

Deep Generative Models

Lecture 7

Roman Isachenko

Ozon Masters

2021

VAE limitations

- ▶ Poor variational posterior distribution (encoder)

$$q(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})).$$

- ▶ Poor prior distribution

$$p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I}).$$

- ▶ Poor probabilistic model (decoder)

$$p(\mathbf{x}|\mathbf{z}, \theta) = \mathcal{N}(\mathbf{x}|\mu_\theta(\mathbf{z}), \sigma_\theta^2(\mathbf{z})).$$

- ▶ Loose lower bound

$$p(\mathbf{x}|\theta) - \mathcal{L}(q, \theta) = (?).$$

VAE prior

ELBO revisiting

$$\mathcal{L}(q, \theta) = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{q(\mathbf{z}_i|\mathbf{x}_i)} \log p(\mathbf{x}_i|\mathbf{z}_i, \theta)}_{\text{Reconstruction loss}} - \underbrace{(\log n - \mathbb{E}_{q(\mathbf{z})} \mathbb{H}[q(i|\mathbf{z})])}_{0 \leq \text{Mutual info} \leq \log N} - \underbrace{KL(q(\mathbf{z})||p(\mathbf{z}))}_{\text{Marginal KL}}$$

How to choose the optimal $p(\mathbf{z})$?

- ▶ SG: $p(\mathbf{z}) = \mathcal{N}(0, I)$ \Rightarrow over-regularization;
- ▶ MoG: $p(\mathbf{z}|\boldsymbol{\lambda}) = \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k^2)$ \Rightarrow (*), (**);
- ▶ $p(\mathbf{z}) = q(\mathbf{z}) = \frac{1}{n} \sum_{i=1}^n q(\mathbf{z}|\mathbf{x}_i)$ \Rightarrow overfitting and highly expensive.

(*) <https://arxiv.org/abs/1611.02648>

(**) <https://pdfs.semanticscholar.org/f6fe/5e8e25994c188ba6a124462e2cc55f2c5a67.pdf>

Variational Mixture of posteriors

$$p(\mathbf{z}|\boldsymbol{\lambda}) = \frac{1}{K} \sum_{k=1}^K q(\mathbf{z}|\mathbf{u}_k),$$

where $\boldsymbol{\lambda} = \{\mathbf{u}_1, \dots, \mathbf{u}_K\}$ are trainable pseudo-inputs.

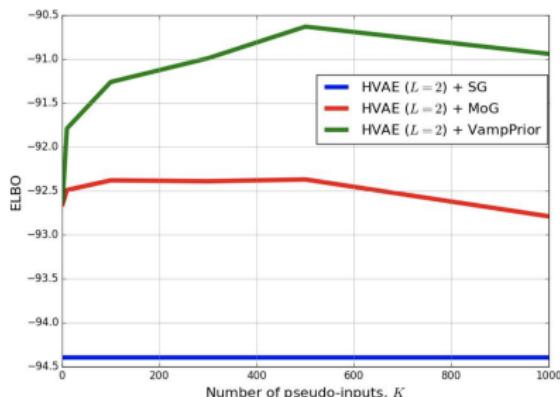
- ▶ Multimodal \Rightarrow prevents over-regularization;.
- ▶ $K \ll n \Rightarrow$ prevents from potential overfitting + less expensive to train.
- ▶ Pseudo-inputs are prior hyperparameters \Rightarrow connection to the Empirical Bayes.

VampPrior, 2017

Do we equally need the multimodal prior?

Is it beneficial to couple the prior with the variational posterior or MoG is enough?

MODEL	LL
VAE ($L = 1$) + NF [32]	-85.10
VAE ($L = 2$) [6]	-87.86
IWAE ($L = 2$) [6]	-85.32
HVAE ($L = 2$) + SG	-85.89
HVAE ($L = 2$) + MoG	-85.07
HVAE ($L = 2$) + VAMPRIOR <i>data</i>	-85.71
HVAE ($L = 2$) + VAMPRIOR	-83.19
AVB + AC ($L = 1$) [28]	-80.20
VLAЕ [7]	-79.03
VAE + IAF [18]	-79.88
CONVHVAE ($L = 2$) + VAMPRIOR	-81.09
PIXELHVAE ($L = 2$) + VAMPRIOR	-79.78

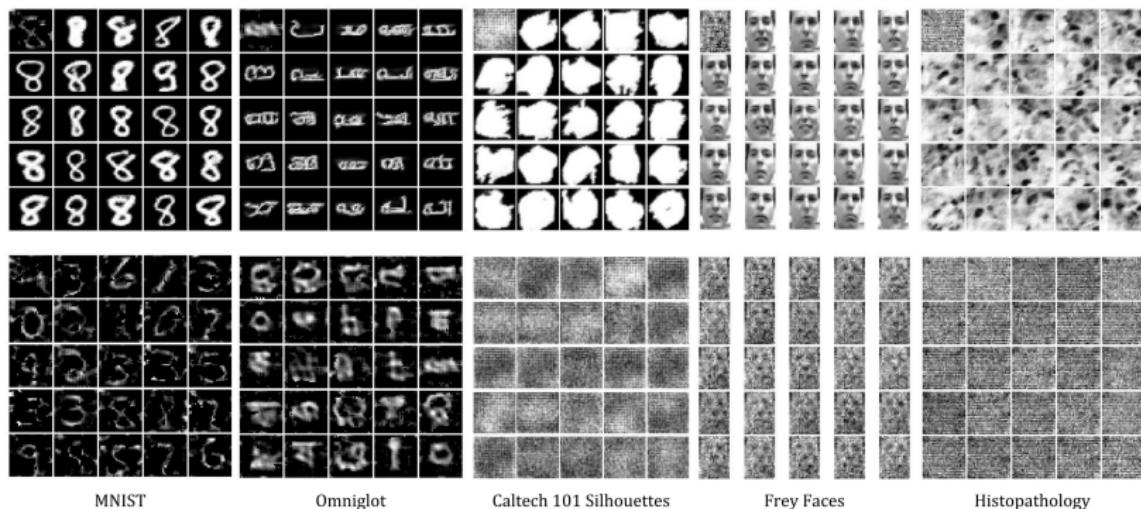


<https://arxiv.org/abs/1705.07120>

VampPrior, 2017

Top row: images generated by PixelHVAE + VampPrior for chosen pseudo-input in the left top corner.

Bottom row: Images represent a subset of trained pseudo-inputs for different datasets.



VAE limitations

- ▶ Poor variational posterior distribution (encoder)

$$q(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})).$$

- ▶ Poor prior distribution

$$p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I}).$$

- ▶ Poor probabilistic model (decoder)

$$p(\mathbf{x}|\mathbf{z}, \theta) = \mathcal{N}(\mathbf{x}|\mu_\theta(\mathbf{z}), \sigma_\theta^2(\mathbf{z})).$$

- ▶ Loose lower bound

$$p(\mathbf{x}|\theta) - \mathcal{L}(q, \theta) = (?).$$

Posterior collapse: toy example

Let define latent variable model in the following way:

$$p(\mathbf{x}|\boldsymbol{\theta}) = \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z}$$

- ▶ prior distribution $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|0, \mathbf{I})$;
- ▶ probabilistic model $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}), \boldsymbol{\sigma}_{\boldsymbol{\theta}}(\mathbf{z}))$ (diagonal covariance);
- ▶ variational posterior $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\phi}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{\boldsymbol{\phi}}(\mathbf{x}), \boldsymbol{\sigma}_{\boldsymbol{\phi}}(\mathbf{x}))$ (diagonal covariance).

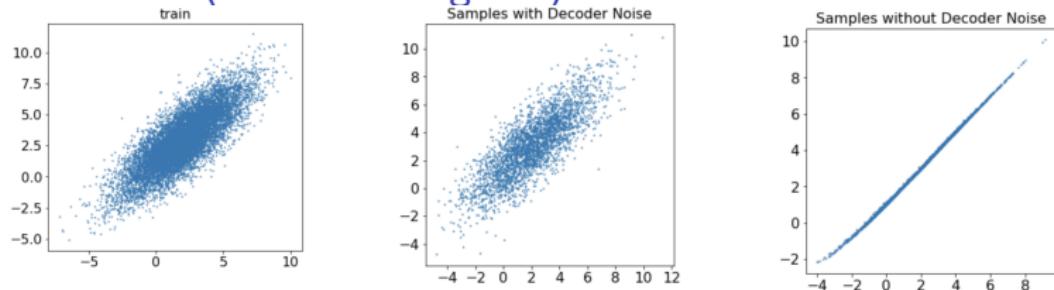
Let data distribution is $\pi(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Possible cases:

- ▶ covariance matrix $\boldsymbol{\Sigma}$ is diagonal (univariate case);
- ▶ covariance matrix $\boldsymbol{\Sigma}$ is **not** diagonal (multivariate case).

What is the difference?

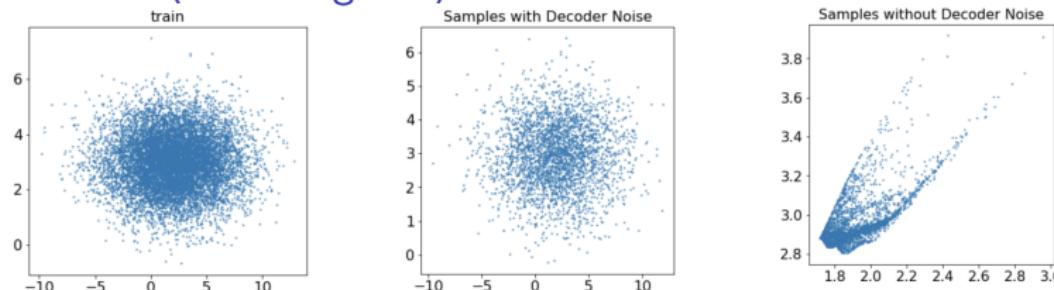
Posterior collapse: toy example

Multivariate (Σ is non-diagonal)



The encoder uses latent variables to model data.

Univariate (Σ is diagonal)



Latent variables are not used, since the decoder could model the data without the encoder.

Posterior collapse

Representation learning

"Identifies and disentangles the underlying causal factors of the data, so that it becomes easier to understand the data, to classify it, or to perform other tasks".

$$p(\mathbf{x}|\theta) = \int p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z}, \theta)p(\mathbf{z})d\mathbf{z}$$

If the decoder model $p(\mathbf{x}|\mathbf{z}, \theta)$ is powerful enough to model $p(\mathbf{x}|\theta)$ the latent variables \mathbf{z} becomes irrelevant.

$$\mathcal{L}(q, \theta) = \frac{1}{n} \sum_{i=1}^n [\mathbb{E}_{q(\mathbf{z}_i|\mathbf{x}_i)} \log p(\mathbf{x}_i|\mathbf{z}_i, \theta) - KL(q(\mathbf{z}_i|\mathbf{x}_i)||p(\mathbf{z}_i))]$$

Early in the training the approximate posterior $q(\mathbf{z}|\mathbf{x})$ carries little information about \mathbf{x} and the model sets the posterior to the prior to avoid paying any cost $KL(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$.

PixelVAE, 2016

$$p(\mathbf{x}|\theta) = \int p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z}, \theta)p(\mathbf{z})d\mathbf{z}$$

- ▶ More powerful $p(\mathbf{x}|\mathbf{z}, \theta)$ leads to more powerful generative model $p(\mathbf{x}|\theta)$.
- ▶ Too powerful $p(\mathbf{x}|\mathbf{z}, \theta)$ could lead to posterior collapse, where variational posterior $q(\mathbf{z}|\mathbf{x})$ will not carry any information about data and close to prior $p(\mathbf{z})$.

How to make the generative model $p(\mathbf{x}|\mathbf{z}, \theta)$ more powerful?

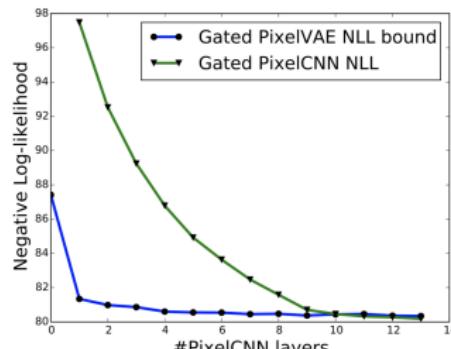
Autoregressive decoder

$$p(\mathbf{x}|\mathbf{z}, \theta) = \prod_{i=1}^n p(x_i|\mathbf{x}_{1:i-1}, \mathbf{z}, \theta)$$

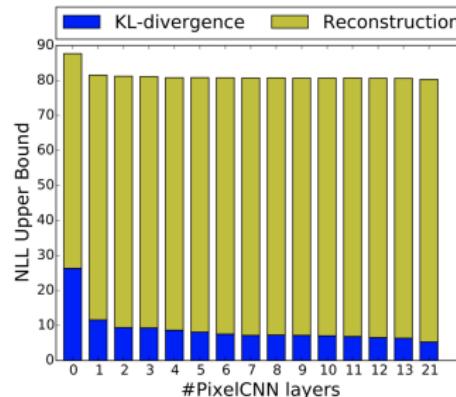
PixelVAE, 2016

VAE model with autoregressive PixelCNN decoder with few autoregressive layers.

- ▶ Global structure is captured by latent variables.
- ▶ Local statistics are captured by limited receptive field autoregressive model.



(a)



(b)

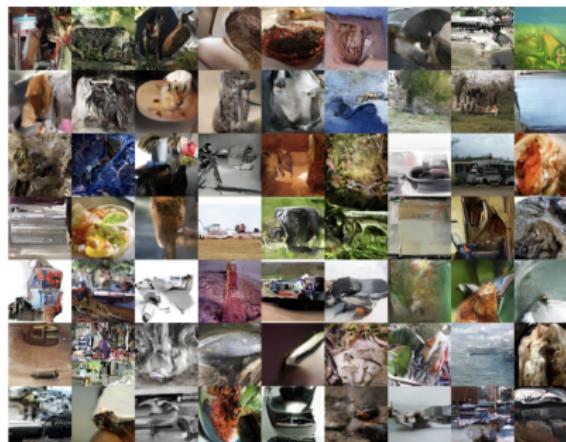
PixelVAE, 2016

MNIST

Model	NLL Test
DRAW (Gregor et al., 2016)	≤ 80.97
Discrete VAE (Rollef, 2016)	$= 81.01$
IAF VAE (Kingma et al., 2016)	≈ 79.88
PixelCNN (van den Oord et al., 2016a)	$= 81.30$
PixelRNN (van den Oord et al., 2016a)	$= 79.20$
Convolutional VAE	≤ 87.41
PixelVAE	≤ 80.64
Gated PixelCNN (our implementation)	$= 80.10$
Gated PixelVAE	$\approx 79.48 (\leq 80.02)$
Gated PixelVAE without upsampling	$\approx \mathbf{79.02} (\leq 79.66)$

ImageNet 64x64

Model	NLL Validation (Train)
Convolutional DRAW (Gregor et al., 2016)	$\leq 4.10 (4.04)$
Real NVP (Dinh et al., 2016)	$= 4.01 (3.93)$
PixelRNN (van den Oord et al., 2016a)	$= 3.63 (3.57)$
Gated PixelCNN (van den Oord et al., 2016b)	$= \mathbf{3.57} (3.48)$
Hierarchical PixelVAE	$\leq 3.66 (3.59)$



Decoder weakening

$$p(\mathbf{x}|\boldsymbol{\theta}) = \int p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z}$$

Powerful decoder $p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})$ makes the model expressive, but posterior collapse is possible.

PixelVAE model uses the autoregressive PixelCNN model with small number of layers to limit receptive field.

How to force the model encode information about \mathbf{x} into \mathbf{z} ?

$$\mathcal{L}(q, \boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \log p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta}) - \beta \cdot KL(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

What we get if $\beta = 1$ ($\beta = 0$)?

KL annealing

- ▶ Start training with $\beta = 0$.
- ▶ Increase it until $\beta = 1$ during training process.

Decoder weakening

Free bits

- ▶ Divide the latent dimensions into the K subsets.
- ▶ Ensure that using less than λ nats of information per subset j .

$$\hat{\mathcal{L}}(q, \theta) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \log p(\mathbf{X}|\mathbf{Z}, \theta) - \sum_{j=1}^K \max(\lambda, KL(q(\mathbf{Z}_j|\mathbf{X})||p(\mathbf{Z}_j))).$$

Increasing the latent information is advantageous for the reconstruction term.

This results in $KL(q(\mathbf{Z}_j|\mathbf{x})||p(\mathbf{Z}_j)) \geq \lambda$ for all j , in practice.

Variational Lossy AutoEncoder, 2016

Lossy code via explicit information placement

$$p(\mathbf{x}|\mathbf{z}, \theta) = \prod_{i=1}^m p(x_i|\mathbf{z}, \mathbf{x}_{\text{WindowAround}(i)}, \theta).$$

- ▶ $\text{WindowAround}(i)$ restricts the receptive field (it forbids to represent arbitrarily complex distribution over \mathbf{x} without dependence on \mathbf{z}).
- ▶ Local statistics of 2D images (texture) will be modeled by a small local window.
- ▶ Global structural information (shapes) is long-range dependency that can only be communicated through latent code \mathbf{z} .

Variational Lossy AutoEncoder, 2016

ELBO revisiting

$$\mathcal{L}(q, \theta) = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{q(\mathbf{z}_i | \mathbf{x}_i)} \log p(\mathbf{x}_i | \mathbf{z}_i, \theta)}_{\text{Reconstruction loss}} - \underbrace{(\log n - \mathbb{E}_{q(\mathbf{z})} \mathbb{H}[q(i|\mathbf{z})])}_{0 \leq \text{Mutual info} \leq \log N} - \underbrace{KL(q(\mathbf{z}) || p(\mathbf{z}))}_{\text{Marginal KL}}$$

VampPrior

$$p(\mathbf{z}|\lambda) = \frac{1}{K} \sum_{k=1}^K q(\mathbf{z}|\mathbf{u}_k),$$

where $\mathbf{u}_1, \dots, \mathbf{u}_K$ is trainable preudo-inputs.

Autoregressive flow prior

$$\log p(\mathbf{z}|\lambda) = \log p(\epsilon) + \log \det \left| \frac{d\epsilon}{d\mathbf{z}} \right|$$

$$\mathbf{z} = g(\epsilon, \lambda) = f^{-1}(\epsilon, \lambda)$$

Variational Lossy AutoEncoder, 2016

Theorem

VAE with AF prior in latent code \mathbf{z} is equivalent to using IAF posterior for latent code ϵ .

Proof

$$\begin{aligned}\mathcal{L}(q, \theta) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z}, \theta) + \log p(\mathbf{z}, \lambda) - q(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}, \theta) + \underbrace{\left(\log p(f(\mathbf{z}, \lambda)) + \log \left| \det \frac{\partial f(\mathbf{z}, \lambda)}{\partial \mathbf{z}} \right| \right)}_{\text{AF prior}} - q(\mathbf{z}|\mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}, \theta) + \log p(f(\mathbf{z}, \lambda)) - \underbrace{\left(q(\mathbf{z}|\mathbf{x}) - \log \left| \det \frac{\partial f(\mathbf{z}, \lambda)}{\partial \mathbf{z}} \right| \right)}_{\text{IAF posterior}} \right]\end{aligned}$$

Variational Lossy AutoEncoder, 2016

Autoregressive flow prior

$$\begin{aligned}\mathcal{L}(q, \theta) &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}, \theta) + \underbrace{\left(\log p(f(\mathbf{z}, \lambda)) + \log \left| \det \frac{\partial f(\mathbf{z}, \lambda)}{\partial \mathbf{z}} \right| \right)}_{\text{AF prior}} - q(\mathbf{z}|\mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left[\log p(\mathbf{x}|\mathbf{z}, \theta) + \log p(f(\mathbf{z}, \lambda)) - \underbrace{\left(q(\mathbf{z}|\mathbf{x}) - \log \left| \det \frac{\partial f(\mathbf{z}, \lambda)}{\partial \mathbf{z}} \right| \right)}_{\text{IAF posterior}} \right]\end{aligned}$$

- ▶ IAF posterior decoder path: $p(\mathbf{x}|\mathbf{z}, \theta)$, $\mathbf{z} \sim p(\mathbf{z})$.
- ▶ AF prior decoder path: $p(\mathbf{x}|\mathbf{z}, \theta)$, $\mathbf{z} = g(\epsilon, \lambda)$, $\epsilon \sim p(\epsilon)$.

AF prior and IAF posterior have the same computation cost, so using AF prior makes the model more expressive at no training time cost.

<https://arxiv.org/abs/1611.02731>

Variational Lossy AutoEncoder, 2016

- ▶ Can VLAE learn lossy codes that encode global statistics?
- ▶ Does using AF priors improves upon using IAF posteriors as predicted by theory?
- ▶ Does using autoregressive decoding distributions improve density estimation performance?

CIFAR10

MNIST

Model	NLL Test
Normalizing flows (Rezende & Mohamed, 2015)	85.10
DRAW (Gregor et al., 2015)	< 80.97
Discrete VAE (Rolfe, 2016)	81.01
PixelRNN (van den Oord et al., 2016a)	79.20
IAF VAE (Kingma et al., 2016)	79.88
AF VAE	79.30
VLAE	79.03

Method	bits/dim \leq
<i>Results with tractable likelihood models:</i>	
Uniform distribution [1]	8.00
Multivariate Gaussian [1]	4.70
NICE [2]	4.48
Deep GMMS [3]	4.00
Real NVP [4]	3.49
PixelCNN [1]	3.14
Gated PixelCNN [5]	3.03
PixelRNN [1]	3.00
PixelCNN++ [6]	2.92
<i>Results with variationally trained latent-variable models:</i>	
Deep Diffusion [7]	5.40
Convolutional DRAW [8]	3.58
ResNet VAE with IAF [9]	3.11
ResNet VLAE	3.04
DenseNet VLAE	2.95

<https://arxiv.org/abs/1611.02731>

VAE limitations

- ▶ Poor variational posterior distribution (encoder)

$$q(\mathbf{z}|\mathbf{x}, \phi) = \mathcal{N}(\mathbf{z}|\mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})).$$

- ▶ Poor prior distribution

$$p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I}).$$

- ▶ Poor probabilistic model (decoder)

$$p(\mathbf{x}|\mathbf{z}, \theta) = \mathcal{N}(\mathbf{x}|\mu_\theta(\mathbf{z}), \sigma_\theta^2(\mathbf{z})).$$

- ▶ Loose lower bound

$$p(\mathbf{x}|\theta) - \mathcal{L}(q, \theta) = (?).$$

Importance Sampling

Generative model

$$\begin{aligned} p(\mathbf{x}|\theta) &= \int p(\mathbf{x}, \mathbf{z}|\theta) d\mathbf{z} = \int \left[\frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] q(\mathbf{z}|\mathbf{x}) d\mathbf{z} \\ &= \int f(\mathbf{x}, \mathbf{z}) q(\mathbf{z}|\mathbf{x}) d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} f(\mathbf{x}, \mathbf{z}) \end{aligned}$$

Here $f(\mathbf{x}, \mathbf{z}) = \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})}$.

ELBO

$$\begin{aligned} \log p(\mathbf{x}|\theta) &= \log \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} f(\mathbf{x}, \mathbf{z}) \geq \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log f(\mathbf{x}, \mathbf{z}) = \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} = \mathcal{L}(q, \theta). \end{aligned}$$

Could we choose better $f(\mathbf{x}, \mathbf{z})$?

Let define

$$f(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K) = \frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}_k | \theta)}{q(\mathbf{z}_k | \mathbf{x})}$$

$$\mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q(\mathbf{z} | \mathbf{x})} f(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K) = p(\mathbf{x} | \theta)$$

ELBO

$$\begin{aligned} \log p(\mathbf{x} | \theta) &= \log \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q(\mathbf{z} | \mathbf{x})} f(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K) \geq \\ &\geq \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q(\mathbf{z} | \mathbf{x})} \log f(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K) = \\ &= \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q(\mathbf{z} | \mathbf{x})} \log \left[\frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}_k | \theta)}{q(\mathbf{z}_k | \mathbf{x})} \right] = \mathcal{L}_K(q, \theta). \end{aligned}$$

VAE objective

$$p(\mathbf{x}|\boldsymbol{\theta}) \geq \mathcal{L}(q, \boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})}{q(\mathbf{z}|\mathbf{x})} \rightarrow \max_{\phi, \boldsymbol{\theta}}$$

$$\mathcal{L}(q, \boldsymbol{\theta}) = \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q(\mathbf{z}|\mathbf{x})} \left(\frac{1}{K} \sum_{k=1}^K \log \frac{p(\mathbf{x}, \mathbf{z}_k|\boldsymbol{\theta})}{q(\mathbf{z}_k|\mathbf{x})} \right) \rightarrow \max_{q, \boldsymbol{\theta}}.$$

IWAE objective

$$\mathcal{L}_K(q, \boldsymbol{\theta}) = \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q(\mathbf{z}|\mathbf{x})} \log \left(\frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}_k|\boldsymbol{\theta})}{q(\mathbf{z}_k|\mathbf{x})} \right) \rightarrow \max_{q, \boldsymbol{\theta}}.$$

If $K = 1$, these objectives coincide.

Theorem

1. $\log p(\mathbf{x}|\theta) \geq \mathcal{L}_K(q, \theta) \geq \mathcal{L}_M(q, \theta)$, for $K \geq M$;
2. $\log p(\mathbf{x}|\theta) = \lim_{K \rightarrow \infty} \mathcal{L}_K(q, \theta)$ if $\frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})}$ is bounded.

Proof of 1.

$$\begin{aligned}
\mathcal{L}_K(q, \theta) &= \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K} \log \left(\frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}_k | \theta)}{q(\mathbf{z}_k | \mathbf{x})} \right) = \\
&= \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K} \log \mathbb{E}_{k_1, \dots, k_M} \left(\frac{1}{M} \sum_{m=1}^M \frac{p(\mathbf{x}, \mathbf{z}_{k_m} | \theta)}{q(\mathbf{z}_{k_m} | \mathbf{x})} \right) \geq \\
&\geq \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K} \mathbb{E}_{k_1, \dots, k_M} \log \left(\frac{1}{M} \sum_{m=1}^M \frac{p(\mathbf{x}, \mathbf{z}_{k_m} | \theta)}{q(\mathbf{z}_{k_m} | \mathbf{x})} \right) = \\
&= \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_M} \log \left(\frac{1}{M} \sum_{m=1}^M \frac{p(\mathbf{x}, \mathbf{z}_m | \theta)}{q(\mathbf{z}_m | \mathbf{x})} \right) = \mathcal{L}_M(q, \theta)
\end{aligned}$$

Theorem

1. $\log p(\mathbf{x}|\theta) \geq \mathcal{L}_K(q, \theta) \geq \mathcal{L}_M(q, \theta)$, for $K \geq M$;
2. $\log p(\mathbf{x}|\theta) = \lim_{K \rightarrow \infty} \mathcal{L}_K(q, \theta)$ if $\frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})}$ is bounded.

Proof of 2.

Consider r.v. $\xi_K = \frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}_k|\theta)}{q(\mathbf{z}_k|\mathbf{x})}$.

If summands are bounded, then (from the strong law of large numbers)

$$\xi_K \xrightarrow[K \rightarrow \infty]{\text{a.s.}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} = p(\mathbf{x}|\theta).$$

Hence $\mathcal{L}_K(q, \theta) = \mathbb{E} \log \xi_K$ converges to $\log p(\mathbf{x}|\theta)$ as $K \rightarrow \infty$.

$$\log p(\mathbf{x}|\theta) \geq \mathcal{L}_K(q, \theta) \geq \mathcal{L}(q, \theta)$$

If $K > 1$ the bound could be tighter.

$$\mathcal{L}(q, \theta) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})};$$

$$\mathcal{L}_K(q, \theta) = \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q(\mathbf{z}|\mathbf{x})} \log \left(\frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}_k|\theta)}{q(\mathbf{z}_k|\mathbf{x})} \right).$$

- ▶ $\mathcal{L}_1(q, \theta) = \mathcal{L}(q, \theta);$
- ▶ $\mathcal{L}_\infty(q, \theta) = \log p(\mathbf{x}|\theta).$

Which $q(\mathbf{z}|\mathbf{x})$ gives $\mathcal{L}(q, \theta) = \log p(\mathbf{x}|\theta)?$

Which $q(\mathbf{z}|\mathbf{x})$ gives $\mathcal{L}(q, \theta) = \mathcal{L}_K(q, \theta)?$

Theorem

The VAE objective is equal to IWAE objective

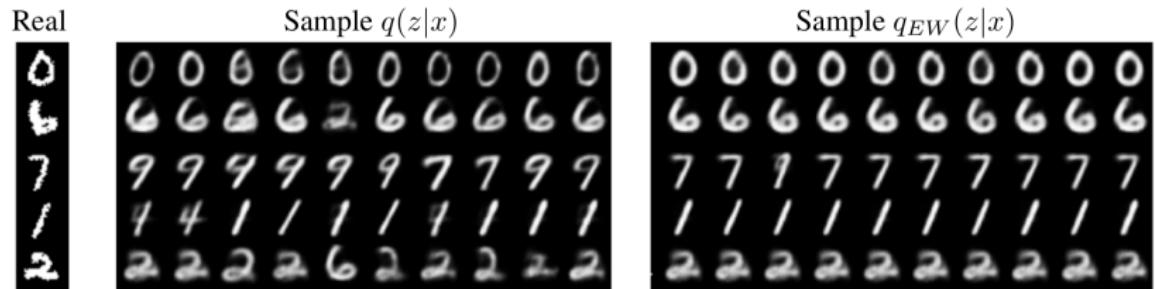
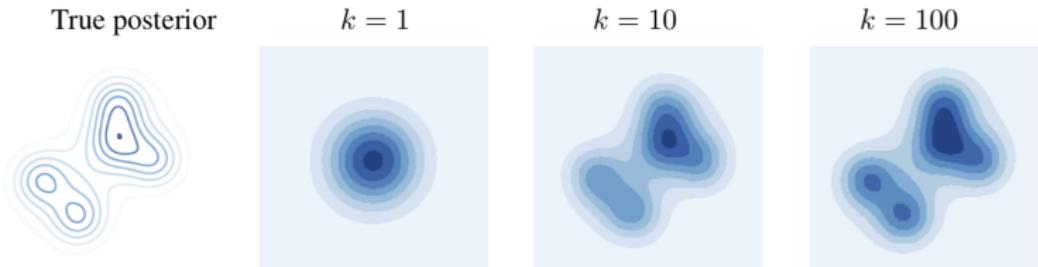
$$\mathcal{L}(q_{EW}, \theta) = \mathcal{L}_K(q, \theta)$$

for the following variational distribution

$$q_{EW}(\mathbf{z}|\mathbf{x}) = \mathbb{E}_{\mathbf{z}_2, \dots, \mathbf{z}_K \sim q(\mathbf{z}|\mathbf{x})} q_{IW}(\mathbf{z}|\mathbf{x}, \mathbf{z}_{2:K}),$$

where

$$q_{IW}(\mathbf{z}|\mathbf{x}, \mathbf{z}_{2:K}) = \frac{\frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})}}{\frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}, \mathbf{z}_k)}{q(\mathbf{z}_k|\mathbf{x})}} q(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{\frac{1}{K} \left(\frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} + \sum_{k=2}^K \frac{p(\mathbf{x}, \mathbf{z}_k)}{q(\mathbf{z}_k|\mathbf{x})} \right)}.$$



Summary

- ▶ ELBO surgery reveals insights about prior distribution in VAE.
The optimal prior is aggregated posterior.
- ▶ VampPrior proposed to use variational mixture of posteriors as prior to approximate the aggregated posterior.
- ▶ More powerful decoder in VAE leads to more expressive generative model. However, too expressive decoder could lead to posterior collapse.
- ▶ Decoder weakening is a set of techniques to avoid posterior collapse.
- ▶ Autoregressive flows could be used in prior. This is equivalent to the use of IAF posterior.
- ▶ Importance sampling could get tighter lower bound to the likelihood.

References

- ▶ **VAE with a VampPrior**
<https://arxiv.org/abs/1705.07120>
Summary: Variational Mixture of Posteriors prior is introduced. The VampPrior components are given by variational posteriors conditioned on learnable pseudo-inputs. Prior is extended to a two layer hierarchical model with a coupled prior and posterior, it learns significantly better models. The model avoids the local optima issues related to useless latent dimensions that plague VAEs. The prior is compared with standard gaussian and mixture of gaussians.
- ▶ **PixelVAE: A Latent Variable Model for Natural Images**
<https://arxiv.org/abs/1611.05013>
Summary: Use autoregressive model (PixelCNN) in VAE decoder. Use hierarchical structure of latent variables. Restrict the number of conv layers in PixelCNN (receptive field) to encode global structure in latent space. The performance is slightly worse than GatedPixelCNN.
- ▶ **Improving Variational Inference with Inverse Autoregressive Flow**
<https://arxiv.org/abs/1606.04934>
Summary: Free bits proposed to weaken autoregressive decoder.
- ▶ **VLAE: Variational Lossy Autoencoder**
<https://arxiv.org/abs/1611.02731>
Summary: Bits-back coding interpretation for posterior collapse. Solving the problem of ignoring latent codes. Reduce receptive field in PixelCNN decoder to encode global information, use learnable AF prior.
- ▶ **IWAE: Importance Weighted Autoencoders**
<https://arxiv.org/abs/1509.00519>
Summary: Propose the version of ELBO which is tighter to the log-density. Sampling of k objects from latent space instead of one. The more samples we sample the more tight the gap is. In the limit we will get the true likelihood. The gradient is given by importance sampling reweighting. Analyze the number of active latent units.
- ▶ **Reinterpreting Importance-Weighted Autoencoders**
<https://arxiv.org/abs/1704.02916>
Summary: IWAE maximizes a tighter ELBO than VAE ELBO. This ELBO is also related to some implicit variational distribution. The analytical form of this distribution is given and visualized.

$$\frac{a_1 + \cdots + a_K}{K} = \mathbb{E}_{i_1, \dots, i_M} \frac{a_{i_1} + \cdots + a_{i_M}}{M}, \quad i_1, \dots, i_M \sim U[1, K]$$