

深入理解风格迁移三部曲(一)--UNIT

标签（空格分隔）： 陈扬

MUNIT: Multimodal Unsupervised Image-to-Image Translation

github:<https://github.com/NVLabs/MUNIT>

作者:[陈扬](#)

深入理解风格迁移三部曲(一)--UNIT

[简介](#)

[前言](#)

[GAN](#)

[VAE](#)

[UNIT](#)

[损失函数](#)

[实验结果](#)

简介

近期我研究的方向转向了GAN的应用, 其中图像的风格迁移是GAN中一个非常有意思的应用,传统的方法基于拉普拉斯金字塔对成对的图像进行纹理上的风格迁移.随着2014年GAN的爆火,研究者发现GAN通过判别器D学习两个图像域的关系,实现了unpaired image-to-image(非成对图像数据集的风格迁移)的功能,其中有两个广为人知的应用分别是pix2pix和cycleGAN,今天我们另辟蹊径,从NVIDIA-Lab提出的UNIT框架来探索image-to-image的实现原理.

前言

在开始说UNIT之前,我们先来简要的回顾一下GAN和VAE,这也是在我之前的[ODOG](#)项目中有过详细的介绍.

GAN

生成对抗网络（英语：**G**enerative **A**dversarial **N**etwork，简称GAN）是[非监督式学习](#)的一种方法，通过让两个[神经网络](#)相互[博弈](#)的方式进行学习。该方法由[伊恩·古德费洛](#)等人于2014年提出。[1]

核心公式:

In other words, D and G play the following two-player minimax game with value function $V(G, D)$:

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

这个公式我们要分成两个部分来看: 先看前半部分:

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

这个公式的意思是, 先看加号前面 $\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$

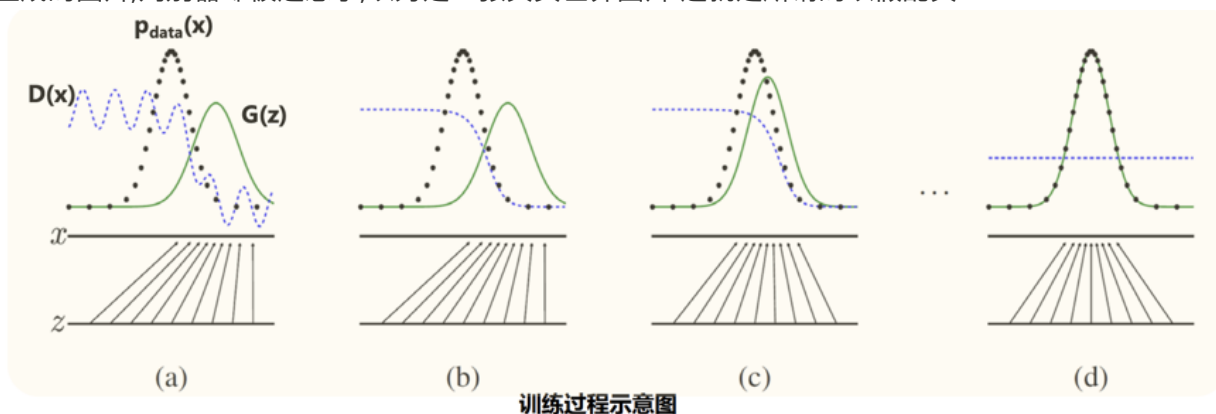
, 我们希望 D 最大, 所以 $\log(D(\mathbf{x}))$ 应该最大, 意味着我的判别器可以很好的识别出, 真实世界图像是 "true", 在看加号后面 $\mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$, 要让 \log 尽可能的大, 需要的是 $D(G(\mathbf{z}))$ 尽可能的小, 意味着我们生成模型的图片应该尽可能的被判别模型视为 "FALSE".

再看后半部分

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

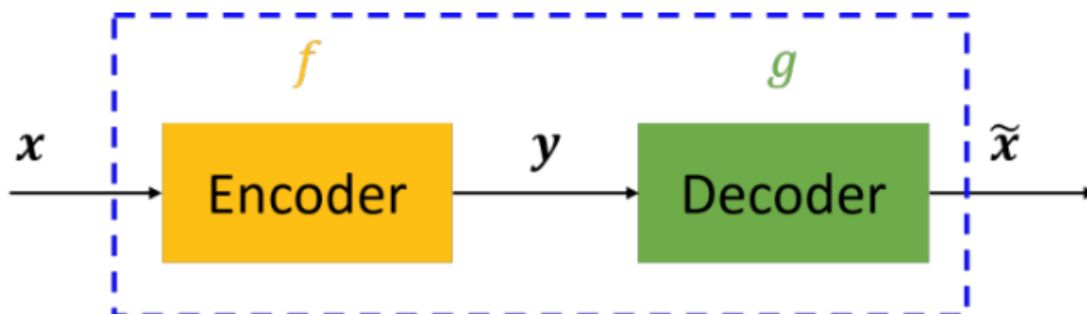
, 我们应该让 G 尽可能的小, 加号前面的式子并没有 G , 所以无关, 在看加号后面的式子

$\mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$, 要让 G 尽可能的小, 就要 $D(G(\mathbf{z}))$ 尽可能的大, 也就是说本来就一张 → 噪声生成的图片, 判别器却被迷惑了, 以为是一张真实世界图片. 这就是所谓的以假乱真.



VAE

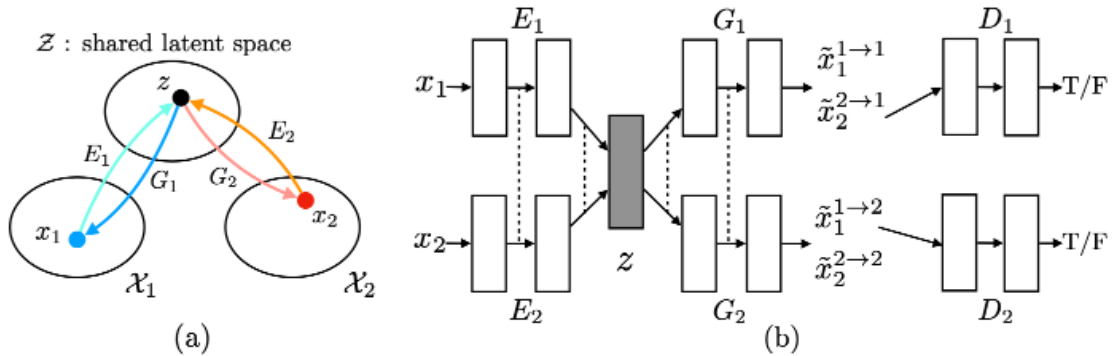
VAE 也叫变分自动编码器, 其基础是自编码器(autoencoder), AE通过对输入 X 进行编码后得到一个低维的向量 y , 然后根据这个向量还原出输入 X 。通过对比 X 与 \tilde{X} 的误差, 再利用神经网络去训练使得误差逐渐减小从而达到非监督学习的目的。



简单点来说,就是先通过 f 把 x 映射到 y ,一般来说我们假定 $y \sim N(0, 1)$,再通过 g 把 y 映射到 \tilde{X} ,在这里不长篇大论Reparametrization tricks.

UNIT

UNIT(Unsupervised Image-to-Image Translation Networks),由NVIDIA-Lab在NIPS 2017年提出,该文章首次提Image-Image Translation这个概念,将计算机视觉和计算机图形学的许多任务总结进去,分为一对多和多对一的两类转换任务,包括CV里的边缘检测,图像分割,语义标签以及CG里的mapping labels or sparse user inputs to realistic images.



该文章定义了 χ_1 和 χ_2 作为两个图像域.传统的supervised Image-to-image 通过对图像域进行采样,求其联合概率分布 $P_{(\chi_1, \chi_2)}(x_1, x_2)$,通过 Encoder-Decoder 的思想,作者定义了两个E和G,希望使得 $z = E(X)$ 在 latent space 上近可能的分布一致.意味着当我们同时对 $Sample(\chi_1, \chi_2)$ 时,我们希望得出:

$$z = E_1^*(x_1) = E_2^*(x_2) \quad (1)$$

这样,我们得到了两个Domain下image的一致表示,再通过令 $G = D$,从latent space中重构 $\hat{x} = G(z)$,

因此,我们两个采样下的 $\{x_1, x_2\}$ 经过 $\{< E_1, G_1 >, < E_2, G_1 >, < E_1, G_2 >, < E_2, G_1 >\}$ 后得到了 $\{\hat{x}_1^{1 \rightarrow 1}, \hat{x}_2^{2 \rightarrow 1}, \hat{x}_1^{1 \rightarrow 2}, \hat{x}_2^{2 \rightarrow 2}\}$,再把:

$$\hat{x}_1^{1 \rightarrow 1}, \hat{x}_2^{2 \rightarrow 1} \rightarrow D_1 \rightarrow T/F \quad (2)$$

$$\hat{x}_1^{1 \rightarrow 2}, \hat{x}_2^{2 \rightarrow 2} \rightarrow D_2 \rightarrow T/F \quad (3)$$

通过Adv_loss对抗学习跨域生成图片的效果.

可能细心的你以及发现了这是不是很类似VAE-GAN吗?是的.

损失函数

作者通过联合训练4个网络 $VAE_1, VAE_2, GAN_1, GAN_2$ 的三个loss function来训练整个网络:

$$\min_{E_1, E_2, G_1, G_2} \max_{D_1, D_2} \mathcal{L}_{VAE_1}(E_1, G_1) + \mathcal{L}_{GAN_1}(E_2, G_1, D_1) + \mathcal{L}_{CC_1}(E_1, G_1, E_2, G_2) \\ \mathcal{L}_{VAE_2}(E_2, G_2) + \mathcal{L}_{GAN_2}(E_1, G_2, D_2) + \mathcal{L}_{CC_2}(E_2, G_2, E_1, G_1)$$

```
1 # -----
2 # Train Encoders and Generators
3 # -----
4 # Total loss
```

```

5  loss_G = (
6      loss_KL_1
7      + loss_KL_2
8      + loss_ID_1
9      + loss_ID_2
10     + loss_GAN_1
11     + loss_GAN_2
12     + loss_KL_1_
13     + loss_KL_2_
14     + loss_cyc_1
15     + loss_cyc_2
16 )
17 loss_D1 = criterion_GAN(D1(X1), valid) +
criterion_GAN(D1(fake_X1.detach()), fake)
18 loss_D2 = criterion_GAN(D2(X2), valid) +
criterion_GAN(D2(fake_X2.detach()), fake)

```

VAE的目标是minimize source domain to latent space's KL diversity and latent space to destination domain's KL diversity(我觉得中文太拗口了,这句话实在是说不来)来最小化变分上界,VAE的定义如下:

$$\begin{aligned}
 \mathcal{L}_{\text{VAE}_1}(E_1, G_1) &= \lambda_1 \text{KL}(q_1(z_1|x_1) \| p_\eta(z)) - \lambda_2 \mathbb{E}_{z_1 \sim q_1(z_1|x_1)} [\log p_{G_1}(x_1|z_1)] \\
 \mathcal{L}_{\text{VAE}_2}(E_2, G_2) &= \lambda_1 \text{KL}(q_2(z_2|x_2) \| p_\eta(z)) - \lambda_2 \mathbb{E}_{z_2 \sim q_2(z_2|x_2)} [\log p_{G_2}(x_2|z_2)]
 \end{aligned} \tag{4}$$

```

1  # Get shared latent representation
2  mu1, z1 = E1(X1)
3  mu2, z2 = E2(X2)
4  # Reconstruct images
5  recon_X1 = G1(z1)
6  recon_X2 = G2(z2)
7  # Translate images
8  fake_X1 = G1(z2)
9  fake_X2 = G2(z1)
10 loss_KL_1 = lambda_1 * compute_kl(mu1)
11 loss_KL_2 = lambda_1 * compute_kl(mu2)
12 loss_KL_1_ = lambda_3 * compute_kl(mu1_)
13 loss_KL_2_ = lambda_3 * compute_kl(mu2_)

```

对抗:GAN_LOSS被用于确保翻译图像类似图像在目标域.定义如下:

$$\begin{aligned}
 \mathcal{L}_{\text{GAN}_1}(E_2, G_1, D_1) &= \lambda_0 \mathbb{E}_{x_1 \sim P_{X_1}} [\log D_1(x_1)] + \lambda_0 \mathbb{E}_{z_2 \sim q_2(z_2|x_2)} [\log(1 - D_1(G_1(z_2)))] \\
 \mathcal{L}_{\text{GAN}_2}(E_1, G_2, D_2) &= \lambda_0 \mathbb{E}_{x_2 \sim P_{X_2}} [\log D_2(x_2)] + \lambda_0 \mathbb{E}_{z_1 \sim q_1(z_1|x_1)} [\log(1 - D_2(G_2(z_1)))]
 \end{aligned} \tag{5}$$

```

1 loss_GAN_1 = lambda_0 * criterion_GAN(D1(fake_X1), valid)
2 loss_GAN_2 = lambda_0 * criterion_GAN(D2(fake_X2), valid)
3 loss_D1 = criterion_GAN(D1(X1), valid) +
  criterion_GAN(D1(fake_X1.detach()), fake)
4 loss_D2 = criterion_GAN(D2(X2), valid) +
  criterion_GAN(D2(fake_X2.detach()), fake)

```

循环一致性:由于shared latent-space假设暗含了循环一致性约束, 因此我们在提出的框架中实施循环一致性约束, 以进一步规范不适定的无监督图像间转换问题。产生的信息处理流称为循环重建流, 定义如下:

$$\begin{aligned} \mathcal{L}_{CC_1}(E_1, G_1, E_2, G_2) &= \lambda_3 \text{KL}(q_1(z_1|x_1) \| p_\eta(z)) + \lambda_3 \text{KL}(q_2(x_1^{1 \rightarrow 2}) \| p_\eta(z)) - \\ &\quad \lambda_4 \mathbb{E}_{z_2 \sim q_2(z_2|x_1^{1 \rightarrow 2})} [\log p_{G_1}(x_1|z_2)] \\ \mathcal{L}_{CC_2}(E_2, G_2, E_1, G_1) &= \lambda_3 \text{KL}(q_2(z_2|x_2) \| p_\eta(z)) + \lambda_3 \text{KL}(q_1(x_2^{2 \rightarrow 1}) \| p_\eta(z)) - \\ &\quad \lambda_4 \mathbb{E}_{z_1 \sim q_1(z_1|x_2^{2 \rightarrow 1})} [\log p_{G_2}(x_2|z_1)] \end{aligned}$$

```

1 # Cycle translation
2 mu1_, z1_ = E1(fake_X1)
3 mu2_, z2_ = E2(fake_X2)
4 cycle_X1 = G1(z2_)
5 cycle_X2 = G2(z1_)
6 loss_ID_1 = lambda_2 * criterion_pixel(recon_X1, X1)
7 loss_ID_2 = lambda_2 * criterion_pixel(recon_X2, X2)
8 loss_cyc_1 = lambda_4 * criterion_pixel(cycle_X1, X1)
9 loss_cyc_2 = lambda_4 * criterion_pixel(cycle_X2, X2)

```

训练好的网络, 我们可以通过对 latent sapce 的 latent variable 重编码, 进而把输入图像迁移到各个域中:

实验结果



Figure 6: Attribute-based face translation results.

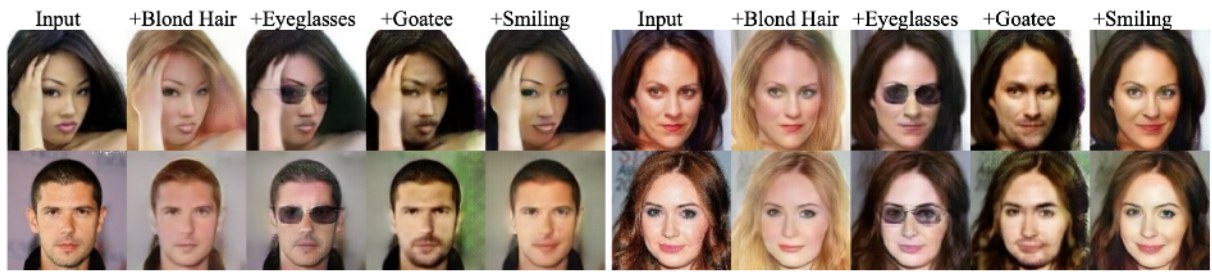


Figure 6: Attribute-based face translation results.

作者在展示的时候看起来好像可以实现一对多的风格转换,实际上这个算法只能实现 1对1的风格迁移,是作者做了 N 对1对1 的实验,所以看起来像 1对 N 的结果.

这算是比较早期的一篇文章,其实现原理也是借鉴很很多前人的工作,实际上我觉得从原创性上来看比不上cycleGAN,不过这这个VAE-GAN延伸的应用性似乎更好一些.接下来会介绍NVIDIA-Lab的FUNIT和MUNIT.