

Learning Hierarchical Features from Generative Models

Shengjia Zhao, Jiaming Song, Stefano Ermon
Stanford University, ICML 2017

Presented by Zhe Gan, Duke University

September 22nd, 2017

- Background:
 - Deep neural nets have been shown to be very successful at learning feature hierarchies in supervised learning tasks.
 - Generative models, on the other hand, have benefited less from hierarchical models with multiple layers of latent variables.
- Contribution:
 - proved that hierarchical latent variable models do not take advantage of the hierarchical structure when trained with existing variational methods.
 - hence provide limitations on the kind of features existing models can learn.
 - proposed [variational ladder autoencoder](#) to learn interpretable and disentangled hierarchical features.

Problem setting: variational autoencoder

- Generative model: $p_{\theta}(\mathbf{x}|\mathbf{z})$
- Inference model: $q_{\phi}(\mathbf{z}|\mathbf{x})$
- Variational lower bound:

$$\log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; \theta, \phi) \quad (1)$$

$$= \mathbb{E}_{p_{data}(\mathbf{x})} [\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - KL(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))] \quad (2)$$

- Hierarchical VAE:

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z}_1) \prod_{\ell=1}^{L-1} p(\mathbf{z}_{\ell}|\mathbf{z}_{\ell+1})p(\mathbf{z}_L) \quad (3)$$

Limitations of Hierarchical VAEs

- Representational Efficiency: one layer model is enough for generative modeling.

Proposition

\mathcal{L}_{ELBO} is globally maximized as a function of $q_{z|x}$ and $p(x|z)$ when $\mathcal{L}_{ELBO} = -H(p_{data}(x))$. If \mathcal{L}_{ELBO} is globally maximized, the following Gibbs sampling chain converges to $p_{data}(x)$ if it is ergodic

$$z_1^{(t)} \sim q(z_1|x^{(t)}) \quad (4)$$

$$x^{(t+1)} \sim p(x|z_1^{(t)}) \quad (5)$$

- Proof: under ideal optimization of \mathcal{L}_{ELBO} , $p(x) = \int_z p(x, z) dz = p_{data}(x)$ and $q(z|x) = p(z|x)$, this also implies $q(x|z_1) = p(x|z_1)$. Because the following Gibbs chain converges to $p_{data}(x)$ when it is ergodic

$$z_1^{(t)} \sim q(z_1|x^{(t)}), x^{(t+1)} \sim q(x|z_1^{(t)}) \quad (6)$$

Limitations of Hierarchical VAEs

- Under the assumptions of the above Proposition, we can sample from $p_{data}(\mathbf{x})$ without using the latent code $(\mathbf{z}_2, \dots, \mathbf{z}_L)$ at all.
- The Gibbs chain generates samples (left figure) with similar visual quality as ancestral sampling with the entire model (right figure), even though the Gibbs chain only used the bottom layer of the model

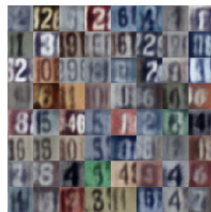
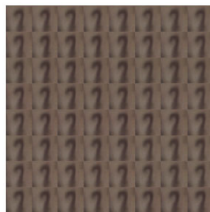
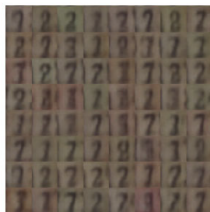
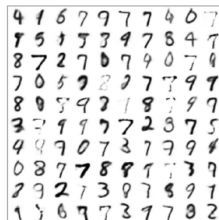
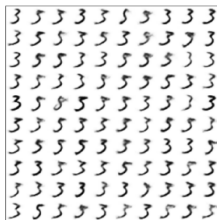
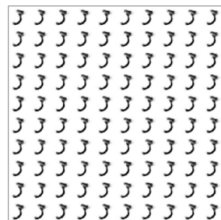


Limitations of Hierarchical VAEs

- **Feature learning:** no hierarchical and disentangled features are learned in a typical HVAE.
- We can consider $q(\mathbf{z}|\mathbf{x})$ as a (probabilistic) feature detector. It is natural to think that q might learn hierarchical features similarly to a feed-forward network $\mathbf{x} \rightarrow \mathbf{z}_1 \rightarrow \dots \rightarrow \mathbf{z}_L$
- $q(\mathbf{z}_{>\ell}|\mathbf{z}_\ell)$ maps low-level features to high-level features, then $q(\mathbf{z}_\ell|\mathbf{z}_{>\ell})$ do the reverse.
- Under ideal optimization of \mathcal{L}_{ELBO} , $q(\mathbf{z}_\ell|\mathbf{z}_{>\ell}) = p(\mathbf{z}_\ell|\mathbf{z}_{>\ell})$
- $p(\mathbf{z}_\ell|\mathbf{z}_{>\ell})$ is typically a simple distribution family such as Gaussians, therefore, under this case, the only type of feature hierarchy we can hope to learn is one under which $q(\mathbf{z}_\ell|\mathbf{z}_{>\ell})$ is Gaussian.

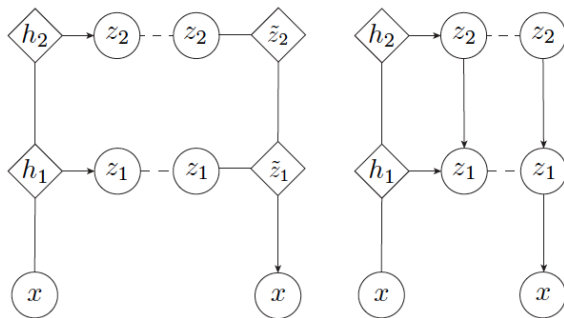
Limitations of Hierarchical VAEs

Only very minor variations correspond to lower layers (left and center), and almost all the variation is represented by the top layer (right).



Variational Ladder Autoencoders

- **Proposed:** one-layer model that learns feature hierarchies
- **Intuition:** If z_i is more abstract than z_j , then the inference mapping $q(z_i|\mathbf{x})$ and generative mapping when other layers are fixed $p(\mathbf{x}|z_i, \mathbf{z}_{-i} = \text{fixed})$ requires a more expressive network to capture.
- (Left) VLAE; (Right) LVAE.



Variational Ladder Autoencoders

- Generative model:

$$p(\mathbf{z}) = p(\mathbf{z}_1, \dots, \mathbf{z}_L) = \mathcal{N}(0, \mathbf{I}).$$

$p(\mathbf{x}|\mathbf{z}_1, \dots, \mathbf{z}_L)$ is defined as

$$\tilde{\mathbf{z}}_L = f_L(\mathbf{z}_L) \quad (7)$$

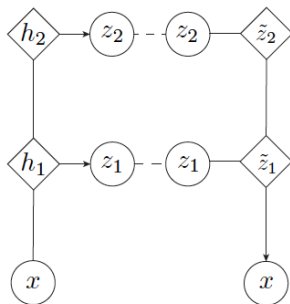
$$\tilde{\mathbf{z}}_\ell = f_\ell(\tilde{\mathbf{z}}_{\ell+1}, \mathbf{z}_\ell) \quad (8)$$

$$\mathbf{x} \sim r(\mathbf{x}; f_0(\tilde{\mathbf{z}}_1)) \quad (9)$$

- Inference network:

$$\mathbf{h}_\ell = g_\ell(\mathbf{h}_{\ell-1}) \quad (10)$$

$$\mathbf{z}_\ell \sim \mathcal{N}(\boldsymbol{\mu}_\ell(\mathbf{h}_\ell), \boldsymbol{\sigma}_\ell(\mathbf{h}_\ell)) \quad (11)$$



Experiments

1st layer encodes stroke width; 2nd layer encodes digit width and tilt; 3rd layer encodes digit identity.

6	2	3	4	3	4	3	7	8	9	7	4	8	5	9	5	5	2
2	8	6	5	8	8	3	4	4	7	4	7	4	8	6	9	5	6
4	4	0	8	4	2	3	4	4	1	0	8	5	5	5	5	6	6
8	7	4	7	7	1	6	2	1	6	5	1	8	2	7	8	1	1
4	4	3	1	4	3	4	3	4	1	0	6	7	6	6	6	6	1
0	3	1	4	3	5	4	7	7	7	1	3	9	7	2	5	9	1
2	4	3	2	0	2	4	4	4	7	9	9	6	9	1	0	1	8
6	2	8	7	4	9	7	6	2	1	4	4	3	1	7	7	1	8
4	0	4	6	2	5	6	4	5	3	9	8	1	0	5	2	2	9
1	3	2	2	4	6	2	4	6	3	4	9	1	7	7	9	9	1
6	7	0	7	2	4	0	4	2	4	1	1	2	1	2	7	6	8
3	6	1	2	7	7	0	4	4	8	4	9	9	9	1	8	1	1
3	4	3	4	2	4	2	2	7	6	5	1	5	5	2	6	5	2
4	0	5	4	6	6	6	7	6	0	2	9	5	9	7	8	5	7
0	0	4	4	3	9	8	2	0	2	1	9	5	1	8	1	6	2
5	0	1	4	7	3	4	4	2	4	9	7	5	5	1	3	5	6
0	6	5	3	4	4	6	4	5	6	1	1	0	9	1	5	0	1
4	5	4	5	4	6	4	5	7	6	9	9	7	1	1	9	6	8
4	4	7	5	4	1	3	2	8	6	9	4	9	9	7	8	5	9
4	5	6	0	8	4	2	6	6	1	2	8	5	5	6	6	8	5

4	1	1	1	6	4	4	4	1	4	5	6	4	4	8
4	1	6	3	1	6	4	8	3	4	6	6	7	5	9
6	4	5	4	4	1	6	6	5	0	9	7	6	8	5
1	4	1	4	1	4	6	1	8	6	4	5	1	4	6
4	6	1	3	1	5	4	4	8	8	9	4	9	5	8
8	1	6	1	1	5	4	2	1	6	3	4	1	4	8
3	6	6	6	4	6	6	4	6	5	4	5	4	2	2
2	3	8	1	7	6	4	2	1	8	5	8	1	3	8
8	1	5	1	6	2	6	4	6	6	8	4	5	4	5
0	6	5	8	6	7	5	4	4	9	2	1	4	8	4
4	6	2	2	4	4	7	8	7	2	2	4	0	6	0
1	4	7	2	3	4	8	2	4	8	5	6	5	7	9
4	2	3	4	4	9	9	8	4	9	9	5	6	8	2
0	7	2	7	5	9	4	4	2	1	4	0	4	8	5
0	9	6	5	6	2	5	6	2	6	5	0	7	0	5
5	2	3	4	4	6	2	0	0	8	7	7	2	5	8
8	7	2	2	6	8	2	2	5	9	4	9	8	2	0
6	4	5	6	9	4	8	5	5	4	5	7	7	5	5
6	0	3	2	2	3	8	7	2	7	5	0	0	0	7
3	0	7	0	0	7	2	2	9	2	0	5	5	2	8

[illegible]

Experiments

1st: color schemes; 2nd: shape variations; 3rd: digit identity; 4th: the general structure of the image.

