

CausalVAE: Disentangled Representation Learning via Neural Structural Causal Models

Mengyue Yang^{1,2}, Furui Liu^{1,*}, Zhitang Chen¹, Xinwei Shen³, Jianye Hao¹, Jun Wang² ¹
Noah's Ark Lab, Huawei, Shenzhen, China ² University College London, London, United
Kingdom ³ The Hong Kong University of Science and Technology, Hong Kong, China
CVPR 2021

Nikoo Naghavian
Winter 2023

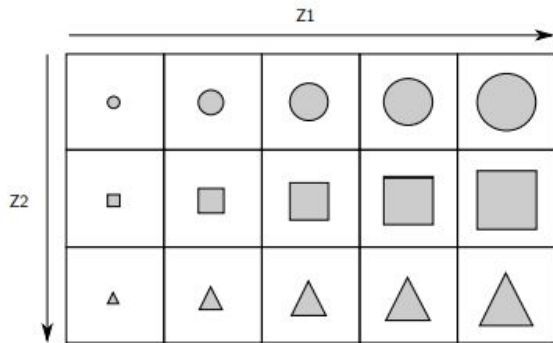
outline

- Introduction
- Model structure of causalVAE
- Learning strategy
- Experiments
- Methods comparison
- Conclusion
- Code explanation



Disentangled Representation

A disentangled representation can be defined as one where single latent units are sensitive to changes in single (independent) generative (ground-truth) factors, while being relatively invariant to changes in other factors





Variational Autoencoder (VAE)

- Setting:
 - Let \mathbf{z} denote the concepts representation, \mathbf{x} be observations
 - Using marginal distribution $p(\mathbf{x})$ to get joint distribution $p(\mathbf{x}, \mathbf{z})$
- Model:
 - Encoder: inference latent code \mathbf{z}
 - Decoder: generate/reconstruct \mathbf{x} from \mathbf{z}
- Objective:
 - Maximize Evidence lower bound
 - $\mathbb{E}_{p_{\theta}(\mathbf{z})}[p_{\theta}(\mathbf{x})] \geq ELBO = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x})] - D_{kl}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$
 - $\mathbb{E}_{p_{\theta}(\mathbf{z})}[p_{\theta}(\mathbf{x}|\mathbf{z})]$: Reconstruct \mathbf{x} from \mathbf{z} in generate process
 - $D_{kl}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}))$: Introduce posterior $q_{\phi}(\mathbf{z}|\mathbf{x})$ and prior $p_{\theta}(\mathbf{z})$ (Multivariate Gaussian) in inference process





Problems in traditional disentanglement works

- The concepts are causally related.
- Unsupervised process could not guarantee the learned representation is identifiable.

Light
Pendulum
Shadow



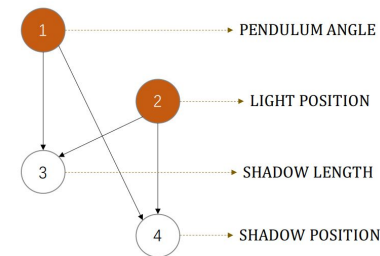
Swing Pendulum



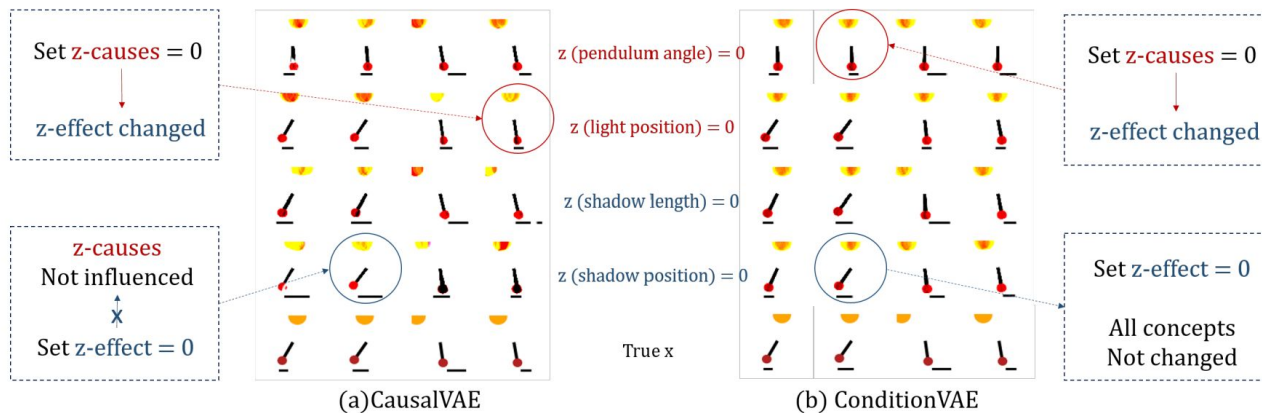
Causal disentanglement

- Learning disentanglement representation which align to real word concepts
- Learning causal graph automatically
- Achieve do-operation on casual representation
- Generate counterfactual images that do not appear in training data

The result of intervention

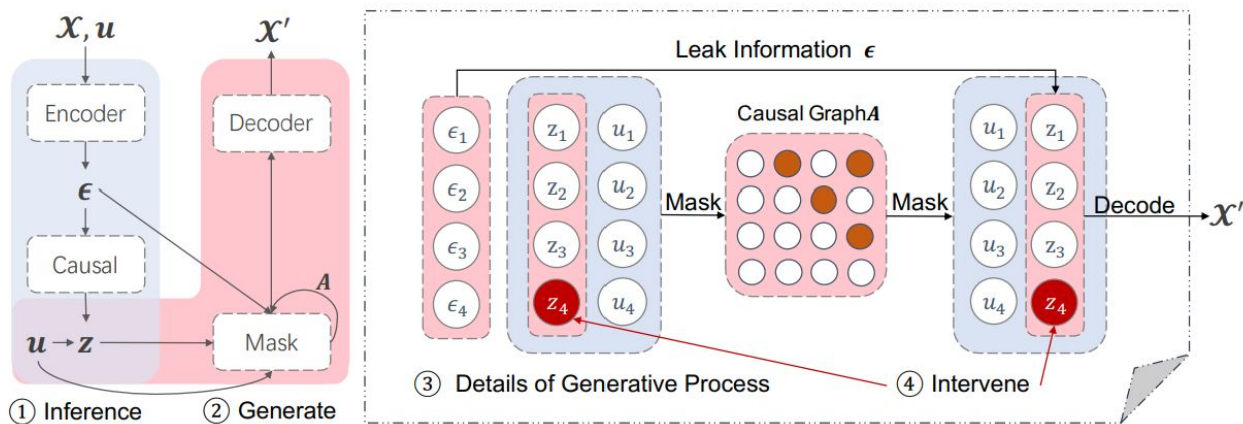


Causal graph



The result of intervention on pendulum dataset

Model structure of CausalVAE



$$\mathbf{z} = \mathbf{A}^T \mathbf{z} + \epsilon = (\mathbf{I} - \mathbf{A}^T)^{-1} \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

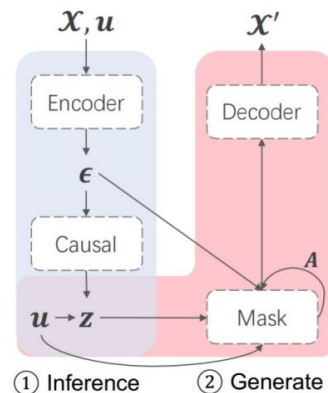
$$z_i = g_i(\mathbf{A}_i \circ \mathbf{z}; \boldsymbol{\eta}_i) + \epsilon_i,$$

Transforming Independent Exogenous Factors into Causal Representations

- Structural Causal Models (SCMs):

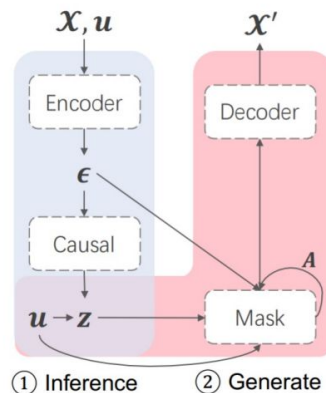
$$\mathbf{z} = \mathbf{A}^T \mathbf{z} + \boldsymbol{\epsilon} = (\mathbf{I} - \mathbf{A}^T)^{-1} \boldsymbol{\epsilon}$$

- ϵ_i 's are jointly independent which satisfy $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
- \mathbf{z}_i is the lower dimensional representation of i-th concept.
- \mathbf{A} is causal graph



VAE-based Model

- Inference Model $q_{\phi}(\mathbf{z}, \epsilon | \mathbf{x}, \mathbf{u})$
 - Encoder transform observations \mathbf{X} into ϵ
 - A Causal layer generate causal representation
 - $\mathbf{z} = (\mathbf{I} - \mathbf{A}^T)^{-1} \epsilon$
 - Introduce additional observation \mathbf{u} :
 - \mathbf{z} satisfy conditional Gaussian $\mathbf{z} \sim \mathcal{N}(\lambda_1(\mathbf{u}), \lambda_2^2(\mathbf{u}))$, where \mathbf{u} is additional observation.
- Generative Model
 - $p_{\theta}(\mathbf{x}, \mathbf{z}, \epsilon | \mathbf{u}) = p_{\theta}(\mathbf{x} | \mathbf{z}, \epsilon, \mathbf{u}) p_{\theta}(\epsilon, \mathbf{z} | \mathbf{u})$
 - Introduce a Mask Layer
 - $z_i = g_i(\mathbf{A}^T \mathbf{z}; \boldsymbol{\eta}_i) + \epsilon_i$, where $\boldsymbol{\eta}_i$ is parameter of g_i
 - Achieve do-operation



Learning Strategy

- Evidence Lower Bound (ELBO) of CausalVAE

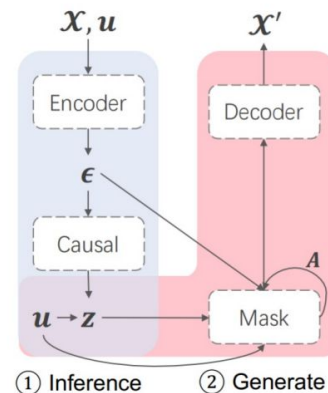
$$\mathbb{E}_{q_{\mathcal{X}}}[\log p_{\theta}(\mathbf{x}|\mathbf{u})] \geq \text{ELBO} = \mathbb{E}_{q_{\mathcal{X}}}[\mathbb{E}_{\epsilon, \mathbf{z} \sim q_{\phi}}[\log p_{\theta}(\mathbf{x}|\mathbf{z}, \epsilon, \mathbf{u})] - \frac{1}{\epsilon} \mathcal{D}(q_{\phi}(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) || p_{\theta}(\epsilon, \mathbf{z}|\mathbf{u}))]$$

- Probabilistic definition:
 - Inference model : $q_{\phi}(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) = q_{\phi}(\epsilon|\mathbf{x}, \mathbf{u})\delta(\mathbf{z} = \mathbf{C}\epsilon) = q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{u})\delta(\epsilon = \mathbf{C}^{-1}\mathbf{z}), \mathbf{C} = (\mathbf{I} - \mathbf{A}^T)^{-1}$
 - Generative model: $p_{\theta}(\epsilon, \mathbf{z}|\mathbf{u}) = p_{\epsilon}(\epsilon)p_{\theta}(\mathbf{z}|\mathbf{u}),$
- Decomposed ELBO

$$\text{ELBO} = \mathbb{E}_{q_{\mathcal{X}}}[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{u})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathcal{D}(q_{\phi}(\epsilon|\mathbf{x}, \mathbf{u}) || p_{\epsilon}(\epsilon)) - \mathcal{D}(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{u}) || p_{\theta}(\mathbf{z}|\mathbf{u}))].$$

$$p_{\epsilon}(\epsilon) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p_{\theta}(\mathbf{z}|\mathbf{u}) = \prod_i^n p_{\theta}(z_i|u_i), p_{\theta}(z_i|u_i) = \mathcal{N}(\lambda_1(u_i), \lambda_2^2(u_i)),$$



Loss function

$$l_u = \mathbb{E}_{q_{\mathcal{X}}} \|\mathbf{u} - \sigma(\mathbf{A}^T \mathbf{u})\|_2^2 \leq \kappa_1, \quad (11)$$

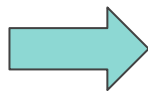
$$l_m = \mathbb{E}_{\mathbf{z} \sim q_{\phi}} \sum_{i=1}^n \|z_i - g_i(\mathbf{A}_i \circ \mathbf{z}; \boldsymbol{\eta}_i)\|^2 \leq \kappa_2, \quad (12)$$

$$H(\mathbf{A}) \equiv \text{tr}\left(\left(\mathbf{I} + \frac{c}{m} \mathbf{A} \circ \mathbf{A}\right)^n\right) - n = 0, \quad (13)$$

Constraints

maximize ELBO,

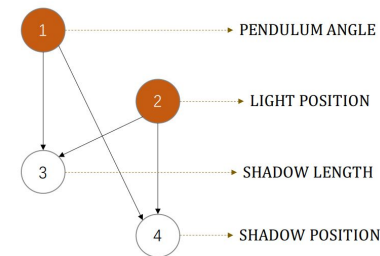
s.t. (11)(12)(13).



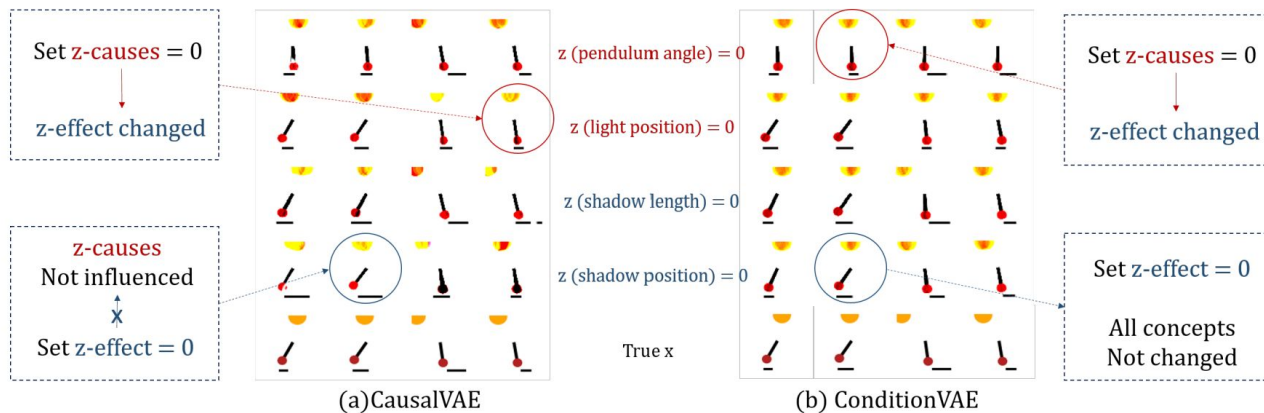
Lagrangian multiplier
method

$$\mathcal{L} = -\text{ELBO} + \alpha H(\mathbf{A}) + \beta l_u + \gamma l_m,$$

The result of intervention

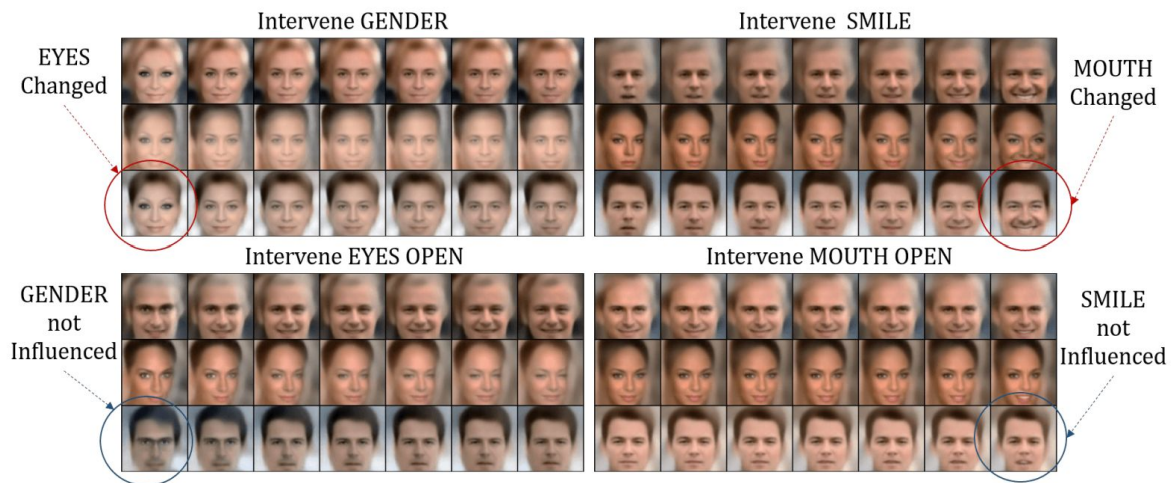


Causal graph

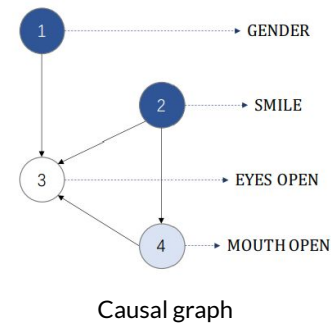


The result of intervention on pendulum dataset

Another example



Intervention results of CausalVAE model on CelebA (SMILE)





methods comparison based on MIC and TIC metrics

- Datasets (Synthetic : pendulum , flow
Real world benchmark : CalebA)
- Supervised : conditionVAE
- Unsupervised : beta_vae , causalvae-unsup , ladderVae

The MIC and TIC between learned representation z and the label u

Metrics(%)	CausalVAE		ConditionVAE		β -VAE		CausalVAE-unsup		LadderVAE	
	MIC	TIC	MIC	TIC	MIC	TIC	MIC	TIC	MIC	TIC
Pendulum	95.1 \pm 2.4	81.6 \pm 1.9	93.8 \pm 3.3	80.5 \pm 1.4	22.6 \pm 4.6	12.5 \pm 2.2	21.2 \pm 1.4	12.0 \pm 1.0	22.4 \pm 3.1	12.8 \pm 1.2
Flow	72.1 \pm 1.3	56.4 \pm 1.6	75.5 \pm 2.3	56.5 \pm 1.8	23.6 \pm 3.2	12.5 \pm 0.6	22.8 \pm 2.7	12.4 \pm 1.4	34.3 \pm 4.3	24.4 \pm 1.5
CelebA(SMILE)	83.7 \pm 6.2	71.6 \pm 7.2	78.8 \pm 10.9	66.1 \pm 12.1	22.5 \pm 1.2	9.92 \pm 1.2	27.2 \pm 5.3	14.6 \pm 4.2	23.5 \pm 3.0	10.3 \pm 1.6
CelebA(BEARD)	92.3 \pm 5.6	83.3 \pm 8.6	89.8 \pm 6.2	78.7 \pm 7.7	22.4 \pm 1.9	9.82 \pm 2.2	11.4 \pm 1.5	20.0 \pm 2.2	23.5 \pm 3.0	8.1 \pm 1.2



Conclusion

According to the authors, CAUSALVAE is the **first work on causal disentanglement**.

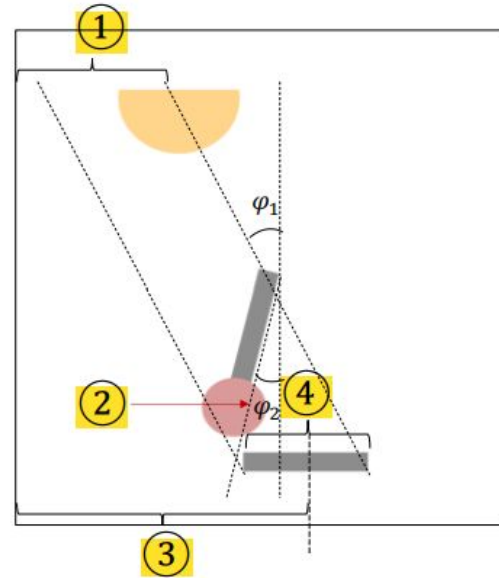
- It allows to discover causal relationships among the ground-truth factors
 - Prior knowledge can be eventually incorporated into the Adjacency matrix
- It is **identifiable**
- It supports the so called **do-operation**
- It considers **linear** causal relationships
- It requires **full knowledge of the ground-truth factors**

Code explanation

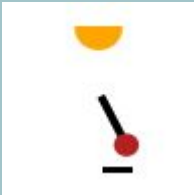




Data Preprocessing



output



Thanks for your attention