

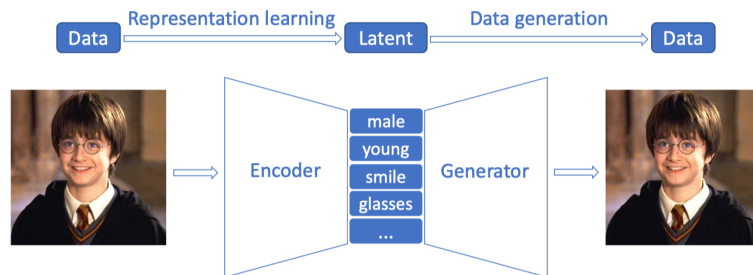
Disentangled Generative Causal Representation Learning

Xinwei Shen, Furui Liu, Hanze Dong, Qing Lian, Zhitang
Chen, Tong Zhang

HKUST

Representation Learning and Disentanglement

Representation learning and generation



- Observed data $x \sim q_x$ on $\mathcal{X} \subseteq \mathbb{R}^d$
- Latent variable $z \sim p_z$ on $\mathcal{Z} \subseteq \mathbb{R}^k$
- Bidirectional generative model: learning an *encoder* $E : \mathcal{X} \rightarrow \mathcal{Z}$ (to learn representations) and a *generator* $G : \mathcal{Z} \rightarrow \mathcal{X}$ (to generate data).
- Example: variational auto-encoder (VAE)

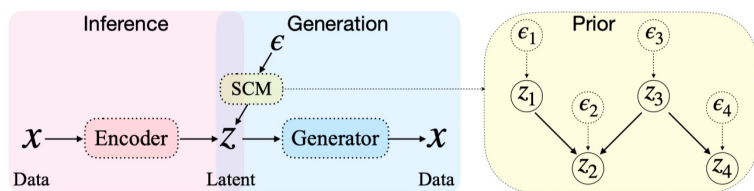
Disentanglement

Disentanglement as a common goal:

- In representation learning, an effective representation for downstream learning tasks should disentangle the underlying factors of variation.
- In generation, it is highly desirable if one can control the semantic generative factors.
- Both goals can be achieved with the *disentanglement* of latent variable z , which informally means that each dimension of z measures a distinct factor of variation in the data (Bengio et al., 2013).

Formulation

Generative model with a causal prior



- We adopt the general nonlinear Structural Causal Model (SCM):

$$f(z) = A^\top f(z) + h(\epsilon), \quad (3)$$

$$z = f^{-1}((I - A^\top)^{-1}h(\epsilon)) := F_\beta(\epsilon), \quad (4)$$

where ϵ denotes the exogenous variables, $A \in \mathbb{R}^{k \times k}$ is the weighted adjacency matrix, f and h are element-wise nonlinear transformations.

- (3) enables intervention; (4) enables generation.

Supervised regularizer

- Let $\xi \in \mathbb{R}^m$ be the underlying factors of x , and y_i be some continuous or discrete observation of factor ξ_i satisfying $\xi_i = \mathbb{E}(y_i|x)$ for $i = 1, \dots, m$.
- Let $\bar{E}(x)$ be the deterministic part of the stochastic transformation $E(x)$, i.e., $\bar{E}(x) = \mathbb{E}(E(x)|x)$, which is used for representation learning.
- We consider the following objective:

$$L(E, G) = L_{\text{gen}}(E, G) + \lambda L_{\text{sup}}(E), \quad (2)$$

where

- $L_{\text{sup}} = \sum_{i=1}^m \mathbb{E}_{(x,y)}[\text{CE}(\bar{E}_i(x), y_i)]$ if y_i is the binary or bounded continuous label of ξ_i ;
- $L_{\text{sup}} = \sum_{i=1}^m \mathbb{E}_{(x,y)}[\bar{E}_i(x) - y_i]^2$ if y_i is the continuous observation of ξ_i .

Algorithm

Algorithm

Algorithm 1: Disentangled generative cAusal Representation (DEAR) Learning

Input: training set $\{x_1, \dots, x_n, y_1, \dots, y_n\}$, initial parameter $\phi, \theta, \beta, \psi$, batch size n

```

1 while not convergence do
2   for multiple steps do
3     Sample  $\{x_1, \dots, x_n\}$  from the training set,  $\{\epsilon_1, \dots, \epsilon_n\}$  from  $\mathcal{N}(0, I)$ 
     Generate from the causal prior  $z_i = F_\beta(\epsilon_i), i = 1, \dots, n$ 
     Update  $\psi$  by descending the stochastic gradient:
      $\frac{1}{n} \sum_{i=1}^n \nabla_\psi \left[ \log(1 + e^{-D_\psi(x_i, E_\phi(x_i))}) + \log(1 + e^{D_\psi(G_\theta(z_i), z_i)}) \right]$ 
4   Sample  $\{x_1, \dots, x_n, y_1, \dots, y_n\}, \{\epsilon_1, \dots, \epsilon_n\}$  as above; generate  $z_i = F_\beta(\epsilon_i)$ 
     Compute  $\theta$ -gradient:  $-\frac{1}{n} \sum_{i=1}^n s(G_\theta(z_i), z_i) \nabla_\theta D_\psi(G_\theta(z_i), z_i)$ 
     Compute  $\phi$ -gradient:  $\frac{1}{n} \sum_{i=1}^n \nabla_\phi D_\psi(x_i, E_\phi(x_i)) + \frac{1}{n_s} \sum_{i=1}^{n_s} \nabla_\phi L_{\text{sup}}(\phi; x_i, y_i)$ 
     Compute  $\beta$ -gradient:  $-\frac{1}{n} \sum_{i=1}^n s(G(z_i), z_i) \nabla_\beta D_\psi(G_\theta(F_\beta(\epsilon_i)), F_\beta(\epsilon_i))$ 
     Update parameters  $\phi, \theta, \beta$  using the gradients

```

Return: ϕ, θ, β

Formulation of DEAR

- Rewrite the generative loss:

$$L_{\text{gen}}(\phi, \theta, \beta) = D_{\text{KL}}(q_\phi(x, z), p_{\theta, \beta}(x, z)). \quad (5)$$

- Formulation to learn disentangled generative causal representations:

$$\min_{\phi, \theta, \beta} L(\phi, \theta, \beta) := L_{\text{gen}}(\phi, \theta, \beta) + \lambda L_{\text{sup}}(\phi). \quad (6)$$

Theory

Identifiability of disentanglement

Theorem

Assume the infinite capacity of E , G and f . Further assume the true binary adjacency matrix can be learned. Then DEAR learns the disentangled encoder E^* . Specifically, we have $g_i(\xi_i) = \sigma^{-1}(\xi_i)$ if CE loss is used in the supervised regularizer, and $g_i(\xi_i) = \xi_i$ if L_2 loss is used.

Optimization

- The SCM prior $p_\beta(z)$ and implicit generated conditional $p_\theta(x|z)$ make L_{gen} in (5) lose an analytic form.
- The lemma gives the gradient.
- We adopt a GAN method to adversarially estimate the gradient of L_{gen} as in Shen et al. (2020).

Lemma (Gradient)

Let $r(x, z) = q(x, z)/p(x, z)$ and $\mathcal{D}(x, z) = \log r(x, z)$. Then we have

$$\nabla_\theta L_{\text{gen}} = -\mathbb{E}_{z \sim p_\beta(z)} [s(x, z) \nabla_x \mathcal{D}(x, z)^\top |_{x=G_\theta(z)} \nabla_\theta G_\theta(z)],$$

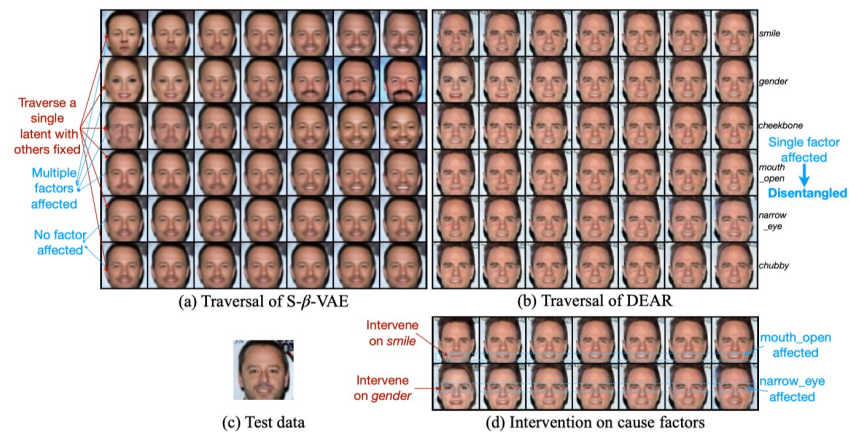
$$\nabla_\phi L_{\text{gen}} = \mathbb{E}_{x \sim q_x} [\nabla_z \mathcal{D}(x, z)^\top |_{z=E_\phi(x)} \nabla_\phi E_\phi(x)],$$

$$\nabla_\beta L_{\text{gen}} = -\mathbb{E}_\epsilon [s(x, z) (\nabla_x \mathcal{D}(x, z)^\top \nabla_\beta G(F_\beta(\epsilon)) + \nabla_z \mathcal{D}(x, z)^\top \nabla_\beta F_\beta(\epsilon)) |_{\substack{x=G(F_\beta(\epsilon)) \\ z=F_\beta(\epsilon)}}],$$

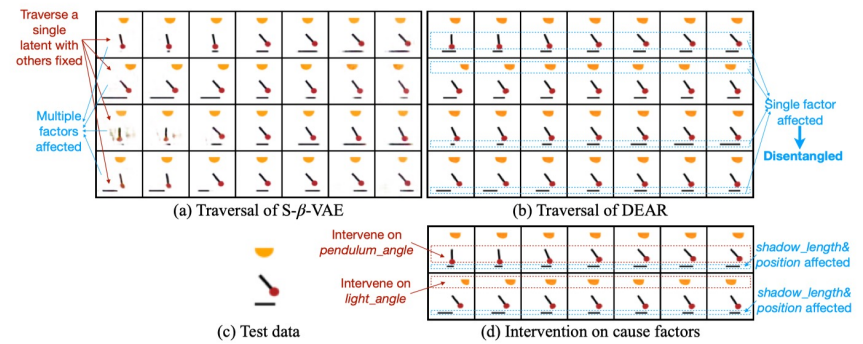
where $s(x, z) = e^{\mathcal{D}(x, z)}$ is the scaling factor.

Controllable Generation

Causal controllable generation (CelebA)



Causal controllable generation (Pendulum)



Better Representations

Distributional robustness

Table: The worst-case and average test accuracy.

Method	(a) CelebA		(b) Pendulum	
	WorstAcc(%)	AvgAcc(%)	WorstAcc(%)	AvgAcc(%)
ERM	59.12 \pm 1.78	82.12 \pm 0.26	60.48 \pm 2.73	87.40 \pm 0.89
DEAR-lin-10%	71.40 \pm 0.47	81.04 \pm 0.14	63.93 \pm 1.33	89.70 \pm 0.63
DEAR-nlr-10%	70.44 \pm 1.02	81.94 \pm 0.31	65.59 \pm 1.90	90.19 \pm 0.63
ERM-multilabel	59.17 \pm 4.02	82.05 \pm 0.25	61.70 \pm 4.02	87.20 \pm 1.00
S-VAE	60.54 \pm 3.48	79.51 \pm 0.58	20.78 \pm 4.45	84.26 \pm 1.31
S- β -VAE	63.85 \pm 2.09	80.82 \pm 0.19	44.12 \pm 9.73	86.99 \pm 1.78
S-TCVAE	64.93 \pm 3.30	81.58 \pm 0.14	35.50 \pm 5.57	86.64 \pm 1.15
DEAR-lin	76.05 \pm 0.70	83.56 \pm 0.09	74.95 \pm 1.26	93.61 \pm 0.13
DEAR-nlr	71.37 \pm 0.66	83.81 \pm 0.08	72.48 \pm 0.74	93.11 \pm 0.14

Sample efficiency

- Statistical efficiency score: the average test accuracy based on 100 samples divided by the average accuracy based on 10,000/all samples (Locatello et al., 2019).

Table: Sample efficiency and test accuracy with different training sample sizes.

Method	(a) CelebA			(b) Pendulum		
	100(%)	10,000(%)	Eff(%)	100(%)	all(%)	Eff(%)
ResNet	68.06 \pm 0.19	79.51 \pm 0.31	85.59 \pm 0.27	79.71 \pm 0.98	90.64 \pm 1.57	87.97 \pm 2.11
DEAR-lin-10%	78.09 \pm 0.59	79.54 \pm 0.41	98.18 \pm 0.49	88.93 \pm 1.40	93.18 \pm 0.18	95.43 \pm 1.33
DEAR-nlr-10%	80.30 \pm 0.24	80.87 \pm 0.12	99.29 \pm 0.23	87.65 \pm 0.46	91.27 \pm 0.21	96.03 \pm 0.29
ResNet-pretrain	76.84 \pm 2.08	83.75 \pm 0.93	91.74 \pm 1.98	79.59 \pm 0.93	89.16 \pm 1.60	89.28 \pm 0.59
S-VAE	77.07 \pm 1.42	79.87 \pm 1.67	96.49 \pm 1.68	84.16 \pm 0.69	90.89 \pm 0.28	92.60 \pm 0.49
S- β -VAE	71.78 \pm 1.99	76.63 \pm 0.24	93.67 \pm 2.41	79.95 \pm 1.65	87.87 \pm 0.52	90.98 \pm 1.47
S-TCVAE	77.10 \pm 2.08	81.63 \pm 0.20	94.45 \pm 2.72	85.36 \pm 1.11	90.33 \pm 0.33	94.51 \pm 1.31
DEAR-lin	83.51 \pm 0.77	84.92 \pm 0.11	98.34 \pm 0.81	90.21 \pm 0.94	93.31 \pm 0.14	96.68 \pm 0.89
DEAR-nlr	84.44 \pm 0.48	85.10 \pm 0.09	99.23 \pm 0.51	90.62 \pm 0.32	92.57 \pm 0.08	97.93 \pm 0.29

Conclusion

- We identified a problem with previous methods using the independent prior assumption, and proved that they fail to disentangle when the underlying factors are causally correlated.
- We proposed a new disentangled learning method, DEAR, which integrates an SCM prior into a bidirectional generative model, trained with a suitable GAN loss.
- We provided theoretical justifications on the identifiability of the formulation and the asymptotic consistency of our algorithm.
- Extensive experiments were conducted to demonstrate the effectiveness of DEAR in causal controllable generation, and the benefits of the learned representations for downstream tasks.