

CausalVAE: Disentangled Representation Learning via Neural Structural Causal Models

Mengyue Yang¹ Furui Liu¹ Zhitang Chen¹ Xinwei Shen² Jianye Hao¹
Jun Wang³

¹Noah's Ark Lab, Huawei, Shenzhen, China

²The Hong Kong University of Science and Technology, Hong Kong, China

³University College London, London, United Kingdom

Introduction

Disentangled Representations

"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" [6].

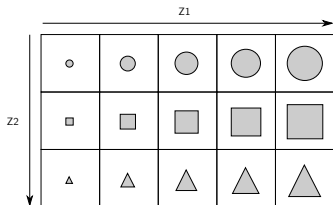


Figure 1: Toy example where each dimension of the latent space z controls a given generative factor: z_1 the size, z_2 the shape of the 2D objects.

A significant amount of work in this area has been carried out for generative models like VAE [13, 21].

- Unsupervised disentanglement learning [8, 2, 11, 17, 5, 26]
- Impossibility result [18]
- Disentanglement through additionally observed variables [12, 7, 23, 20, 15, 1, 9, 19, 4, 3, 10]

What if ground-truth factors are not independent?

Disentangled representation learning builds on the assumption that the **ground-truth factors are independent**.

In many real world applications, they might be causally related.



Figure 2: A swinging pendulum: an illustrative example.

Causal Graphs in one slide

Causal graphs (also known as causal Bayesian networks or **DAGs**) are probabilistic graphical models used to encode assumptions about the data-generating process¹.

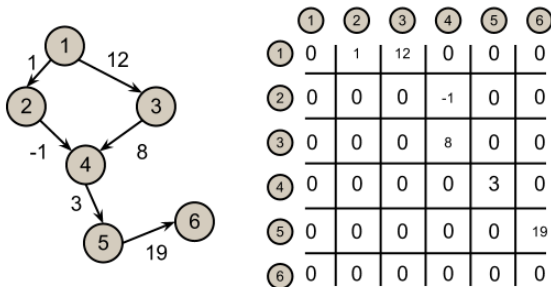


Figure 3: Example of causal graph (left) and its adjacency matrix (right).

¹We focus on linear causal models.

CausalVAE

Model structure of CausalVAE

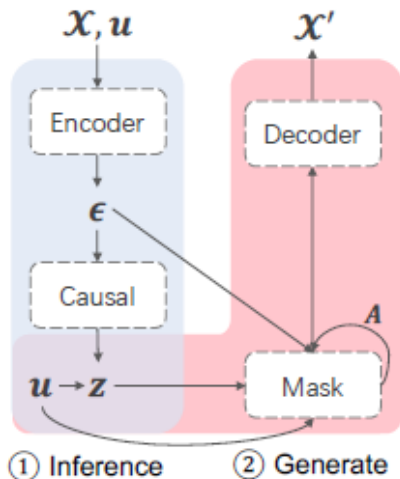


Figure 4: Model structure of CAUSALVAE

The Causal Layer

The key element of CAUSALVAE is the presence of a **Causal Layer**.

- The Causal Layer describes a Structural Causal Model (SCM) [22].
- The Causal Layer transforms **independent exogenous factors** ϵ into **causal endogenous factors** z corresponding to the ground-truth factors of interest.
 - The causal structure is defined by a DAG with adjacency matrix A .
- The Causal Layer is here restricted to implement a linear SCM:

$$z = A^T z + \epsilon = B^T \epsilon \quad (1)$$

A Probabilistic Generative Model for CausalVAE (1)

Let $x \in \mathbb{R}^d$ and $u \in \mathbb{R}^n$ be two observed variables². Let $\epsilon \in \mathbb{R}^n$ be the latent exogenous independent variables and $z \in \mathbb{R}^n$ be the latent endogenous variables. Equation (1) defines the relation between z and ϵ . Then, consider the following generative model:

$$p_{\theta}(x, z, \epsilon|u) = p_{\theta}(x|z, \epsilon, u)p_{\theta}(z, \epsilon|u), \quad (2)$$

where:

$$p_{\theta}(z, \epsilon|u) = p_{\theta}(z|u)p_{\theta}(\epsilon) \quad (3)$$

²In CAUSALVAE u is restricted to be the explicit label associated to each ground-truth factor.

A Probabilistic Generative Model for CausalVAE (2)

Let introduce the generative and inference models:

$$p_{\theta}(x|z, \epsilon, u) = p_{\theta}(x|z), \quad (4)$$

$$q_{\phi}(z, \epsilon|x, u) \quad (5)$$

corresponding to the following decoding and encoding processes:

$$x = f(z) + \xi, \quad (6)$$

$$z = h(x, u) + \zeta, \quad (7)$$

where ξ and ζ are the vectors of independent noise with probability densities p_{ξ} and p_{ζ} ³.

³When ξ and ζ are infinitesimal, the encoder and decoder can be regarded as deterministic ones.

ELBO of CausalVAE

Given data set \mathcal{D} with the empirical data distribution $q_{\mathcal{D}}(x, u)$, the parameters θ and ϕ are learned by optimizing the following evidence lower bound (ELBO):

$$ELBO = \mathbb{E}_{q_{\mathcal{D}}} [\mathbb{E}_{q_{\phi}(z, \epsilon|x, u)} [\log p_{\theta}(x|z, \epsilon, u)] - KL(q_{\phi}(z, \epsilon|x, u) || p_{\theta}(z, \epsilon|u))]. \quad (8)$$

Combining eqs. (1), (3) and (4) into eq. (8), by abuse of notation, we can write:

$$\begin{aligned} ELBO &= \mathbb{E}_{q_{\mathcal{D}}} [\mathbb{E}_{q_{\phi}(z|x, u)} [\log p_{\theta}(x|z)] \\ &\quad - KL(q_{\phi}(\epsilon|x, u) || p_{\theta}(\epsilon)) \\ &\quad - KL(q_{\phi}(z|x, u) || p_{\theta}(z|u))]. \end{aligned} \quad (9)$$

Learning the Causal Structure of Latent Codes

We need some additional constraints:

$$H(\mathbf{A}) = \text{tr}((\mathbf{I} + \mathbf{A} \circ \mathbf{A})^n) - n = 0, \quad (10)$$

$$l_u = \mathbb{E}_{q_D} \|\mathbf{u} - \sigma \mathbf{A}^T \mathbf{u}\|_2^2 \leq k_1. \quad (11)$$

Equation (10) is a DAG constraint [25].

Then:

$$\mathcal{L}_{\text{CAUSALVAE}} = \text{ELBO} - \alpha H(\mathbf{A}) - \beta l_u. \quad (12)$$

Experiments

Datasets PENDULUM (pendulum angle, light angle, shadow length, shadow location), CELEBA (gender, smile, eyes open, mouth open).

Metrics Maximal Information Coefficient (MIC) and Total Information Coefficient (TIC) [14].

Comparing Methods β -VAE [8], LADDERVAE [24], DC-IGN [16].

MIC and TIC Results

Metrics(%)	CausalVAE		DC-IGN		β -VAE		CausalVAE-unsup		LadderVAE	
	MIC	TIC	MIC	TIC	MIC	TIC	MIC	TIC	MIC	TIC
Pendulum	96.3 \pm 3.6	89.0 \pm 2.9	61.8 \pm 8.7	48.1 \pm 7.3	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
CelebA	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

Figure 5: The MIC and TIC between learned representation z and the label u .

The Learning Process of a Causal Matrix A

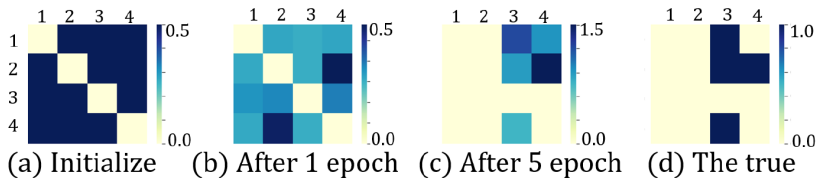


Figure 6: The learning process of causal matrix A for CELEBA. (1=gender, 2=smile, 3=eyes open, 4=mouse open).

Intervention Experiments

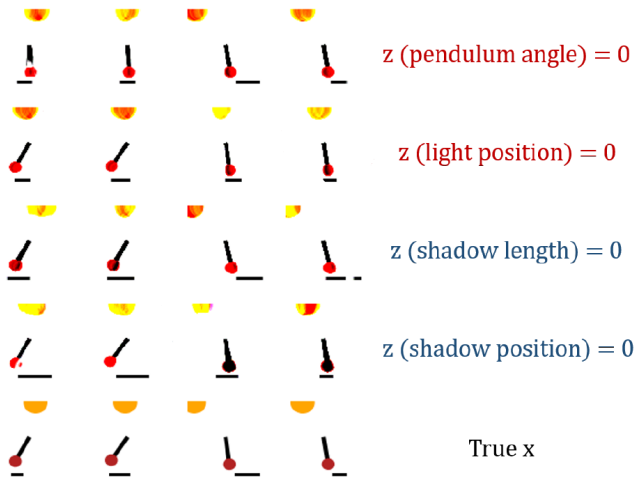


Figure 7: The results of Intervention experiments on the pendulum dataset.

Conclusions

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**

- [1] D. Bouchacourt, R. Tomioka, and S. Nowozin.
Multi-level variational autoencoder: Learning disentangled representations from grouped observations.
In Proc. of the 32nd AAAI Conf. on Artif. Intel., AAAI, 2018.
- [2] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner.
Understanding disentangling in β -vae.
In Proc. of the 30th Int. Conf. on Neural Inf. Proc. Sys., NeurIPS, 2017.
- [3] J. Chen and K. Batmanghelich.
Robust ordinal vae: Employing noisy pairwise comparisons for disentanglement.
ArXiv, 2020.

- [4] J. Chen and K. Batmanghelich.
Weakly supervised disentanglement by pairwise similarities.
In Proc. of the 34th AAAI Conf. on Artif. Intel., AAAI, 2020.
- [5] T. Q. Chen, X. Li, R. B. Grosse, and D. K. Duvenaud.
Isolating sources of disentanglement in variational autoencoders.
In Proc. of the 31st Int. Conf. on Neural Inf. Proc. Sys., NeurIPS, 2018.
- [6] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel.
Infogan: Interpretable representation learning by information maximizing generative adversarial nets.
In Proc. of the 29th Int. Conf. on Neural Inf. Proc. Sys., NeurIPS, 2016.

- [7] B. Cheung, J. A. Livezey, A. K. Bansal, and B. A. Olshausen.
Discovering hidden factors of variation in deep networks.
In *CoRR*, 2015.
- [8] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. M. Botvinick, S. Mohamed, and A. Lerchner.
beta-vae: Learning basic visual concepts with a constrained variational framework.
In *Proc. of the 5th Int. Conf. on Learn. Repr.*, ICLR, 2017.
- [9] H. Hosoya.
Group-based learning of disentangled representations with generalizability for novel contents.
In *Proc. of the 28th Int. Joint Conf. on Artif. Intel.*, IJCAI, 2019.

- [10] I. Khemakhem, D. P. Kingma., R. P. Mont, and A. Hyvärinen.
Variational autoencoders and nonlinear ica: A unifying framework.
In *Proc. of the 23rd Int. Conf. on Artif. Intel. and Stat.*, AISTATS, 2020.
- [11] H. Kim and A. Mnih.
Disentangling by factorising.
In *Proc. of the 35th Int. Conf. on Mach. Learn.*, ICML, 2018.
- [12] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling.
Semi-supervised learning with deep generative models.
In *Proc. of the 27th Int. Conf. on Neural Inf. Proc. Sys.*, NeurIPS, 2014.

- [13] D. P. Kingma and M. Welling.
Auto-encoding variational bayes.
In Proc. of the 2nd Int. Conf. on Learn. Repr., ICLR, 2014.
- [14] J. B. Kinney and G. S. Atwal.
Equitability, mutual information, and the maximal information coefficient.
Proc. of the Nat. Acad. of Sc. of the Unit. Stat. of Americ., 111, 02 2014.
- [15] J. Klys, J. Snell, and R. Zemel.
Learning latent subspaces in variational autoencoders.
In Proc. of the 31st Int. Conf. on Neural Inf. Proc. Sys., NeurIPS, 2018.

- [16] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum.
Deep convolutional inverse graphics network.
In *Proc. of the 28th Int. Conf. on Neural Inf. Proc. Sys.*, NeurIPS, 2015.
- [17] A. Kumar, P. Sattigeri, and A. Balakrishnan.
Variational inference of disentangled latent concepts from unlabeled observations.
In *Proc. of the 6th Int. Conf. on Learn. Repr.*, ICLR, 2018.
- [18] F. Locatello, S. Bauer, M. Lucic, S. Gelly, B. Schölkopf, and O. Bachem.
Challenging common assumptions in the unsupervised learning of disentangled representations.
In *Proc. of the 36th Int. Conf. on Mach. Learn.*, ICML, 2019.

- [19] F. Locatello, B. Poole, G. Rätsch, B. Schölkopf, O. Bachem, and M. Tschannen.
Weakly-supervised disentanglement without compromises.
In *Proc. of the 37th Int. Conf. on Mach. Learn.*, ICML, 2020.
- [20] F. Locatello, M. Tschannen, S. Bauer, G. Rätsch, B. Schölkopf, and O. Bachem.
Disentangling factors of variation using few labels.
In *Proc. of the 8th Int. Conf. on Learn. Repr.*, ICLR, 2020.
- [21] D. J. Rezende, S. Mohamed, and D. Wierstra.
Stochastic backpropagation and approximate inference in deep generative models.
In *Proc. of the 31st Int. Conf. on Mach. Learn.*, ICML, 2014.

- [22] S. Shimizu, P. O. Hoyer, A. Hyvärinen, and A. Kerminen.
A linear non-gaussian acyclic model for causal discovery.
JMLR, 7(72):2003–2030, 2006.
- [23] N. Siddharth, B. Paige, J. van de Meent, A. Desmaison,
N. Goodman, P. Kohli, F. Wood, and P. Torr.
**Learning disentangled representations with semi-supervised
deep generative models.**
In *Proc. of the 30th Int. Conf. on Neural Inf. Proc. Sys.*, NeurIPS,
2017.
- [24] C. K. Sønderby, T. Raiko, L. Maaløe, S. r. K. Sønderby, and
O. Winther.
Ladder variational autoencoders.
In *Proc. of the 29th Int. Conf. on Neural Inf. Proc. Sys.*, NeurIPS,
2016.

- [25] Y. Yu, J. Chen, T. Gao, and M. Yu.
Dag-gnn: Dag structure learning with graph neural networks.
In *Proc. of the 36th Int. Conf. on Mach. Learn.*, ICML, 2019.
- [26] S. Zhao, J. Song, and S. Ermon.
Infovae: Balancing learning and inference in variational autoencoders.
In *Proc. of the 33rd AAAI Conf. on Artif. Intel.*, AAAI, 2019.