

(NIPS 2015)

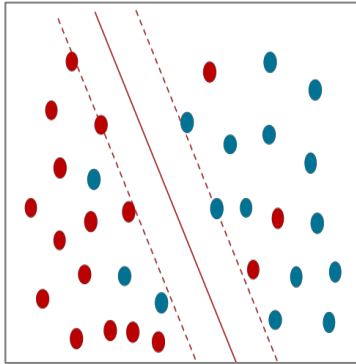
# Learning Structured Output Representation Using Deep Conditional Generative Models

Sohn Kihyuk, Honglak Lee, Xinchun Yan

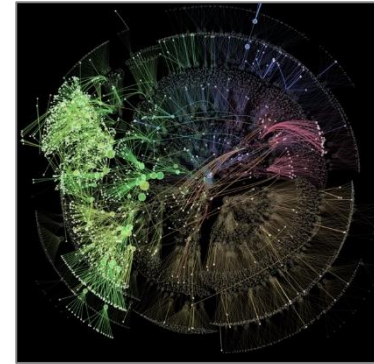
---

박정수

Data Mining & Information Systems Lab.  
Department of Computer Science and Engineering,  
College of Informatics, Korea University



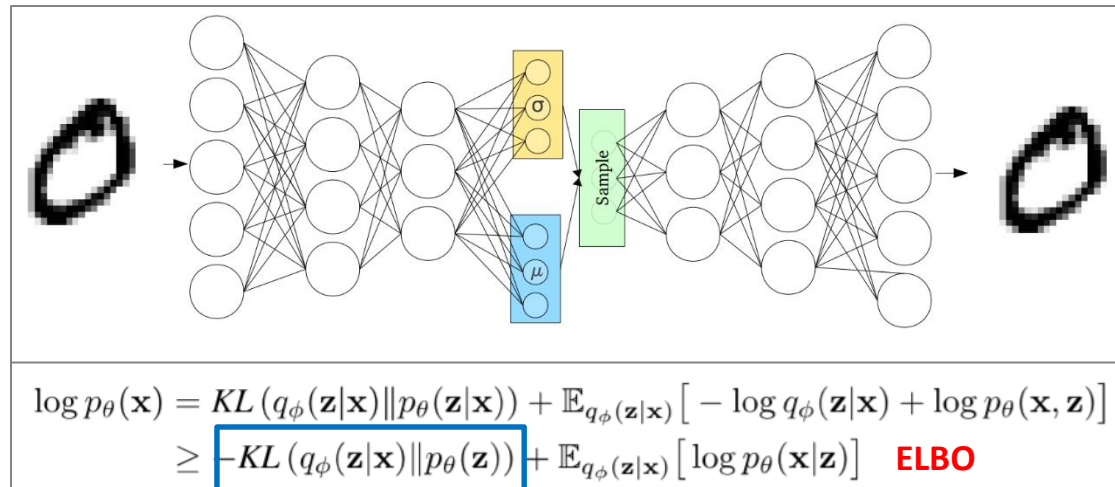
Output Space : 1D



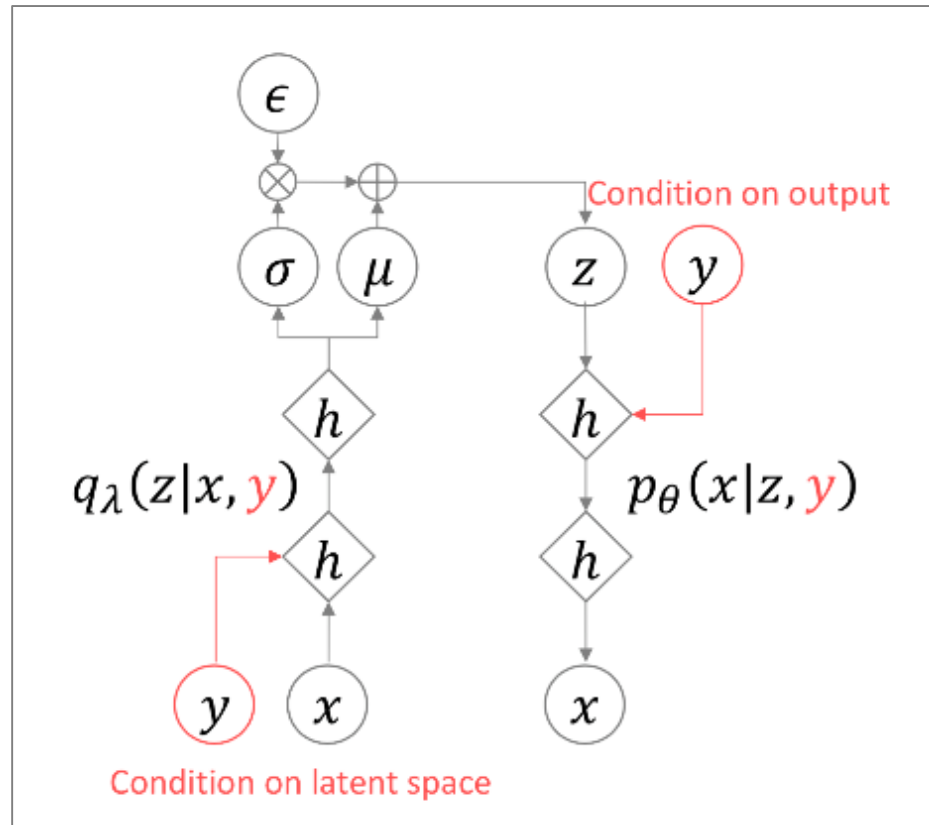
Output Space : High Dimension

- Output representation learning and structure(manifold) prediction required
- Deterministic NN algorithms(i.e. CNN) not suitable for modeling a distribution with multiple modes

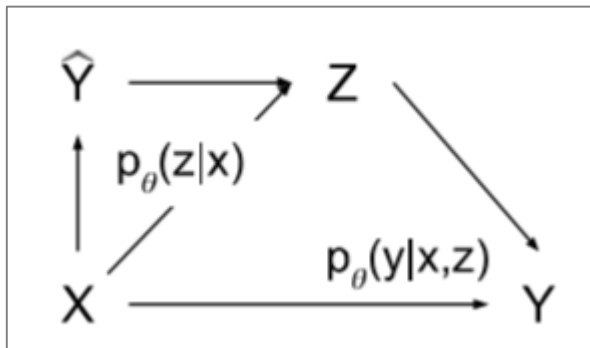
## Variational Autoencoder



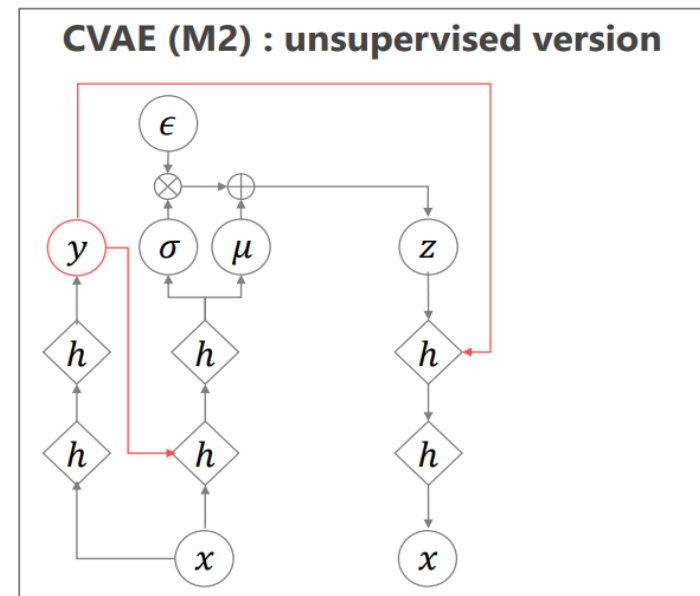
- Using the hypothesis that given data  $\mathbf{X}$  can be generated from the normal distribution( $\mathbf{Z}$ ), VAE uses variational approach for latent representation learning in the form of AE
- Another Loss Term(KL Divergence) was introduced as a result of variational approach



For modeling the high dimensional output space conditioned on the input observation, CVAE was proposed



In the case of semi or unsupervised learning, we use the predicted  $y$  value as a condition label information



$$\log p_{\theta}(\mathbf{y}|\mathbf{x}) \geq -KL(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})||p_{\theta}(\mathbf{z}|\mathbf{x})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} [\log p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z})]$$

$$\tilde{\mathcal{L}}_{\text{CVAE}}(\mathbf{x}, \mathbf{y}; \theta, \phi) = \underline{-KL(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})||p_{\theta}(\mathbf{z}|\mathbf{x}))} + \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z}^{(l)}),$$

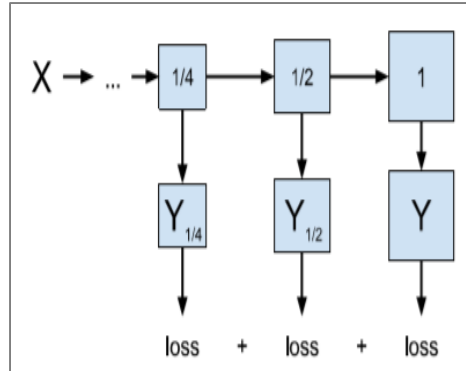
In the phase of training, we used the recognition model, but as for the testing phase, we only use prior distribution.



Therefore, we set the prior same as recognition model and derive the loss function as below.

$$\tilde{\mathcal{L}}_{\text{GSNN}}(\mathbf{x}, \mathbf{y}; \theta, \phi) = \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z}^{(l)}) \quad \Bigg| \quad \tilde{\mathcal{L}}_{\text{hybrid}} = \alpha \tilde{\mathcal{L}}_{\text{CVAE}} + (1 - \alpha) \tilde{\mathcal{L}}_{\text{GSNN}},$$

## Multi scale prediction



- Predict output at different scales.
- Being able to achieve global to local prediction of pixel level segmentation

## Input Omission noise



- Random block omission noise
- Providing more challenging output prediction task during training thus model becomes more robust

As for the evaluation of trained model, we can make a prediction accuracy by performing a deterministic inference

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z}^*), \mathbf{z}^* = \mathbb{E}[\mathbf{z}|\mathbf{x}]$$

Alternatively, we can compare the conditional likelihoods

$$p_{\theta}(\mathbf{y}|\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z}^{(s)}), \quad \mathbf{z}^{(s)} \sim p_{\theta}(\mathbf{z}|\mathbf{x})$$

Since MC algorithm requires a large number of samples, we can use importance sampling

$$p_{\theta}(\mathbf{y}|\mathbf{x}) \approx \frac{1}{S} \sum_{s=1}^S \frac{p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z}^{(s)})p_{\theta}(\mathbf{z}^{(s)}|\mathbf{x})}{\underline{q_{\phi}(\mathbf{z}^{(s)}|\mathbf{x}, \mathbf{y})}}, \quad \mathbf{z}^{(s)} \sim q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})$$

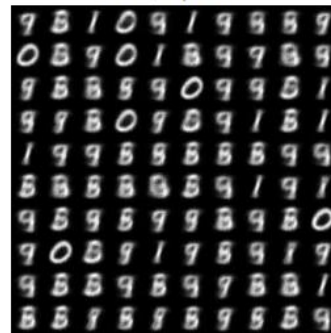




input



CVAE, epoch 1



VAE, epoch 1



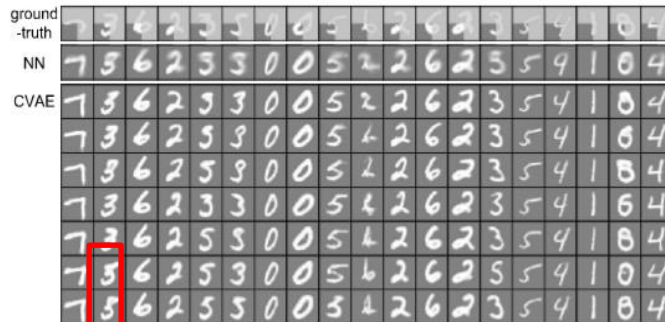
CVAE, epoch 20



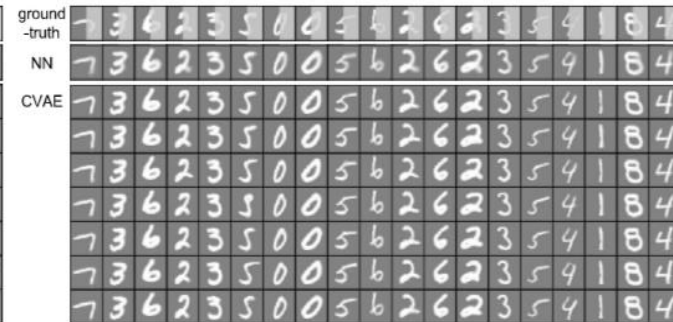
VAE, epoch 20

Converges much faster than VAE for the same epoch

## Imputation of MNIST images



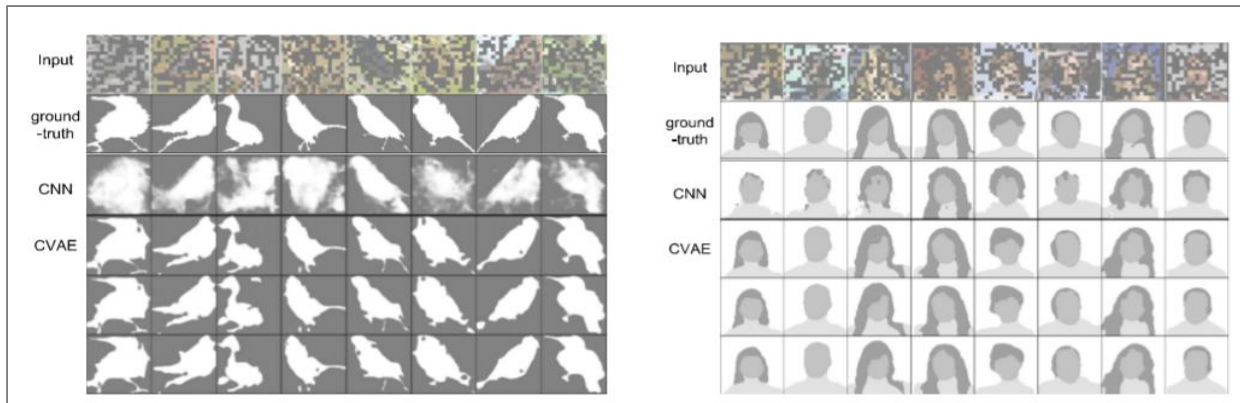
1 quadrant



2 quadrant

negative CLL	1 quadrant		2 quadrants		3 quadrants	
	validation	test	validation	test	validation	test
NN (baseline)	100.03	99.75	62.14	62.18	26.01	25.99
GSNN (Monte Carlo)	100.03	99.82	62.48	62.41	26.20	26.29
CVAE (Monte Carlo)	68.62	68.39	45.57	45.34	20.97	20.96
CVAE (Importance Sampling)	<b>64.05</b>	<b>63.91</b>	<b>44.96</b>	<b>44.73</b>	<b>20.97</b>	<b>20.95</b>
Performance gap	35.98	35.91	17.51	17.68	5.23	5.33
- per pixel	0.061	0.061	0.045	0.045	0.027	0.027

## Object segmentation with Partial Observations



Dataset		CUB (IoU)		LFW (pixel)	
noise level	block size	GDNN	CVAE	GDNN	CVAE
25%	1	89.37	<b>98.52</b>	96.93	<b>99.22</b>
	4	88.74	<b>98.07</b>	96.55	<b>99.09</b>
	8	90.72	<b>96.78</b>	97.14	<b>98.73</b>
50%	1	74.95	<b>95.95</b>	91.84	<b>97.29</b>
	4	70.48	<b>94.25</b>	90.87	<b>97.08</b>
	8	76.07	<b>89.10</b>	92.68	<b>96.15</b>
70%	1	62.11	<b>89.44</b>	85.27	<b>89.71</b>
	4	57.68	<b>84.36</b>	85.70	<b>93.16</b>
	8	63.59	<b>76.87</b>	87.83	<b>92.06</b>

## Proof of ELBO in CVAE

$$\log p_{\theta}(\mathbf{y}|\mathbf{x}) = KL(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})||p_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{y})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} [-\log q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y}) + \log p_{\theta}(\mathbf{y}, \mathbf{z}|\mathbf{x})] \quad (\text{S1})$$

$$\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} [-\log q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y}) + \log p_{\theta}(\mathbf{y}, \mathbf{z}|\mathbf{x})] \quad (\text{S2})$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} [-\log q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y}) + \log p_{\theta}(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} [\log p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z})] \quad (\text{S3})$$

$$= -KL(q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})||p_{\theta}(\mathbf{z}|\mathbf{x})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} [\log p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{z})] \quad (\text{S4})$$