

VAE는 MLL loss를 복잡하게 만들었지만, autoregressive 보다 성능이 좋지는 않다.

IZIIC flow-based는 high-dim으로 확장된다.

IZIIC는 hierarchical로 complex covariance structure model로 VAE의 성능 개선을 줄인다.

↳ skip-connection과 bidirectional stochastic을 가진 BIVA 소거.

↳ semantic latent

* Introduction.

Implicit(GAN)과 explicit(MLL) model과 explicit 선택

Autoregressive는 성능은 좋지만, run-time이 좋지가 않다.

flow-based는 자연스럽지만, 그동안 널리에서 autoregressive 보다 잘.

VAE는 posterior를 근사화하는데 low-dim manifold에서 blurred poor MLL로 모시됨.

flow로 covariance structure를 보완하거나, stacked 가우스나지만, 여전히 성능이 낮음

논문에서는 충분한 representation의 latent에 autoregressive의 장점들이 있다.

latent의 inactive를 활성화시키고, 정렬의 흐름을 확장시키기 위해, skip-connection을 사용하는 BIVA 소거.

flexible posterior approximation을 위해, bottom-up과 up-bottom의 양방향 inference를 구성

Contribution

1. ablation

2. 성능 개선

3. semi-supervised 가능

4. 이상 탐지 가능

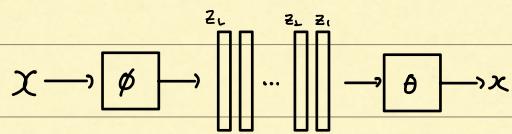
※ Variational Autoencoder.

Hierarchy 구조는 보통

$$p_\theta(x, z) = p_\theta(x|z_1)p_\theta(z_L) \prod_{i=1}^{L-1} p_\theta(z_i|z_{i+1})$$

posterior approximation은 bottom-up으로 이루어짐. ($1 \rightarrow L$)

$$q_\phi(z|x) = q_\phi(z_1|x) \prod_{i=1}^{L-1} q_\phi(z_{i+1}|z_i)$$



$$\log p_\theta(x) \geq \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] \equiv \mathcal{L}(\theta, \phi).$$

Hierarchical은 더 expressive 표현, Latent variables가太多的 collapsed 일어나는 경향이 있다.

Ladder VAE는 G의 bottom-up과 inference top-down을 공유한다.

$$q_{\phi,\theta}(z|x) = q_\phi(z_L|x) \prod_{i=1}^{L-1} q_{\phi,\theta}(z_i|z_{i+1}, x)$$

LVAE에서는 bottom-up(?) 모양이 모든

z가 x의 영향을 받고, 이를 통해 collapse되는

일 일어난다.

하지만 너무 깊으면 여전히 일어남.

∴ BIVAE는 top-bottom path를 추가하고,

Bottom-up stochastic inference를 구현하도록

각 z_i간 latent의 factorization을 정의.

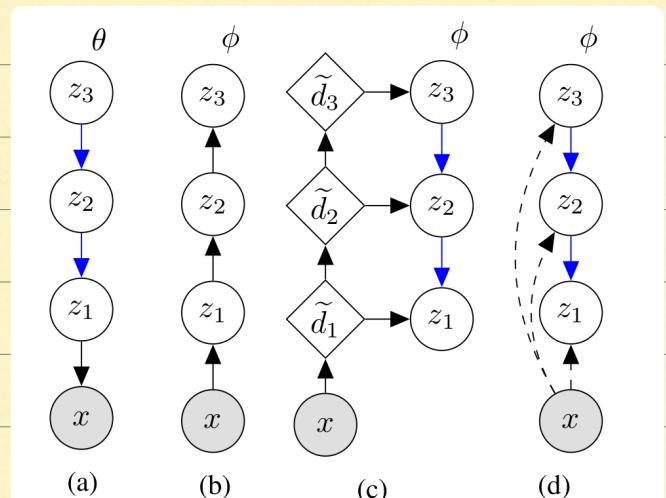


Figure 5: (a) Generative model of a VAE/LVAE with $L = 3$ stochastic variables, (b) VAE inference model, (c) LVAE inference model, and (d) skip connections among stochastic variables in the LVAE where dashed lines denote a skip-connection. Blue arrows indicate that there are shared parameters between the inference and generative model.

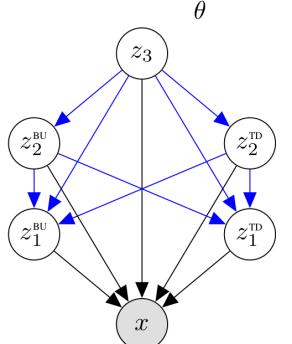
* Bidirectional - Inference VAE.

* Model Architecture.

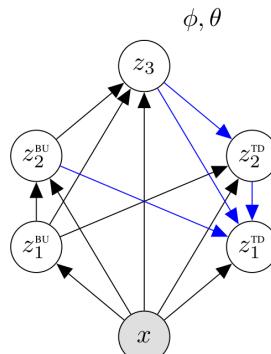
- Generative model

$BU = \text{Bottom-up}$, $TD = \text{Top-down}$.

latent \mathbf{z}_i 를 $\mathbf{z}_i = [z_i^{BU}, z_i^{TD}]$ 로 분해



(a) Generative model



(b) Inference model

Figure 1: A $L = 3$ layered BIVA with (a) the generative model and (b) inference model. Blue arrows indicate that the deterministic parameters are shared between the inference and generative models. See Appendix B for a detailed explanation and a graphical model that includes the deterministic variables.

\mathbf{z}_i 를 input으로 했을 때 $\mathcal{N}(\mu, \sigma^2)$ parameterization 된 deterministic TD가 더 있다.

↳ \mathbf{z}_i 를 포함하는 latent

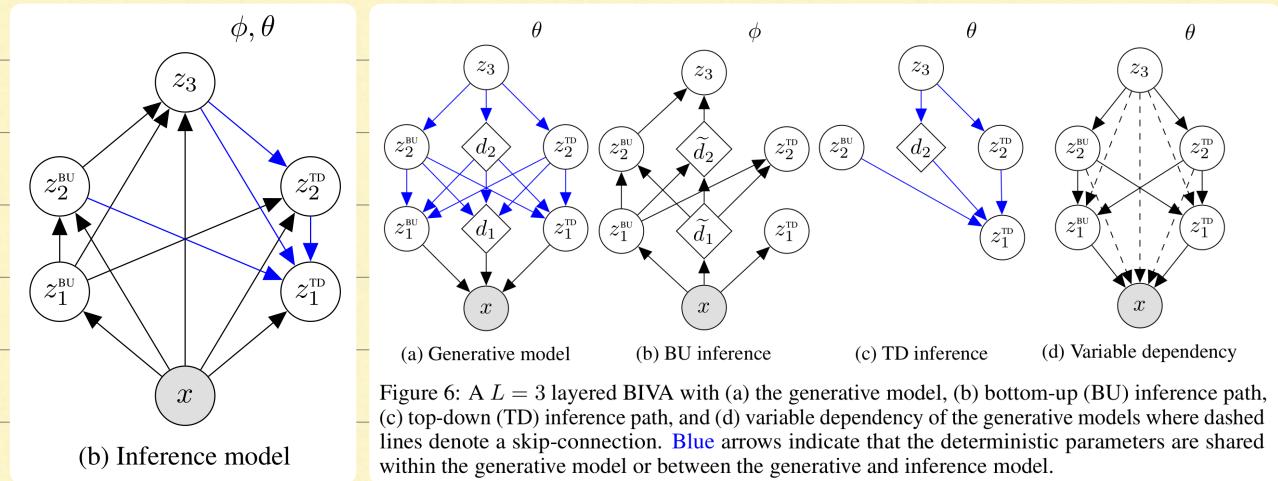
$\mathbf{z}_i = [z_i^{BU}, z_i^{TD}]$ 는 concat,

$$p_{\theta}(x, \mathbf{z}) = p_{\theta}(x|\mathbf{z})p_{\theta}(\mathbf{z}_L) \prod_{i=1}^{L-1} p_{\theta}(z_i^{BU}|z_{>i})p_{\theta}(z_i^{TD}|z_{>i}),$$

$p_{\theta}(x, \mathbf{z})$ 는 \mathbf{z}_i 와 각각으로 dependent이고, \mathbf{z}_i 가 $\mathbf{z}_{>i}$ 에 dependent이다.

즉 \mathbf{z}_i ($i \sim L$)는 μ 와 σ 의 대신 각각의 선형합의 mean, var parameterization.

- Bidirectional inference network.



BU는 모든 bottom의 차원 영향을 받는다.

z_i^{bu} 는 $z_{<i}^{\text{bu}}$ 에 걸쳐적으로 dependency, $z_{<i}^{\text{bu}}$ 는 d_i 에 서로 dependent이다.

z_i^{td} 는 $z_{<i}^{\text{td}}$ 에 $z_{>i}^{\text{td}}$. 그리고 data의 영향.

그리고 TD path는 Get parameter sharing.

$$q_\phi(\mathbf{z}|\mathbf{x}) = q_\phi(z_L|\mathbf{x}, z_{<L}^{\text{BU}}) \prod_{i=1}^{L-1} q_\phi(z_i^{\text{BU}}|\mathbf{x}, z_{<i}^{\text{BU}}) q_{\phi,\theta}(z_i^{\text{TD}}|\mathbf{x}, z_{<i}^{\text{BU}}, z_{>i}^{\text{BU}}, z_{>i}^{\text{TD}}).$$

↳ ELBO의 경우 Training.

⇒ Motivation

- Deterministic top-down path.

Skip connection은 간단하지만, deep MV가 효과적이다.

(CLSTM VAE)→ deep hierarchy의 처리성이 좋기 때문에 deterministic TD path를 생각한다.

But latent → skip connection 추가, latent → 높은 hierarchy의 dependent 하도록 함.

↳ 이는 information flow, collapse 방지

- Bidirectional Inference

Auxiliary VAE(AVAE)와 영감을 받음.

↳ 높은 레벨 차를 추적하고, covariance structure 추가.

BIVAE는 latent를 BU, TD로 분리하여 통합한 효과

BU는 모든 bottom의 차원 factorize SDI 때문

TD는 복잡한 cov를 학습할 수 있도록 하며, high-level semantic feature의 유통

⇒ Anomaly detection with BIVAE

BIVAE의 high-hierarchy는 high-level semantics을 capture하는 데 좋다고 주장.

ELBO는 alternative log likelihood lower bound로 claim.

↳ 예측을 영역으로 사용.

$$\mathcal{L}^{>k} = \mathbb{E}_{p_\theta(z_{\leq k}|z_{>k})q_\phi(z_{>k}|x)} \left[\log \frac{p_\theta(x|\mathbf{z})p_\theta(z_{>k})}{q_\phi(z_{>k}|x)} \right]$$

* Experiments

▲ Ablation.

Table 1: A comparison of the LVAE with no skip-connections and no bottom-up inference, the LVAE+ with skip-connections and no bottom-up inference, and BIVA. All models are trained on the CIFAR-10 dataset.

	PARAM.	BITS/DIM
LVAE $L=15, \mathcal{L}_1$	72.36M	≤ 3.60
LVAE+ $L=15, \mathcal{L}_1$	73.35M	≤ 3.41
LVAE+ $L=29, \mathcal{L}_1$	119.71M	≤ 3.45
BIVA $L=15, \mathcal{L}_1$	102.95M	≤ 3.12

LVAE+는 BIVA보다 BUD가 낮다.

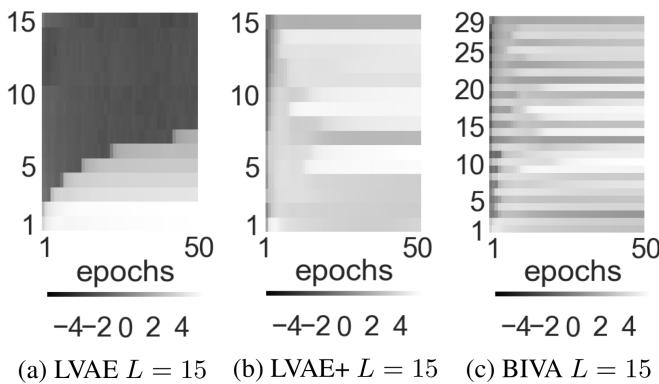


Figure 2: The $\log KL(q||p)$ for each stochastic latent variable as a function of the training epochs on CIFAR-10. (a) is a $L = N = 15$ stochastic latent layer LVAE with no skip-connections and no bottom-up inference. (b) is a $L = N = 15$ LVAE+ with skip-connections and no bottom-up inference. (c) is a $L = 15$ stochastic latent layer ($N = 29$ latent variables) BIVA for which $1, 2, \dots, N$ denotes the stochastic latent variables following the order $z_1^{\text{BU}}, z_1^{\text{TD}}, z_2^{\text{BU}}, z_2^{\text{TD}}, \dots, z_L$.

LVAE는 BUD가 높은 편이다.

∴ BUDA latent 변수의 더 강한 연결을 갖는다 (BUD)

▲ Binary Image.

Table 2: Test log-likelihood on statically binarized MNIST for different number of importance weighted samples. The finetuned models are trained for an additional number of epochs with no free bits, $\lambda = 0$. For testing resiliency we trained 4 models and evaluated the standard deviations to be ± 0.031 for \mathcal{L}_1 .

	$-\log p(x)$
<i>With autoregressive components</i>	
PIXELCNN [57]	= 81.30
DRAW [13]	< 80.97
IAFVAE [23]	≤ 79.88
PIXELVAE [14]	≤ 79.66
PIXELRNN [57]	= 79.20
VLAЕ [5]	≤ 79.03
<i>Without autoregressive components</i>	
DISCRETE VAE [42]	≤ 81.01
BIVA, \mathcal{L}_1	≤ 81.20
BIVA, \mathcal{L}_{1e3}	≤ 78.67
BIVA FINETUNED, \mathcal{L}_1	≤ 80.47
BIVA FINETUNED, \mathcal{L}_{1e3}	≤ 78.59

Table 3: Semi-supervised test error for BIVA on MNIST for 100 randomly chosen and evenly distributed labelled samples.

	ERROR %
M1+M2 [22]	3.33% (± 0.14)
VAT [32]	2.12%
CATGAN [51]	1.91% (± 0.10)
SDGM [31]	1.32% (± 0.07)
LADDERNET [38]	1.06% (± 0.37)
ADGM [31]	0.96% (± 0.02)
IMPGAN [44]	0.93% (± 0.07)
TRIPLEGAN [29]	0.91% (± 0.58)
SSLGAN [6]	0.80% (± 0.10)
BIVA	0.83% (± 0.02)

* Natural Image.

Table 4: Test log-likelihood on CIFAR-10 for different number of importance weighted samples. We evaluated two different BIVA with various number of layers (L). For testing resiliency we trained 3 models and evaluated the standard deviations to be ± 0.013 for \mathcal{L}_1 and $L = 15$.
BITS/DIM

With autoregressive components	
CONVDRAW [12]	< 3.58
IAFVAE \mathcal{L}_1 [23]	≤ 3.15
IAFVAE \mathcal{L}_{1e3} [23]	≤ 3.12
GATEDPIXELCNN [56]	= 3.03
PIXELRNN [57]	= 3.00
VLAE [5]	≤ 2.95
PIXELCNN++ [45]	= 2.92
Without autoregressive components	
NICE [8]	= 4.48
DEEPGMMS [58]	= 4.00
REALNVP [9]	= 3.49
DISCRETEVAE++ [54]	≤ 3.38
GLOW [21]	= 3.35
FLOW++ [16]	= 3.08
BIVA L=10, \mathcal{L}_1	≤ 3.17
BIVA L=15, \mathcal{L}_1	≤ 3.12
BIVA L=15, \mathcal{L}_{1e3}	≤ 3.08

ImageNet 을 훈련한 모델의 성능

Autoregressive 확장 가능



Figure 3: (left) images from the CelebA dataset preprocessed to 64x64 following [27]. (right) $\mathcal{N}(0, I)$ generations of BIVA with $L = 20$ layers that achieves a $\mathcal{L}_1 = 2.48$ bits/dim on the test set.

나 훈련하지 않은 이미지 생성

* Conclusion.

BIVA를 hierarchical, semi-supervised로 확장한 성능