX_loss = self.loss(D_X_fake, torch.ones_like(
128                 D_X_fake))
129             G_XY_loss = self.loss(D_Y_fake, torch.ones_like(
130                 D_Y_fake))
131
132             # Cycle Loss
133             cycle_X = self.G_YX(fake_Y)
134             cycle_Y = self.G_XY(fake_X)
135
136             cycle_X_loss = self.l1_loss(image_X, cycle_X)
137             cycle_Y_loss = self.l1_loss(image_Y, cycle_Y)
138
139             G_loss = (
140                 G_XY_loss + G_YX_loss
141                 + constants.L1_LAMBDA_CYCLE*(cycle_X_loss +
142                     cycle_Y_loss)
143             )
144
145             # Identity Loss
146             if constants.L1_LAMBDA_IDENTITY > 0:
147                 identitiy_loss_X = self.l1_loss(image_X, self.
148                     G_YX(image_X))
149                 identitiy_loss_Y = self.l1_loss(image_Y, self.
150                     G_XY(image_Y))
151                 G_loss += constants.L1_LAMBDA_IDENTITY*(
152                     identitiy_loss_X+identitiy_loss_Y)
153
154                 self.opt_G.zero_grad()
155                 self.G_scaler.scale(G_loss).backward()
156                 self.G_scaler.step(self.opt_G)

```

```

149         self.G_scaler.update()
150
151     if idx % 200 == 0:
152         self.save_examples(f'{epoch}_{idx}')
153
154     if idx % 500 == 0 and idx > 0:
155         self.save_checkpoint()
156
157     if idx % 10 == 0:
158         pbar.set_postfix(
159             D_X_real=torch.sigmoid(D_X_real).mean().item(),
160             D_Y_real=torch.sigmoid(D_Y_real).mean().item(),
161             D_X_fake=torch.sigmoid(D_X_fake).mean().item(),
162             D_Y_fake=torch.sigmoid(D_Y_fake).mean().item(),
163         )
164
165     def save_examples(self, epoch, n = None):
166         folder = self.examples_folder
167         i = iter(self.test_loader)
168         for n in range(self.num_examples if n is None else n):
169             X, Y = next(i)
170             X, Y = X.to(constants.DEVICE), Y.to(constants.DEVICE)
171
172             self.G_XY.eval()
173             self.G_YX.eval()
174             with torch.no_grad():
175                 fake_Y = self.G_XY(X)
176                 fake_X = self.G_YX(Y)
177                 save_image(Y*0.5 + 0.5, os.path.join(folder, f'{epoch}_{n}_Y.png'))
178                 save_image(X*0.5 + 0.5, os.path.join(folder, f'{epoch}_{n}_X.png'))
179                 save_image(fake_Y*0.5 + 0.5, os.path.join(folder, f'{epoch}_{n}_Y_fake.png'))
180                 save_image(fake_X*0.5 + 0.5, os.path.join(folder, f'{epoch}_{n}_X_fake.png'))
181
182             self.G_XY.train()
183             self.G_YX.train()
184
185     def save_checkpoint(self):
186         filename = self.checkpoint_filename
187         checkpoint = {
188             "D_X_state_dict": self.D_X.state_dict(),
189             "D_Y_state_dict": self.D_Y.state_dict(),
190             "opt_D_state_dict": self.opt_D.state_dict(),
191
192             "G_XY_state_dict": self.G_XY.state_dict(),
193             "G_YX_state_dict": self.G_YX.state_dict(),
194             "opt_G_state_dict": self.opt_G.state_dict(),
195         }
196         torch.save(checkpoint, filename)
197         print(">> Checkpoint Saved")
198
199     def load_checkpoint(self):

```

---

```

201     checkpoint = torch.load(self.checkpoint_filename,
202         map_location=constants.DEVICE)
203     self.G_XY.load_state_dict(checkpoint['G_XY_state_dict'])
204     self.G_YX.load_state_dict(checkpoint['G_YX_state_dict'])
205     self.opt_G.load_state_dict(checkpoint['opt_G_state_dict'])
206
207     self.D_X.load_state_dict(checkpoint['D_X_state_dict'])
208     self.D_Y.load_state_dict(checkpoint['D_Y_state_dict'])
209     self.opt_D.load_state_dict(checkpoint['opt_D_state_dict'])
210
211     for optimizer in (self.opt_G, self.opt_D):
212         for param_group in optimizer.param_groups:
213             param_group["lr"] = constants.LEARNING_RATE
214     print('>> Checkpoint Loaded')

```

## A.7 models.py

```

1 import os
2 from networks import ResNetGenerator
3
4 from pix2pix import Pix2Pix
5 from cyclegan import CycleGAN
6 from datasets import SideBySideDataset, UnpairedDataset
7
8 def satellite_to_map_model():
9     folder = os.path.abspath('data/satellite_to_map')
10    train_dataset = SideBySideDataset(os.path.join(folder, 'train'),
11        , output='right')
11    test_dataset = SideBySideDataset(os.path.join(folder, 'val'),
12        , output='right')
12
13    model = Pix2Pix(
14        train_dataset,
15        test_dataset,
16        os.path.join(folder, 'model.pth.tar'),
17        os.path.join(folder, 'example_outputs'),
18        num_examples=3
19    )
20
21    # model.train()
22    model.save_examples('final', 200)
23
24 def segmentation_to_facade_model():
25    folder = os.path.abspath('data/segmentation_to_facade')
26    train_dataset = SideBySideDataset(os.path.join(folder, 'train'),
27        , output='left')
27    test_dataset = SideBySideDataset(os.path.join(folder, 'test'),
28        , output='left')
28
29    model = Pix2Pix(
30        train_dataset,
31        test_dataset,
32        os.path.join(folder, 'model.pth.tar'),
33        os.path.join(folder, 'example_outputs'),
34        num_examples=3,
35    )

```

```

36
37     # model.train()
38     model.save_examples('final', 200)
39
40 def horse_zebra():
41     folder = os.path.abspath('data/horse_zebra')
42     train_dataset = UnpairedDataset(
43         root_dir_X = os.path.join(folder, 'trainA'),
44         root_dir_Y = os.path.join(folder, 'trainB')
45     )
46     test_dataset = UnpairedDataset(
47         root_dir_X = os.path.join(folder, 'testA'),
48         root_dir_Y = os.path.join(folder, 'testB')
49     )
50
51     model = CycleGAN(
52         train_dataset,
53         test_dataset,
54         os.path.join(folder, 'model.pth.tar'),
55         os.path.join(folder, 'example_outputs'),
56         num_examples=3,
57     )
58
59     # model.train()
60
61     model.save_examples('final', 200)
62
63 def apple_orange():
64     folder = os.path.abspath('data/apple_orange')
65     train_dataset = UnpairedDataset(
66         root_dir_X = os.path.join(folder, 'trainA'),
67         root_dir_Y = os.path.join(folder, 'trainB')
68     )
69     test_dataset = UnpairedDataset(
70         root_dir_X = os.path.join(folder, 'testA'),
71         root_dir_Y = os.path.join(folder, 'testB')
72     )
73
74     model = CycleGAN(
75         train_dataset,
76         test_dataset,
77         os.path.join(folder, 'model.pth.tar'),
78         os.path.join(folder, 'example_outputs'),
79         num_examples=3,
80     )
81
82     model.train()
83
84 def ukiyoe():
85     folder = os.path.abspath('data/ukiyoe2photo')
86     train_dataset = UnpairedDataset(
87         root_dir_X = os.path.join(folder, 'trainA'),
88         root_dir_Y = os.path.join(folder, 'trainB')
89     )
90     test_dataset = UnpairedDataset(
91         root_dir_X = os.path.join(folder, 'testA'),
92         root_dir_Y = os.path.join(folder, 'testB')

```

```
93
94
95     model = CycleGAN(
96         train_dataset,
97         test_dataset,
98         os.path.join(folder, 'model.pth.tar'),
99         os.path.join(folder, 'example_outputs'),
100        num_examples=3,
101    )
102
103    model.train()
104
105
106 def tester(p, out):
107     folder = os.path.abspath(p)
108     train_dataset = SideBySideDataset(os.path.join(folder, 'test'),
109         output=out)
110     test_dataset = SideBySideDataset(os.path.join(folder, 'test'),
111         output=out)
112
113     model = Pix2Pix(
114         train_dataset,
115         test_dataset,
116         os.path.join(folder, 'model.pth.tar'),
117         os.path.join(folder, 'example_outputs'),
118        num_examples=3,
119    )
120
121     # model.train()
122     model.save_exmaples('final', 200)
123
124 if __name__ == '__main__':
125     # satellite_to_map_model()
126     # segmentation_to_facade_model()
127     # horse_zebra()
128     # ukiyoe()
129     tester('data/shoes', 'right')
130     # apple_orange()
```

# Bibliography

- [1] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [3] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [4] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.
- [5] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.
- [6] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [7] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2017.
- [8] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [9] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [10] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015.

---

*ference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

- [11] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [12] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016.
- [13] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.