#### Disentangled Representation Learning for Controllable Image Synthesis: an Information-Theoretic Perspective

Shichang Tang ShanghaiTech University SIMIT, Chinese Academy of Sciences University of Chinese Academy of Sciences

Xu Zhou

Xuming He ShanghaiTech University Yi Ma University of California, Berkeley

Speaker: Shichang Tang

tangshch@shanghaitech.edu.cn

tsc2017@gmail.com





# Introduction

#### **A** Limitations of generative models

01

02

03

VAE and GAN have been developed for a few years, but we still do not know much about them

Intro

The latent space of GAN and VAE are usually uninterpretable, and thus uncontrolled

While GAN can synthesize better images, it does not have an encoder and cannot learn an latent representation given an input image



- ✓ Construct a flexible, interpretable generative model that can generate realistic images
- ✓ And be able to learn a disentangled representation so that the output images can be manipulated by controlling the latents

# Background

Intro Ba

2

## **M** Background



# ∧ Background: GAN



Disentangled

Conclusion

Related

Intro



# Disentangled Representation Learning for Controllable Image Synthesis

#### 3 Intro Related Disentangled Conclusion **A VAE with partitioned latent code** x $z_1$ Ex'feature extractor concatenate Z $E_2$ $z_2$ f(x)x (a)(b)

## **VAE** with partitioned latent code

Different uses of generative models:

#### **1.Feature learning:**

Auto-encoding: 
$$x \to (z_1, z_2) \to x$$

#### 2.Generate images from scratch:

Sampling: 
$$(\bar{z_1}, \bar{z_2}) \rightarrow \bar{x}$$

**3.Fusion (editing)**: Sample images independently:  $x \sim P_r$ ,  $\tilde{x} \sim P_r$ Encode:  $z_1 \sim Q_1(\cdot|x)$   $\tilde{z}_2 \sim Q_2(\cdot|\tilde{x})$ 

Fusion:

$$(z_1, \tilde{z}_2) \rightarrow \hat{x}$$



Disentangled

Conclusion

Intro

Related

## **VAE** with partitioned latent code

#### **Entanglement of information**

mutual information between  $z_2$  and x' can be shared by  $z_1$ :

$$I(z_2; x') = I(z_2; x'|z_1) + I(z_1; z_2; x')$$
  
e.g.,:

$$z_1 = E_1(x), \ z_2 = E_2(x), \ x' = E_1^{-1}(z_1)$$

mutual information between  $z_2$ and x' is maximized, but  $z_2$  is not utilized by G at all

Therefore, mutual information between z<sub>2</sub> and x' should be maxminzed, but not shared by z<sub>1</sub>!

Approach: maximize the conditional mutual inforamtion:

$$I(z_2; x'|z_1)$$



Related

Disentangled

Conclusion

Intro

### **A VAE with partitioned latent code**

Intro Related

Optimizable lower bounds: maximizing the conditional MI is to minimize the reconstruction loss when the latents are independent

Maximize:

 $I(\tilde{z}_2; \hat{x} | z_1) \ge \mathbb{E}[\log q(\tilde{z}_2 | z_1, \hat{x})] + H(\tilde{z}_2 | z_1).$ 

If  $\tilde{z}_2$  and  $z_1$  are independent, and  $H(\tilde{z}_2)$  is constant, then one can

minimize

 $L_{f1} = -\mathbb{E}[\log q(\tilde{z}_2 | z_1, \hat{x})] \quad (\text{Reconstruction loss of } \tilde{z}_2)$ 

If  $Q(\cdot|z_1, \hat{x})$  is a Benoulli distribution, it becomes the cross-entropy loss



Or maximize

 $I(\tilde{z}_2; \hat{x} | z_1) \geq \mathbb{E}[\log q^*(f(\tilde{x}) | z_1, \hat{x})] + H(f(x)) \in \text{Reconstruction loss of } f(\tilde{x}) > 0$ 

Under a Gaussian assumption, we have

$$L_{f2} = \mathbb{E}[\|f(G(z_1, \tilde{z}_2)) - f(\tilde{x})\|_2^2]$$

## **VAE** with partitioned latent code

Intro R

#### The use of CGAN

**Conditional mutual information maximization does not guarantee fidelity** 

e.g., samples overfitting the classifier can be unrealistic

D estimates the Jensen-Shannon Divergence (JS-Divergence) between  $P(\tilde{z}_2, \hat{x})$   $\Re P(z_2, x')$ 

$$L_D = -2JSD(P(z_2, x') || P(\tilde{z}_2, \hat{x})) + \log(4)$$
  
=  $\mathbb{E}[\log D(z_2, x') + \log(1 - D(\tilde{z}_2, \hat{x}))]$ 

While G tries to minimize:

$$L_{GAN} = -\mathbb{E}[\log(D(\tilde{z}_2, \hat{x}))]$$

# Intro Related Disentangled Conclusion Comparison methods: VAE VAE+Ls(InfoGAN)

VAE+GAN

VAE+L<sub>f2</sub>  $L_{f2} = \mathbb{E}[\|f(G(z_1, \tilde{z}_2)) - f(\tilde{x})\|_2^2]$ 

VAE+GAN+Lf2





## **Controlling** attributes of CelebA

#### **Comparison methods:**



Intro

Related

Disentangled

Conclusion

## **A** Flipping attributes of CelebA

#### **Qualitative comparison:**

Without the GAN loss, the output can be unrealistic; without the information loss, the output can be unchanged



Fig. 4. Flipping  $z_2$  when it represents "Smiling".



Fig. 5. Flipping  $z_2$  when it represents "Young".



3

#### **A** Controlling attributes of CelebA

**Interpolations and extrapolations** (when z<sub>2</sub> represents "young") :







(e)  $z_2 = 2$ 

(f)  $z_2 = 3$ 

### **A** Controlling attributes of CelebA

**Control multiple attributes at the same time:** 

Related

Intro



3 Disentan<u>gled</u>



Conclusion

(a) real images

(b) not smiling, not (c) not smiling, not young, female young, male



(d) not smiling, young, (e) not smiling, young, (f) smiling, not young, female male female



(g) smiling, not young, (h) smiling, young, fe- (i) smiling, young, male male

# Conclusions

Λ	
Conclusions	<ul> <li>We give a rigorous derivation of a variant of VAE with partitioned code, and analyze the phenomenon of entangled information from an information-theoretic perspective</li> <li>Derive optimizable lower bounds of conditional mutual information</li> <li>Use the proposed method to learn disentangled representation and perform controllable image synthesis</li> </ul>
Future Work	<ul> <li>Use some recent state-of-the-art GANs such as <u>StyleGAN2-ADA</u> or <u>MIX-GAN</u></li> <li>Resolve the difficulty in optimization when there are too many loss terms</li> <li>Experiment on some more attributions and ethnically diverse datasets</li> </ul>

For more theoretical and experimental results, please refer to the paper and the supplementary material!

# **Thanks!**