

Variational Autoencoders

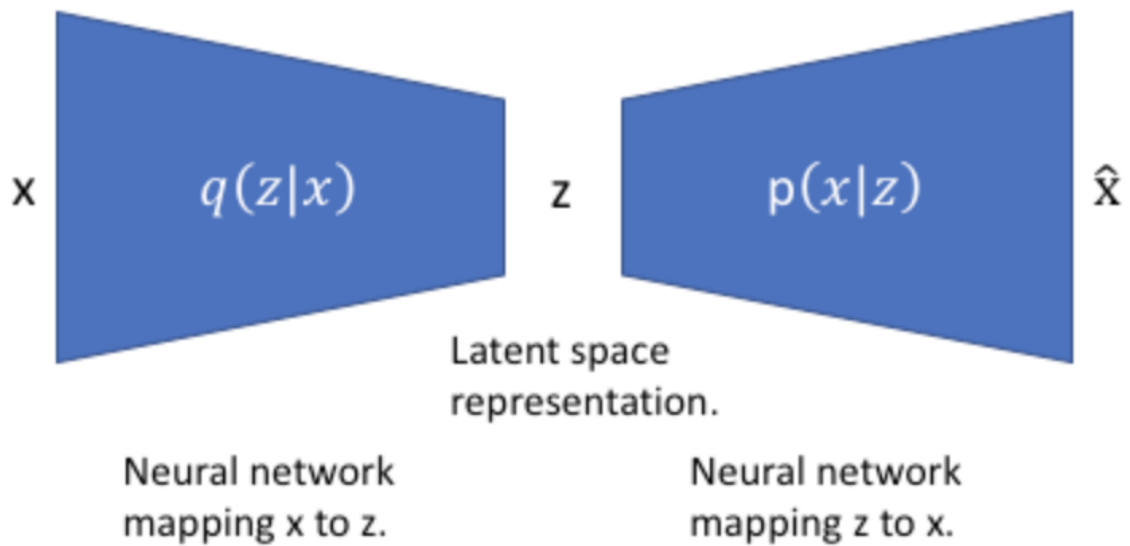
Beta VAE, Ladder VAE, Causal VAE

Presenter: Arshdeep Sekhon
<https://qdata.github.io/deep2Read>

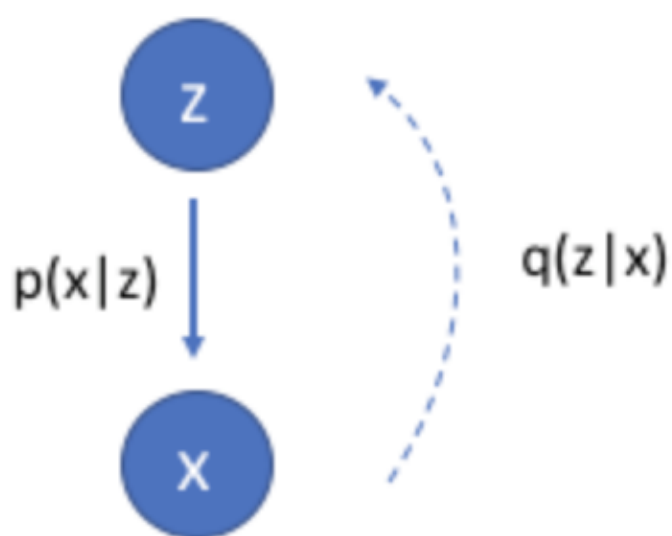
β -VAE: Learning Basic Visual Concepts with a Constrained Variational Framework

Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess,
Xavier Glorot, Matthew Botvinick, Shakir Mohamed, Alexander
Lerchner

Variational Auto-Encoder



Variational Auto-Encoder



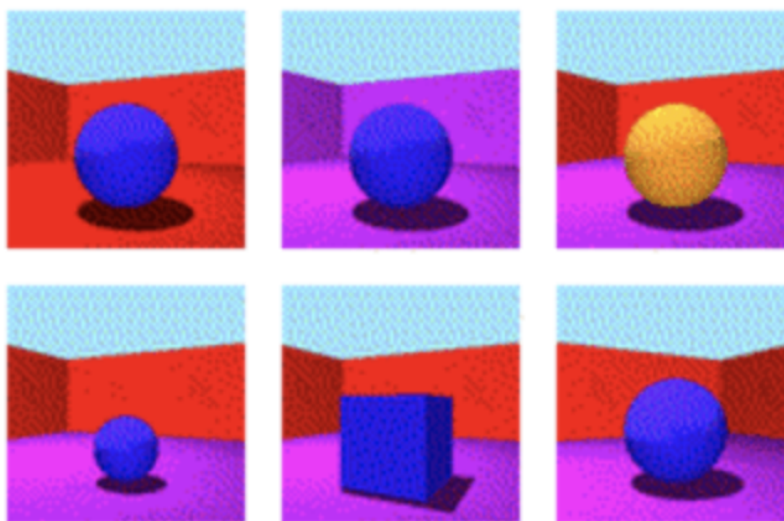
Maximize ELBO:

$$E_{q(z|x)} \log(p(x|z)) - KL(q(z|x) || p(z)) \quad (1)$$

- ▶ Term 1: Reconstruction
- ▶ Term 2: prior

Objective

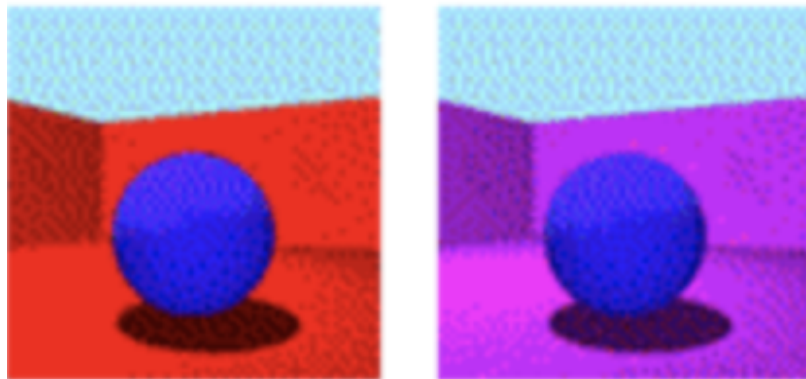
Learning independent factors of generation in an unsupervised manner



- In the above image, generated image depends on : color of walls, size, shape, color of object

Disentangled representations

- ▶ If we change one latent factor in a disentangled representation, it corresponds to changes in only one generative factor
- ▶ For example, in the above 3D image, change z corresponding to wall colors, image remains same, and only background changes.



Why disentangled representations

- ▶ generalize better to unseen situations: useful in zero-shot or knowledge transfer
- ▶ boost AI performance: (Lake 2016)

β -VAE

- ▶ generative model to learn a disentangled \mathbf{z} in an *unsupervised* manner
- ▶ no prior information regarding number of factors, or correspondence

Method

- ▶ Given dataset $D = \{\mathbf{X}, \mathbf{V}, \mathbf{W}\}$
- ▶ images $\mathbf{x} \in R^N$
- ▶ conditionally independent factors $\mathbf{v} \in R^K$
- ▶ conditionally dependent factors $\mathbf{w} \in R^H$
- ▶ these factors are independent $\log(p(\mathbf{v}|\mathbf{x})) = \sum_k \log p(\mathbf{v}_k|\mathbf{x})$
- ▶ assumption: $p(\mathbf{x}|\mathbf{v}, \mathbf{w}) = \text{Sim}(\mathbf{v}, \mathbf{w})$, where Sim is the true simulator that generates images given \mathbf{v}, \mathbf{w}

Goal

Learn a joint distribution of $p(\mathbf{x}, \mathbf{z})$ where $\mathbf{z} \in R^M$, $M \geq K$, such that

$$p(\mathbf{x}|\mathbf{z}) \sim p(\mathbf{x}|\mathbf{v}, \mathbf{w}) = \text{Sim}(\mathbf{v}, \mathbf{w}) \quad (2)$$

Maximize the marginal log likelihood of observed data \mathbf{x} :

$$\max_{\theta} E_{p_{\theta}(\mathbf{z})}[p_{\theta}(\mathbf{x}|\mathbf{z})] \quad (3)$$

Posterior of the inferred latent factors \mathbf{z} :

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \quad (4)$$

Goal is to ensure $q_{\phi}(\mathbf{z}|\mathbf{x})$ corresponds to \mathbf{v} in a disentangled manner.

β -VAE

To achieve disentanglement or statistical Independence, set the prior to be an isotropic Gaussian :

$$p(\mathbf{z}) = N(0, \mathbf{I}) \quad (5)$$

$$\max_{\theta, \phi} E_{\mathbf{x} \sim D} E_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) \quad (6)$$

$$s.t. D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) < \epsilon \quad (7)$$

Rewrite as Lagrangian under the KKT conditions:

$$\max_{\theta, \phi} E_{\mathbf{x} \sim D} E_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \beta(D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) - \epsilon) \quad (8)$$

$$\max_{\theta, \phi} E_{\mathbf{x} \sim D} E_{q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) - \beta(D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))) \quad (9)$$

Trade off

- ▶ Term 1 encourages better representations for Reconstruction fidelity
- ▶ Term 2 or high beta values try to make the dimensions as independent of each other as possible

Likelihood is a poor metric to measure disentanglement

- ▶ Disentangled representations emerge when the right balance is found between reconstruction cost as regularisation and latent channel capacity restriction ($\beta > 1$).
- ▶ $\beta > 1$ can lead to poorer reconstructions due to the loss of high frequency details when passing through a constrained latent bottleneck.
- ▶ need of a new metric to measure disentanglement

New Disentanglement Metric

- ▶ disentangled must be interpretable: can generate images of *small, green* apples, and *large, green* apples by varying the *small* latent
- ▶ independence can be obtained using PCA or ICA, but not interpretable!
- ▶ cross correlation is not a good metric
- ▶ target is to measure both Independence and interpretability

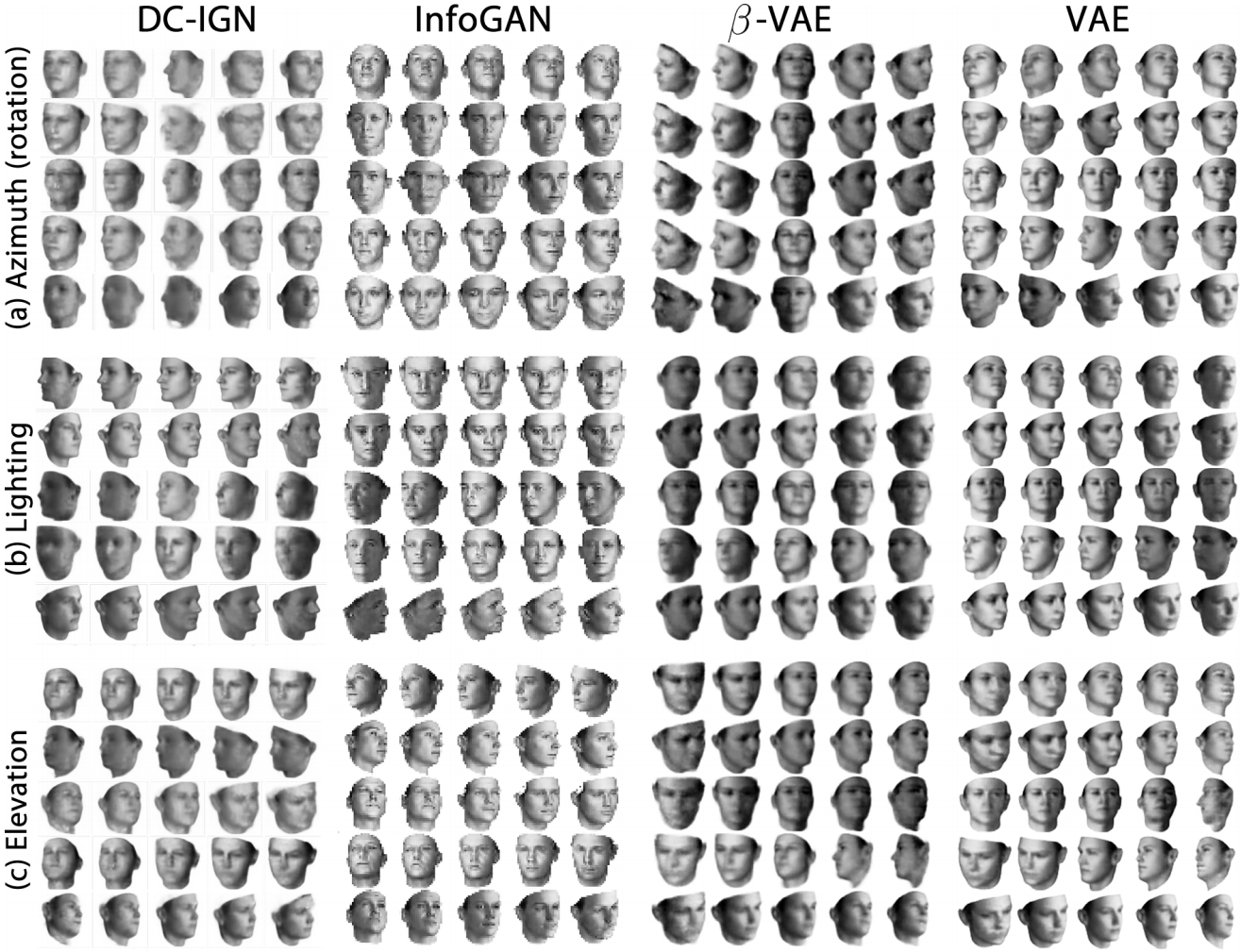
Disentanglement Metric

- ▶ we have labels of the generative factors $v \in V$ for some examples
- ▶ Choose an independent factor randomly $y \sim Unif[1, \dots, K]$
- ▶ For a batch of L samples
 - ▶ sample two sets of latent representations $\mathbf{v}_{1,l}$ and $\mathbf{v}_{2,l}$ enforcing $[\mathbf{v}_{1,l}]_k = [\mathbf{v}_{2,l}]_k$ if $k = y$ the value of the factor is fixed
 - ▶ Get $\mathbf{x}_{1,l} \sim Sim(\mathbf{v}_{1,1})$ then infer $\mathbf{z}_{1,l} = \mu(\mathbf{x}_{1,l})$ using encoder $q(\mathbf{z}|\mathbf{x}) \sim N(\mu(\mathbf{x}), \sigma(\mathbf{x}))$
 - ▶ similarly for the second batch $\mathbf{v}_{2,l}$
 - ▶ $\mathbf{z}_{diff}^b = |\mathbf{z}_{1,l} - \mathbf{z}_{2,l}|$ the linear difference between the two latent
 - ▶ take average across the batch $\mathbf{z}_{diff}^b = \frac{1}{L} \sum_{l=1}^L (\mathbf{z}_{diff}^l)$, and $p(y|\mathbf{z}_{diff}^b)$, this score is disentanglement metric

Disentanglement Metric

- ▶ accuracy of this classifier over multiple batches is used as disentanglement metric score.
- ▶ $p(y|\mathbf{z}_{diff}^b)$: linear classifier with low VC-dimension

Experiments



Experiments: Simulation Dataset

- ▶ 2D shapes
- ▶ Cartesian product of the shape and four independent generative factors defined in vector graphics: position X (32 values), position Y (32 values), scale (6 values) and rotation (40 values over the 2π range).

Experiments

Model	Disentanglement metric score
<i>Ground truth</i>	<i>100%</i>
Raw pixels	$45.75 \pm 0.8\%$
PCA	$84.9 \pm 0.4\%$
ICA	$42.03 \pm 10.6\%$
DC-IGN	$99.3 \pm 0.1\%$
InfoGAN	$73.5 \pm 0.9\%$
VAE untrained	$44.14 \pm 2.5\%$
VAE	$61.58 \pm 0.5\%$
β-VAE	$99.23 \pm 0.1\%$

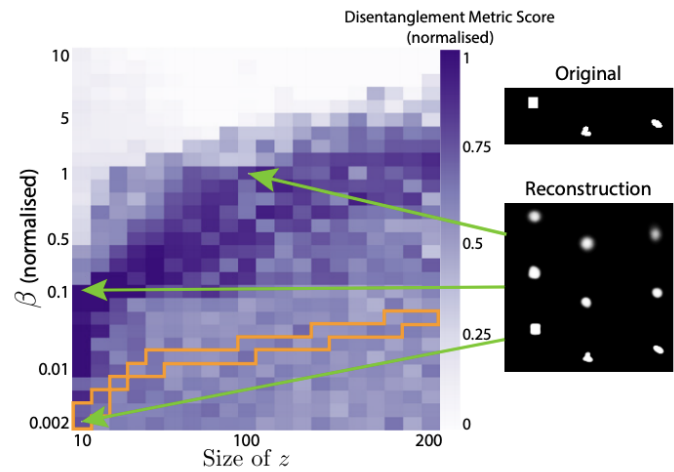
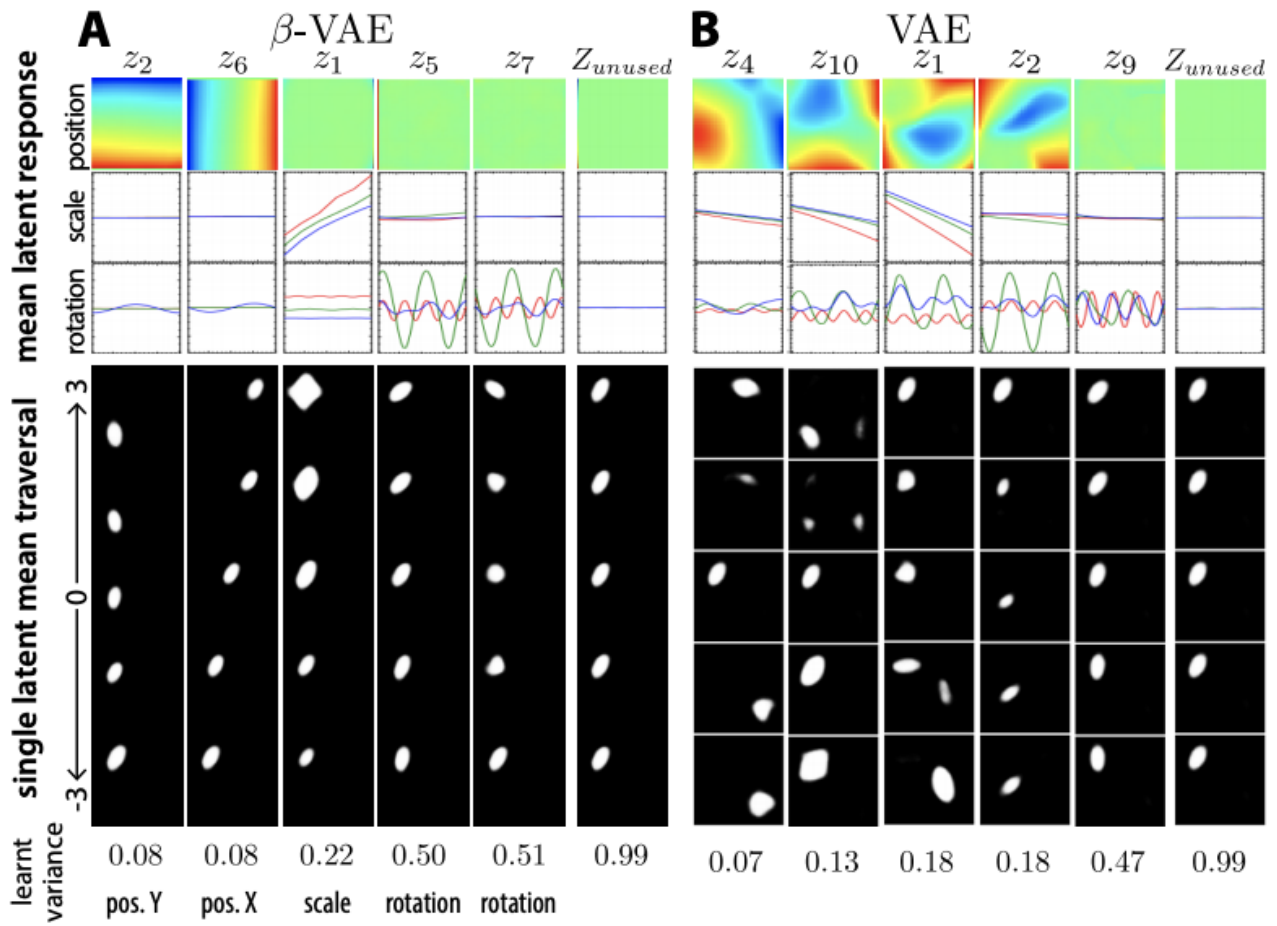
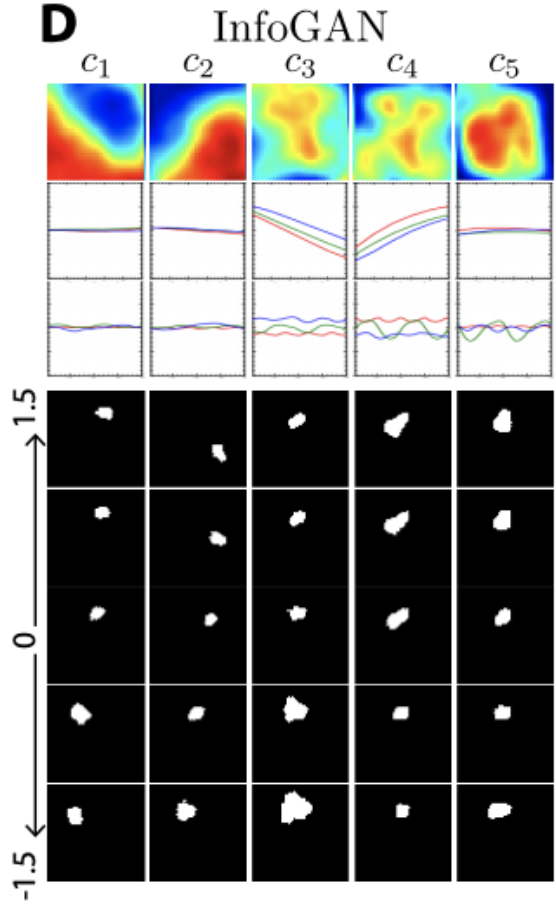
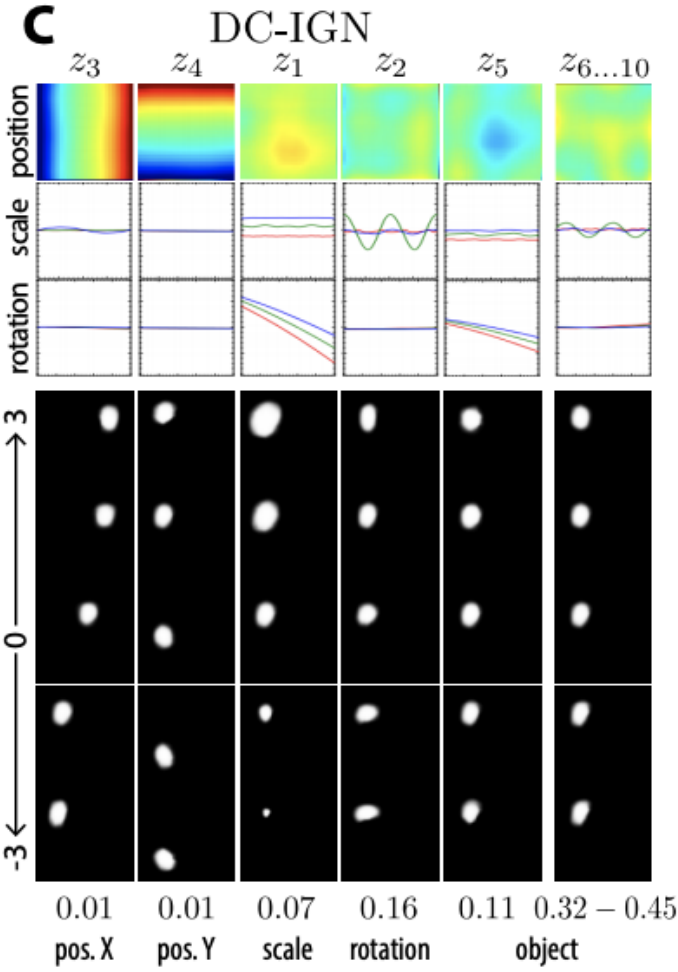


Figure 6: Disentanglement metric classification accuracy for 2D shapes dataset. **Left:** Accuracy for different models and training regimes **Right:** Positive correlation is present between the size of z and the optimal *normalised* values of β for disentangled factor learning for a fixed β -VAE architecture. β values are normalised by latent z size m and input x size n . Note that β values are not uniformly sampled. Orange approximately corresponds to *unnormalised* $\beta = 1$. Good reconstructions are associated with entangled representations (lower disentanglement scores). Disentangled representations (high disentanglement scores) often result in blurry reconstructions.

Experiments



Experiments



Conclusion

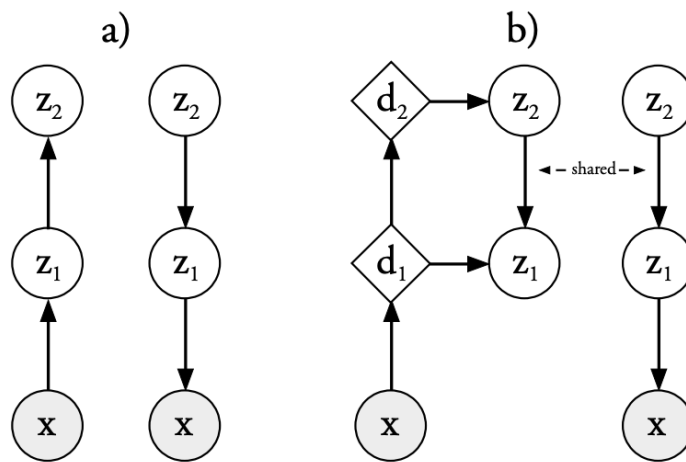
- ▶ reformulated the standard VAE framework as a constrained optimisation problem with strong latent capacity constraint and independence prior pressures.
- ▶ covers a wider range of factor values and is disentangled more cleanly than other benchmarks, all in a completely unsupervised manner

Ladder Variational Autoencoders

Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, Ole Winther

Motivation

- ▶ hierarchies of conditional stochastic variables
- ▶ structured inference model using the same top-down dependency structure both in the inference and generative models.



Ladder VAE

- ▶ VAEs and LVAEs simultaneously train a generative model $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})$ for data \mathbf{x} using latent variables \mathbf{z}
- ▶ inference model $q_{\phi}(\mathbf{z}|\mathbf{x})$
- ▶ variational lower bound to the likelihood $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z})d\mathbf{z}$.

VAE: generative model

$$p_{\theta}(\mathbf{z}) = p_{\theta}(\mathbf{z}_L) \prod_{i=1}^{L-1} p_{\theta}(\mathbf{z}_i | \mathbf{z}_{i+1}) \quad (10)$$

$$p_{\theta}(\mathbf{z}_i | \mathbf{z}_{i+1}) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_{p,i}(\mathbf{z}_{i+1}), \boldsymbol{\sigma}_{p,i}^2(\mathbf{z}_{i+1})), \quad p_{\theta}(\mathbf{z}_L) = \mathcal{N}(\mathbf{z}_L | \mathbf{0}, \mathbf{I}) \quad (11)$$

$$p_{\theta}(\mathbf{x} | \mathbf{z}_1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{p,0}(\mathbf{z}_1), \boldsymbol{\sigma}_{p,0}^2(\mathbf{z}_1)) \text{ or } P_{\theta}(\mathbf{x} | \mathbf{z}_1) = \mathcal{B}(\mathbf{x} | \boldsymbol{\mu}_{p,0}(\mathbf{z}_1)) \quad (12)$$

- ▶ p indicates generative model parameters, q indicates inference model parameters
- ▶ observation models is matching either continuous-valued (Gaussian \mathcal{N}) or binary-valued (Bernoulli \mathcal{B}) data
- ▶ The hierarchical specification allows the lower layers of the latent variables to be highly correlated

VAE inference model: bottom-up

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = q_{\phi}(\mathbf{z}_1|\mathbf{x}) \prod_{i=2}^L q_{\phi}(\mathbf{z}_i|\mathbf{z}_{i-1}) \quad (13)$$

$$q_{\phi}(\mathbf{z}_1|\mathbf{x}) = \mathcal{N}(\mathbf{z}_1|\boldsymbol{\mu}_{q,1}(\mathbf{x}), \boldsymbol{\sigma}_{q,1}^2(\mathbf{x})) \quad (14)$$

$$q_{\phi}(\mathbf{z}_i|\mathbf{z}_{i-1}) = \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_{q,i}(\mathbf{z}_{i-1}), \boldsymbol{\sigma}_{q,i}^2(\mathbf{z}_{i-1})), \quad i = 2 \dots L. \quad (15)$$

$$\mathbf{d}(\mathbf{y}) = \text{MLP}(\mathbf{y}) \quad (16)$$

$$\boldsymbol{\mu}(\mathbf{y}) = \text{Linear}(\mathbf{d}(\mathbf{y})) \quad (17)$$

$$\boldsymbol{\sigma}^2(\mathbf{y}) = \text{Softplus}(\text{Linear}(\mathbf{d}(\mathbf{y}))), \quad (18)$$

Ladder VAE

$$\mathbf{d}_n = \text{MLP}(\mathbf{d}_{n-1}) \quad (19)$$

$$\hat{\mu}_{q,i} = \text{Linear}(\mathbf{d}_i), i = 1 \dots L \quad (20)$$

$$\hat{\sigma}_{q,i}^2 = \text{Softplus}(\text{Linear}(\mathbf{d}_i)), i = 1 \dots L \quad (21)$$

where $\mathbf{d}_0 = \mathbf{x}$. Recursive downward pass:

$$q_\phi(\mathbf{z}|\mathbf{x}) = q_\phi(\mathbf{z}_L|\mathbf{x}) \prod_{i=1}^{L-1} q_\phi(\mathbf{z}_i|\mathbf{z}_{i+1}) \quad (22)$$

$$\sigma_{q,i} = \frac{1}{\hat{\sigma}_{q,i}^{-2} + \sigma_{p,i}^{-2}} \quad (23)$$

$$\mu_{q,i} = \frac{\hat{\mu}_{q,i} \hat{\sigma}_{q,i}^{-2} + \mu_{p,i} \sigma_{p,i}^{-2}}{\hat{\sigma}_{q,i}^{-2} + \sigma_{p,i}^{-2}} \quad (24)$$

$$q_\phi(\mathbf{z}_i|\cdot) = \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_{q,i}, \boldsymbol{\sigma}_{q,i}^2), \quad (25)$$

where $\mu_{q,L} = \hat{\mu}_{q,L}$ and $\sigma_{q,L}^2 = \hat{\sigma}_{q,L}^2$.

Ladder VAE

- ▶ precision-weighted combination of
 - ▶ $\hat{\mu}_q$ and $\hat{\sigma}_q^2$ carrying bottom-up information and
 - ▶ μ_p and σ_p^2 from the generative distribution carrying *top-down* prior information.
- ▶ $\hat{\mu}_q$ and $\hat{\sigma}_q^2$ as the approximate gaussian likelihood that is combined with a gaussian prior μ_p and σ_p^2 from the generative distribution.
- ▶ Together these form the approximate posterior distribution $q_\theta(\mathbf{z}|\mathbf{z}, \mathbf{x})$ using the same top-down dependency structure both in the inference and generative model.

Warm up from deterministic to VAE

- ▶ The variational regularization term causes some of the latent units to become inactive during training
- ▶ the approximate posterior for unit k , $q(z_{i,k}|\dots)$ is regularized towards its own prior $p(z_{i,k}|\dots)$, a phenomenon also recognized in the VAE setting
- ▶ presumably trapped in a local minima or saddle point at $KL(q_{i,k}|p_{i,k}) \approx 0$, with the optimization algorithm unable to re-activate them.

Warm up from deterministic to VAE

initializing training using the reconstruction error only (corresponding to training a standard deterministic auto-encoder), and then gradually introducing the variational regularization term:

$$\mathcal{L}(\theta, \phi; \mathbf{x})_T = -\beta KL(q_\phi(z|x) || p_\theta(\mathbf{z})) + E_{q_\phi(z|x)} [\log p_\theta(\mathbf{x}|\mathbf{z})], \quad (26)$$

where β is increased linearly from 0 to 1 during the first N_t epochs of training.

Experiments

- ▶ Datasets: MNIST, Omniglot, NORB
- ▶ $L = 5$ with sizes 64, 32, 16, 8, 4
- ▶ MNIST: bernoulli with sigmoid output layer

MNIST

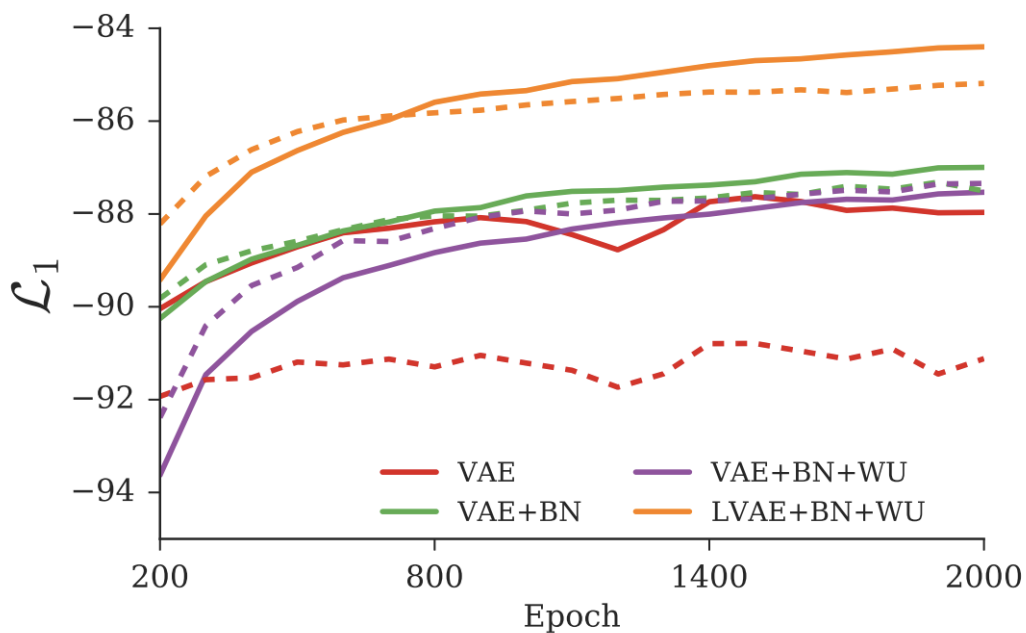


Figure 2: MNIST train (*full lines*) and test (*dashed lines*) set log-likelihood using one importance sample during training. The LVAE improves performance significantly over the regular VAE.

Experiments: MNIST

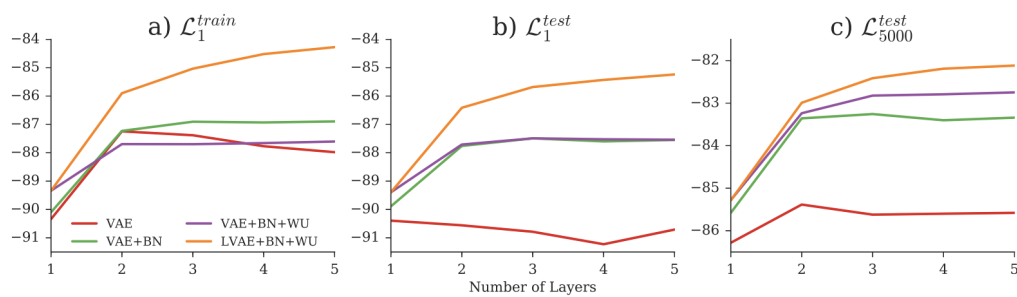
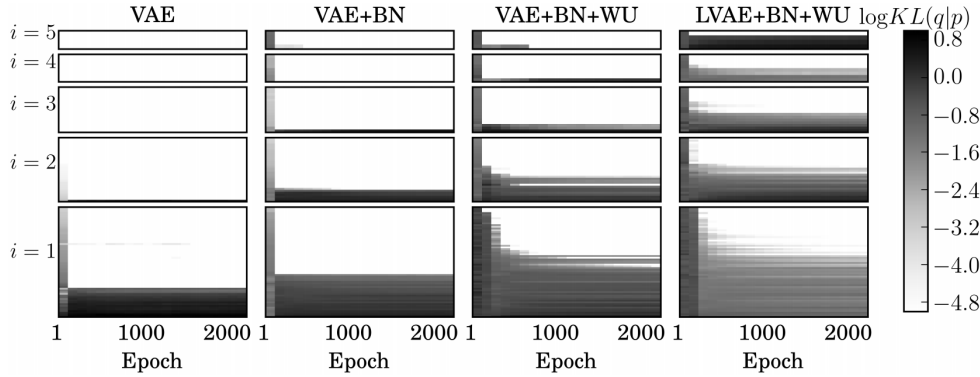


Figure 3: MNIST log-likelihood values for VAEs and the LVAE model with different number of latent layers

Experiments: MNIST



$\log KL(q|p)$ for each latent unit is shown at different training epochs. Low KL (white) corresponds to an inactive unit.

Experiments: MNIST

	$\leq \log p((x))$
VAE 1-layer + NF [17]	-85.10
IWAE, 2-layer + IW=1 [2]	-85.33
IWAE, 2-layer + IW=50 [2]	-82.90
VAE, 2-layer + VGP [20]	-81.90
LVAE, 5-layer	-82.12
LVAE, 5-layer + finetuning	-81.84
LVAE, 5-layer + finetuning + IW=10	-81.74

Experiments: MNIST

	VAE	VAE +BN	VAE +BN +WU	LVAE +BN +WU
OMNIGLOT				
64	-111.21	-105.62	-104.51	-
64-32	-110.58	-105.51	-102.61	-102.63
64-32-16	-111.26	-106.09	-102.52	-102.18
64-32-16-8	-111.58	-105.66	-102.66	-102.21
64-32-16-8-4	-110.46	-105.45	-102.48	-102.11
NORB				
64	2741	3198	3338	-
64-32	2792	3224	3483	3272
64-32-16	2786	3235	3492	3519
64-32-16-8	2689	3201	3482	3449
64-32-16-8-4	2654	3198	3422	3455

Experiments: MNIST

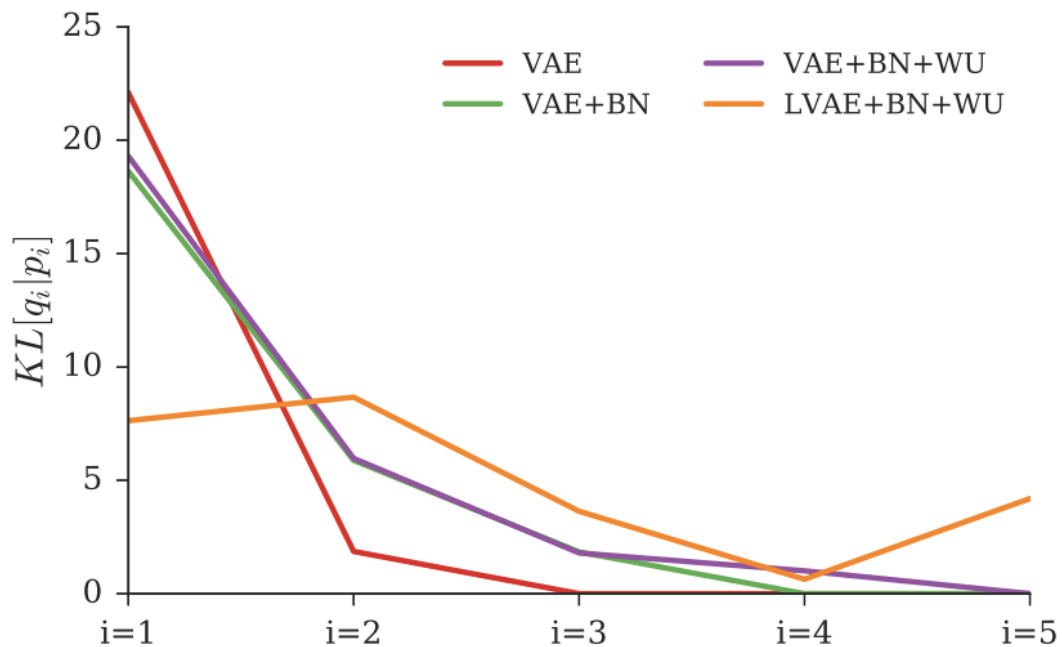


Figure 5: Layer-wise $KL[q|p]$ divergence going from the lowest to the highest layers. In the VAE models the KL divergence is highest in the lowest layers whereas it is more distributed in the LVAE model

Experiments: MNIST

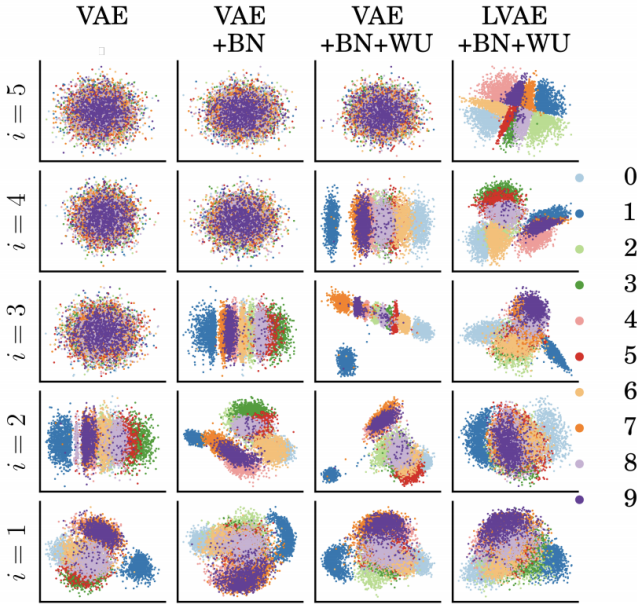


Figure 6: PCA-plots of samples from $q(z_i|z_{i-1})$ for 5-layer VAE and LVAE models trained on MNIST. Color-coded according to true class label

PCA plots of samples from $q(z_i|z_{i-1})$ for 5-layer VAE and LVAE models trained on MNIST

Conclusion

- ▶ new inference model for VAEs combining a bottom-up data-dependent approximate likelihood term with a prior information from the generative distribution
- ▶ learns a deeper and qualitatively different latent representation of data
- ▶ this parameterization makes the optimization easier since the inference is simply correcting the generative distribution instead of fitting the two models separately.

CausalVAE: Disentangled Representation Learning via Neural Structural Causal Models

Mengyue Yang, Furui Liu, Zhitang Chen, Xinwei Shen, Jianye
Hao, Jun Wang

Motivation

- ▶ learning disentangled representations
- ▶ assumption: the data is indeed generated by countable independent factors
- ▶ this paper: the independent factors are causally related



Figure 1: A swinging pendulum: an illustrative example

pendulum angle and light position are the causes of (l, x) of shadow

Advantages of causal disentanglement

- ▶ above are not independent, independence based disentanglement cant extract these factors
- ▶ will still disentangle light and pnedulum?
- ▶ generating counterfactual data– do operation
- ▶ example: $\text{do}(\textit{shadow} = 0)$

Structural Causal model

- ▶ a causal model is an ordered triple $\langle \epsilon, X, F \rangle$,
- ▶ where ϵ : exogenous variables whose values are determined by factors outside the model;
- ▶ X : set of endogenous variables whose values are determined by factors within the model;
- ▶ F : structural equations that express the value of each endogenous variable as a function of the values of the other variables in X and ϵ
- ▶ $x_i = f_i(x_{pa_i}, \epsilon_i)$

Causal VAE

- ▶ Encoder $\mathbf{x} \rightarrow \mathbf{z}$
- ▶ SCM Layer
- ▶ Decoder $\mathbf{z} \rightarrow \mathbf{x}$

Causal VAE: Structural Causal Model

- ▶ Step 1: identify the exogeneous factors ϵ
- ▶ step 2: A 'Causal Structure Layer' that relates the exogeneous factors
- ▶ *the causal structure is learnt, not prespecified*

Method: Causal VAE: Structural Causal Model Layer

- ▶ Consider a Linear Structural Causal Model
- ▶ disentangled factors correspond to \mathbf{u}
- ▶ and the causal graph among them corresponds to \mathbf{A}
- ▶ $\mathbf{z} = \mathbf{A}^T \mathbf{z} + \boldsymbol{\epsilon} = (\mathbf{I} - \mathbf{A}^T)^{-1} \boldsymbol{\epsilon}$
- ▶ $\mathbf{z} \in \mathbf{R}^n$ corresponding to n concepts
- ▶ $\boldsymbol{\epsilon} = N(0, \mathbf{I})$
- ▶ For example, chain graph $\mathbf{z}_1 \rightarrow \mathbf{z}_2 \rightarrow \mathbf{z}_3$

- ▶
$$\begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \mathbf{z}_3 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \\ \mathbf{z}_3 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \end{bmatrix}$$

Causal VAE: Step 1: \mathbf{z} latent representations



$$\mathbf{z} = \mathbf{A}^T \mathbf{z} + \boldsymbol{\epsilon} = (\mathbf{I} - \mathbf{A}^T)^{-1} \boldsymbol{\epsilon} \quad (27)$$

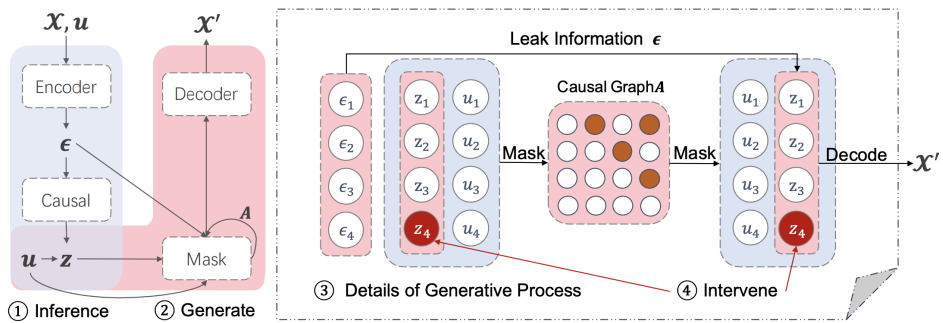
$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (28)$$

- ▶ ?? Only Markov Equivalence Graph if the variables are continuous, and the equations are linear with Gaussian (normal, or bell-shaped) errors.

Causal VAE: Step 2: Structural Causal Model Layer

- ▶ $\mathbf{z}_i = g_i(A_i \odot \mathbf{z}_i; \eta_i) + \epsilon_i$
- ▶ Mask layer that mimics generating children from parents

Causal VAE



Figure

Generative Model

- ▶ $\mathbf{x} \in \mathbf{R}^d$, known concepts $\mathbf{u} \in \mathbf{R}^n$, $\epsilon \in \mathbf{R}^n$
- ▶ $\mathbf{z} = (\mathbf{I} - \mathbf{A}^T)^{-1}\epsilon = \mathbf{C}\epsilon$
- ▶ Conventional VAE: $p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
- ▶ Conditional VAE: $p_\theta(\mathbf{x}, \mathbf{z}|\mathbf{u}) = p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{u})p(\mathbf{z}|\mathbf{u})$
- ▶ Causal VAE: $p_\theta(\mathbf{x}, \mathbf{z}, \epsilon|\mathbf{u}) = p_\theta(\mathbf{x}|\epsilon, \mathbf{z}, \mathbf{u})p(\mathbf{z}, \epsilon|\mathbf{u})$

Generative Model

- ▶ Decoder: $f(z)$
- ▶ Encoder: $h(x, u)$

$$p_{\theta}(x|z, \epsilon, u) = p_{\theta}(x|z) \equiv p_{\xi}(x - f(z)) \quad (29)$$

$$q_{\phi}(z, \epsilon|x, u) \equiv q(z|\epsilon)q_{\zeta}(\epsilon - h(x, u)), \quad (30)$$

$$x = f(z) + \xi, \quad \epsilon = h(x, u) + \zeta, \quad (31)$$

- ▶ ξ and ζ are the vectors of independent noise with probability densities p_{ξ} and q_{ζ} .
- ▶ When ξ and ζ are infinitesimal, the encoder and decoder can be regarded as deterministic ones.

Generative Model

- ▶
- ▶ the joint prior $p_{\theta}(\epsilon, z|u)$ for latent variables z and ϵ as

$$p_{\theta}(\epsilon, z|u) = p_{\epsilon}(\epsilon)p_{\theta}(z|u), \quad (32)$$

- ▶ $p_{\epsilon}(\epsilon) = \mathcal{N}(0, I)$
- ▶ the prior of latent endogenous variables $p_{\theta}(z|u)$ is a factorized Gaussian distribution conditioning on the additional observation u , i.e.

$$p_{\theta}(z|u) = \mathcal{N}(\lambda_1(u_i), \lambda_2^2(u_i)), \quad (33)$$

- ▶ λ_1 and λ_2 are an arbitrary functions. let $\lambda_1(u) = u$ and $\lambda_2(\mathbf{u}) \equiv 1$.

Training the Generative Model

$$\mathbb{E}_{q_{\mathcal{X}}}[\log p_{\theta}(\mathbf{x}|\mathbf{u})] \geq \text{ELBO} = \mathbb{E}_{q_{\mathcal{X}}}[\mathbb{E}_{\epsilon, \mathbf{z} \sim q_{\phi}}[\log p_{\theta}(\mathbf{x}|\mathbf{z}, \epsilon, \mathbf{u})] - \mathcal{D}(q_{\phi}(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u})||p_{\theta}(\epsilon, \mathbf{z}|\mathbf{u}))], \quad (34)$$

$$q_{\phi}(\epsilon, \mathbf{z}|\mathbf{x}, \mathbf{u}) = q_{\phi}(\epsilon|\mathbf{x}, \mathbf{u})\delta(\mathbf{z} = \mathbf{C}\epsilon) = q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{u})\delta(\epsilon = \mathbf{C}^{-1}\mathbf{z}), \quad (35)$$

$$\text{ELBO} = \mathbb{E}_{q_{\mathcal{X}}}[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{u})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \quad (36)$$

$$\mathcal{D}(q_{\phi}(\epsilon|\mathbf{x}, \mathbf{u})||p_{\epsilon}(\epsilon)) - \mathcal{D}(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{u})||p_{\theta}(\mathbf{z}|\mathbf{u}))]. \quad (37)$$

Additional Constraints

$$l_u = \mathbb{E}_{q_{\mathcal{X}}} \|\mathbf{u} - \sigma(\mathbf{A}^T \mathbf{u})\|_2^2 \leq \kappa_1, \quad (38)$$

$$l_m = \mathbb{E}_{\mathbf{z} \sim q_{\phi}} \sum_{i=1}^n \|z_i - g_i(\mathbf{A}_i \circ \mathbf{z}; \boldsymbol{\eta}_i)\|^2 \leq \kappa_2, \quad (39)$$

DAG constraint:

$$H(\mathbf{A}) \equiv \text{tr}((\mathbf{I} + \mathbf{A} \circ \mathbf{A})^n) - n = 0. \quad (40)$$

Training Loss Function:

$$\mathcal{L} = -\text{ELBO} + \alpha H(\mathbf{A}) + \beta l_u + \gamma l_m, \quad (41)$$

Experiments

- ▶ Two datasets: Synthetic (Pendulum) and CelebA
- ▶ Pendulum:
 - ▶ 3 entities (pendulum, light, shadow)
 - ▶ 4 concepts ((pendulum angle, light angle) → (shadow location, shadow length)).
- ▶ CelebA
 - ▶ 4 causally related concepts (gender, smile, eyes open, mouth open), where gender and smile cause eyes open, and smile causes mouth open.
- ▶ Evaluation Criteria:
 - ▶ MIC: Maximal Information Criterion (MIC)
 - ▶ Total Information Criterion (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

Intervention Experiments

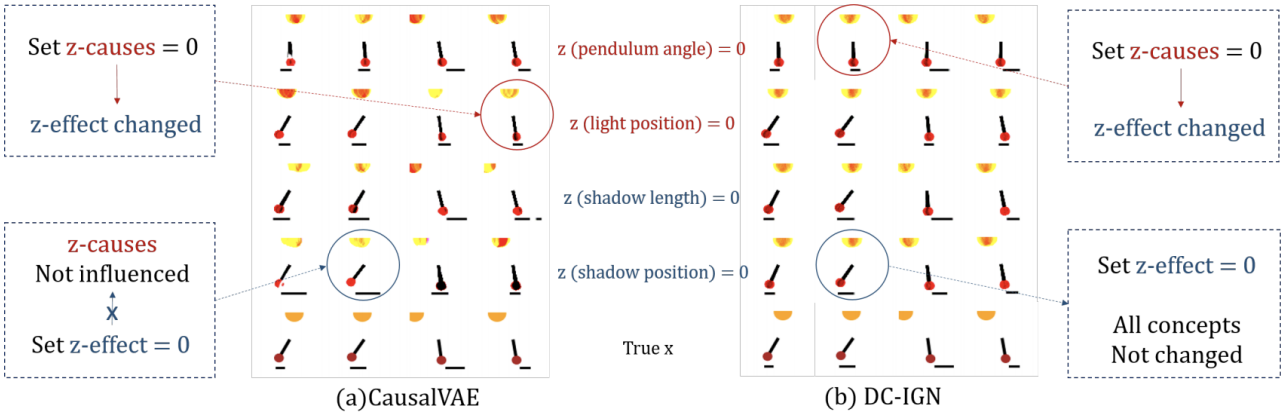


Figure 3: The results of Intervention experiments on the pendulum dataset. Each row shows the result of controlling the PENDULUM ANGLE, LIGHT ANGLE, SHADOW LENGTH, and SHADOW LOCATION respectively. The bottom row is the original input image.

Intervention Experiments

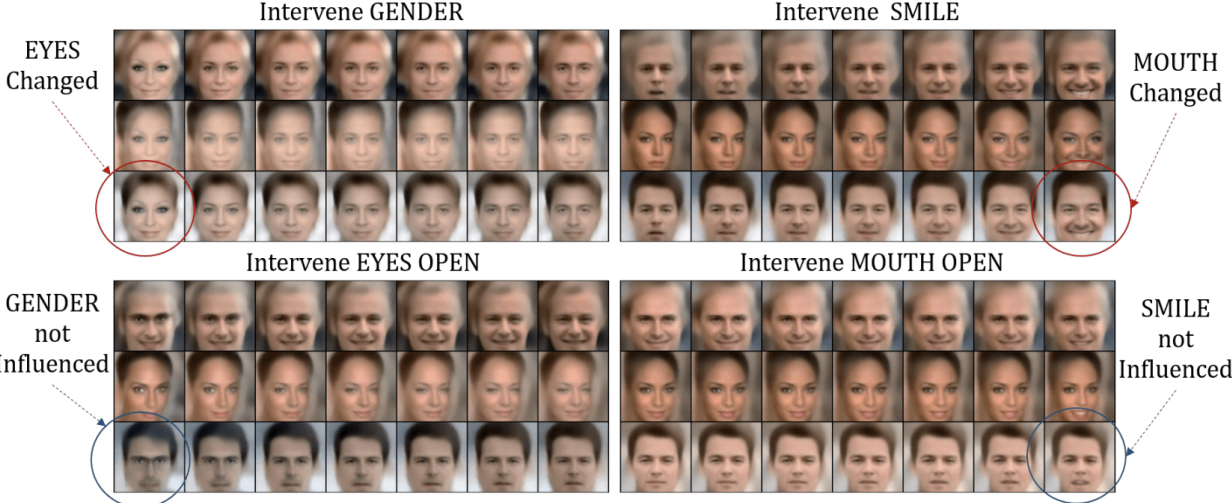


Figure 4: Results of CausalVAE model on CelebA. The controlled factors are GENDER, SMILE, EYES OPEN and MOUTH OPEN respectively.

Conclusion

- ▶ learning disentangled representations of causally related concepts in data
- ▶ allows intervention to generate counterfactual outputs as expected according to our understanding of the causal system.