

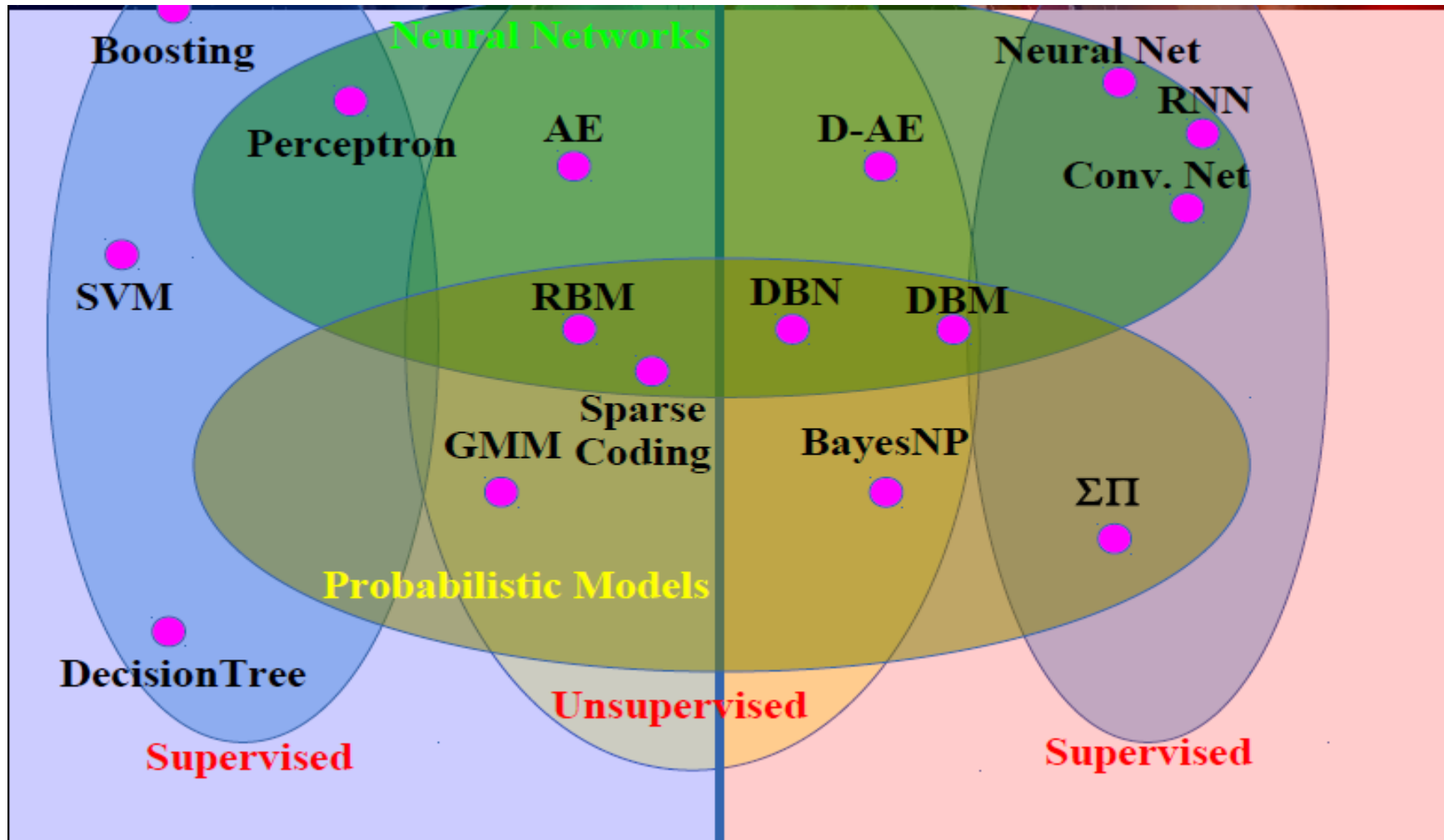
Marriage of Graphic Model and Neural Model

Dong Wang

2016/5/23

Content

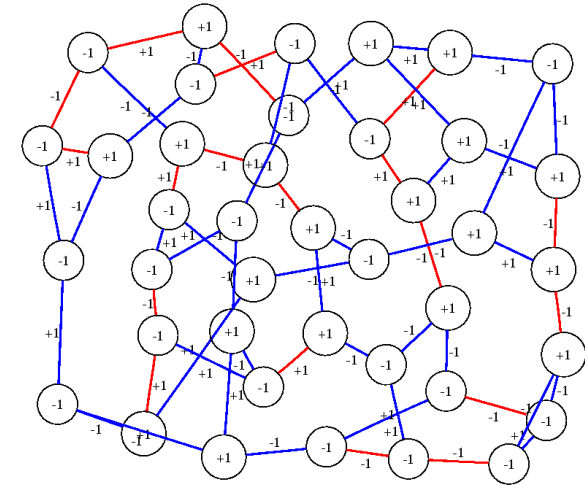
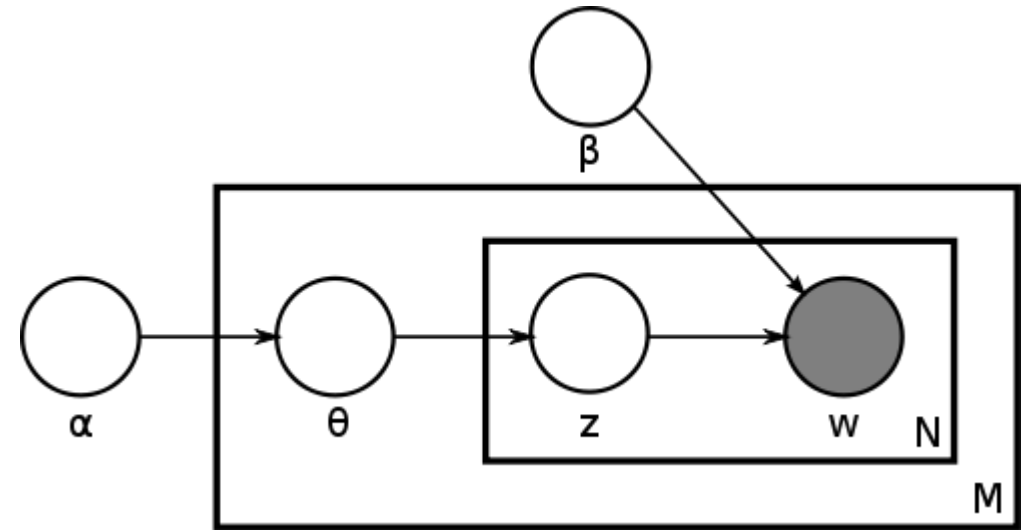
- What are they?
- Marriage 1: Variational AE
- Marriage 2: Denoise AE



Yann LeCun, Marc'Aurelio Ranzato, Deep Learning Tutorial,
 ICML, Atlanta, 2013-06-16

Graphic models

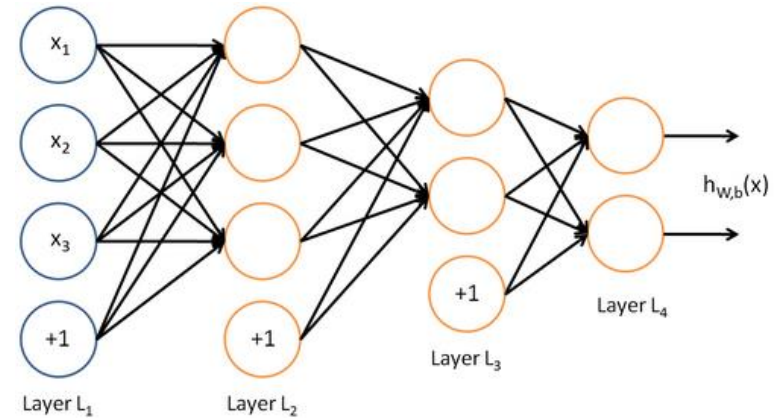
- $G=(V,E)$ represent joint probabilities of V , with conditionals or potentials represented by E
- Probabilistic variables
- Probabilistic inference



Graphical models and variational methods: Message-passing and relaxations , ICML-2008 Tutorial

Neural models

- $G=(V,E)$ represents deterministic inference
- Probabilistic interpolation: Gaussian, Binomial, or MDN.
- Some 'randomness' on input, label, hidden units



Respective cons and pros

- Graphical model
 - Clear definition of facts and their relations
 - Easy to grow
 - Difficult in inference
- Neural models
 - Simple and Homogeneous units
 - Quick inference
 - Difficulty in training
 - Less probabilistic

Some models are in both...

- RBM, DBN, DBM, SGN...
- Clear probabilistic interpolation
- Homogeneous units

Content

- What are they?
- Marriage 1: Variational AE
- Marriage 2: Denoise AE

How to marriage them in more depth?

- For Bayesian models, hope simpler inference
- For neural models, hope more randomness
- These two directions seem prefer the same architecture: stochastic NN.

Variational Bayesian

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) = D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})) + \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)})$$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})]$$

- Mean-field variational approximation
- $q(\mathbf{z}|\mathbf{x})=q(z_1|\mathbf{x})q(z_2|\mathbf{x})\dots$

Variational Bayesian with Auto-encoder

$$\log p_{\theta}(\mathbf{x}^{(i)}) = D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})) + \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$$

$$\log p_{\theta}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{\theta}(\mathbf{x}, \mathbf{z})]$$

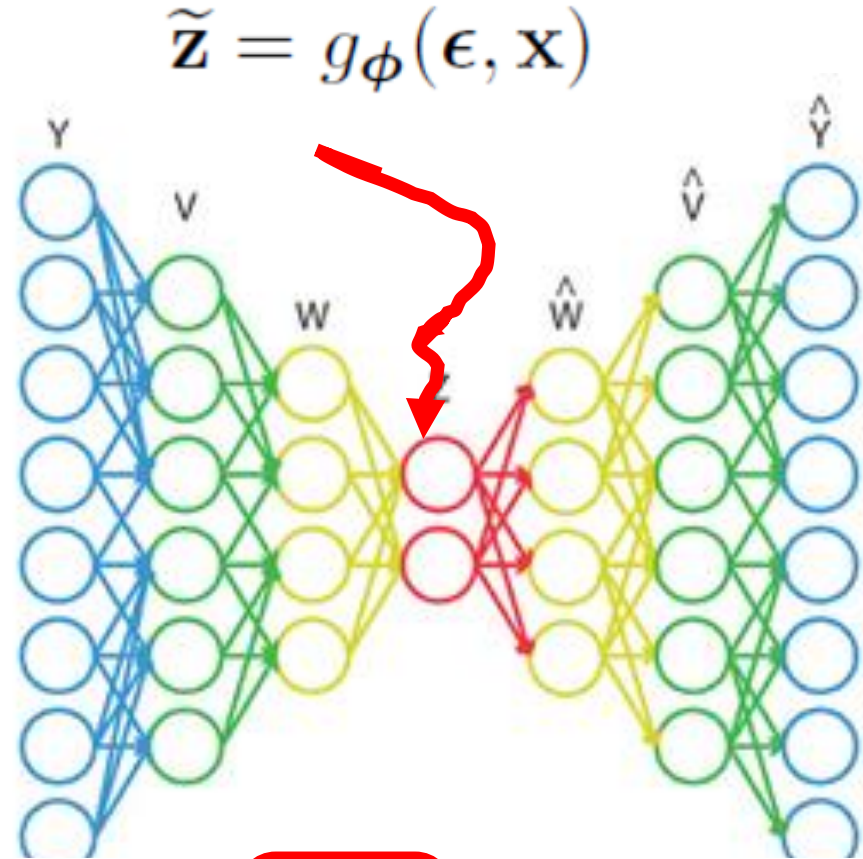
$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})]$$

$$\tilde{\mathbf{z}} = g_{\phi}(\epsilon, \mathbf{x}) \quad \text{with} \quad \epsilon \sim p(\epsilon)$$

- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. In ICLR.
- Danilo J. Rezende, Shakir Mohamed, Daan Wierstra, Stochastic Backpropagation and Approximate Inference in Deep Generative Models, ICMS 2014.

Variational Auto-encoder

$$\tilde{\mathbf{z}} = g_{\phi}(\boldsymbol{\epsilon}, \mathbf{x}) \quad \text{with} \quad \boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$$



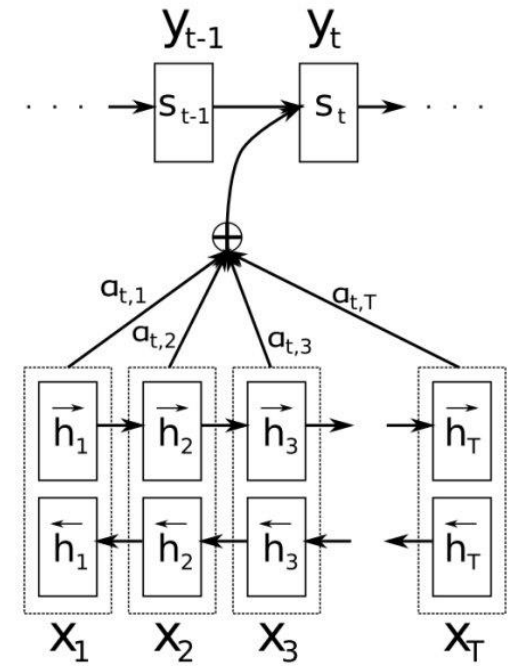
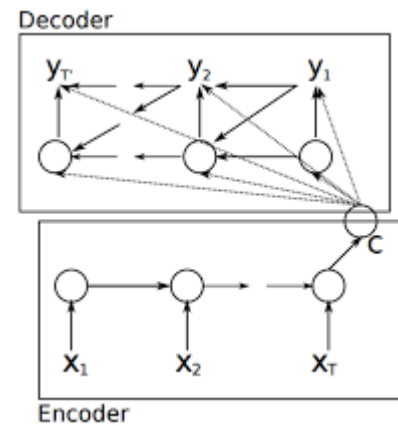
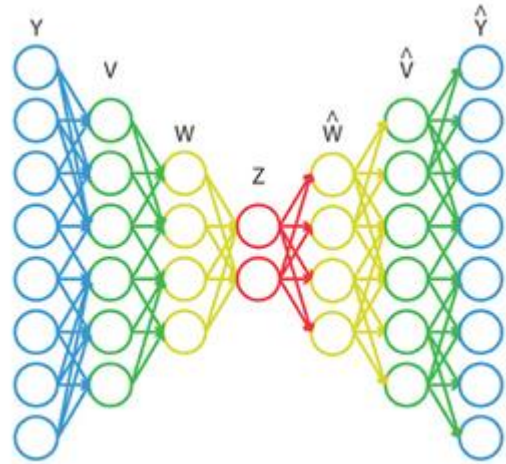
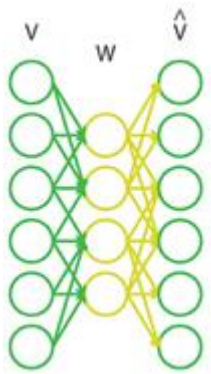
$$\tilde{\mathcal{L}}^B(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z})) + \frac{1}{L} \sum_{l=1}^L (\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)}))$$

where $\mathbf{z}^{(i,l)} = g_{\phi}(\boldsymbol{\epsilon}^{(i,l)}, \mathbf{x}^{(i)})$ and $\boldsymbol{\epsilon}^{(l)} \sim p(\boldsymbol{\epsilon})$

What have been changed?

- Bayesian perspective
 - A encoder (parametric function) is used to map input x to code z , where the variation $p(z)$ is simpler than $p(x)$.
 - With x , $p(z|x)$ keeps simple
 - With z , conditional probability $p(x|z)$ is simpler than x .
 - All seems simpler!
 - Model training becomes parameter adjustment, using BP.
- Neural model perspective
 - Randomness
 - Can we BP? Using MCMC, on the simple $p(x|z)$.
 - Seems a variational + MCMC

Extend to other encoder-decoder models



RNN with latent variable

$$\mathbf{z}_t \sim \mathcal{N}(\boldsymbol{\mu}_{0,t}, \text{diag}(\boldsymbol{\sigma}_{0,t}^2)) \quad (5)$$

$$[\boldsymbol{\mu}_{0,t}, \boldsymbol{\sigma}_{0,t}] = \varphi_{\tau}^{\text{prior}}(\mathbf{h}_{t-1})$$

$$\mathbf{x}_t \mid \mathbf{z}_t \sim \mathcal{N}(\boldsymbol{\mu}_{x,t}, \text{diag}(\boldsymbol{\sigma}_{x,t}^2))$$

$$[\boldsymbol{\mu}_{x,t}, \boldsymbol{\sigma}_{x,t}] = \varphi_{\tau}^{\text{dec}}(\varphi_{\tau}^{\mathbf{z}}(\mathbf{z}_t), \mathbf{h}_{t-1})$$

$$\mathbf{h}_t = f_{\theta}(\varphi_{\tau}^{\mathbf{x}}(\mathbf{x}_t), \varphi_{\tau}^{\mathbf{z}}(\mathbf{z}_t), \mathbf{h}_{t-1})$$

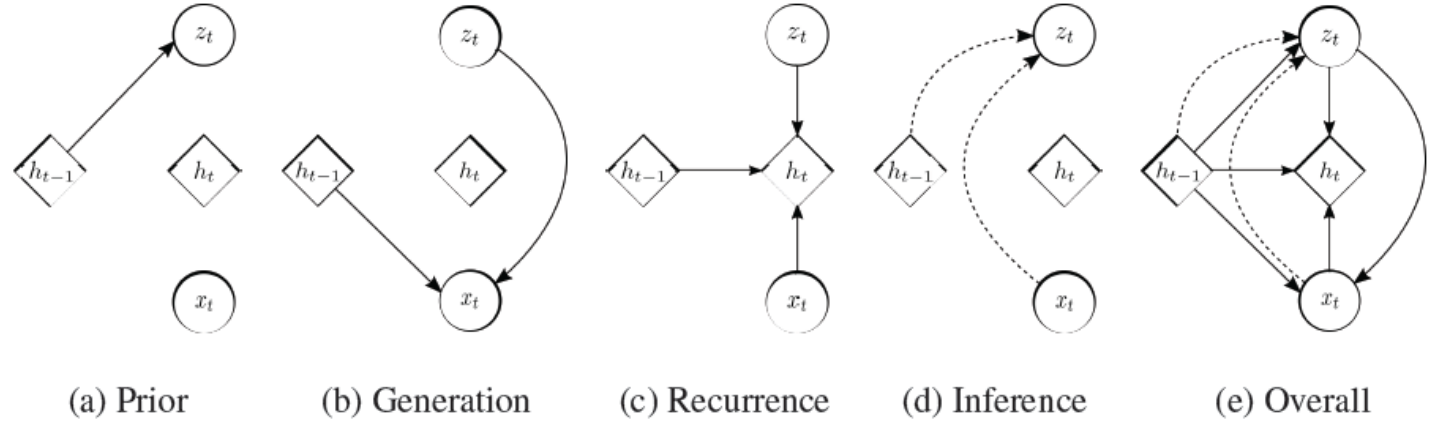
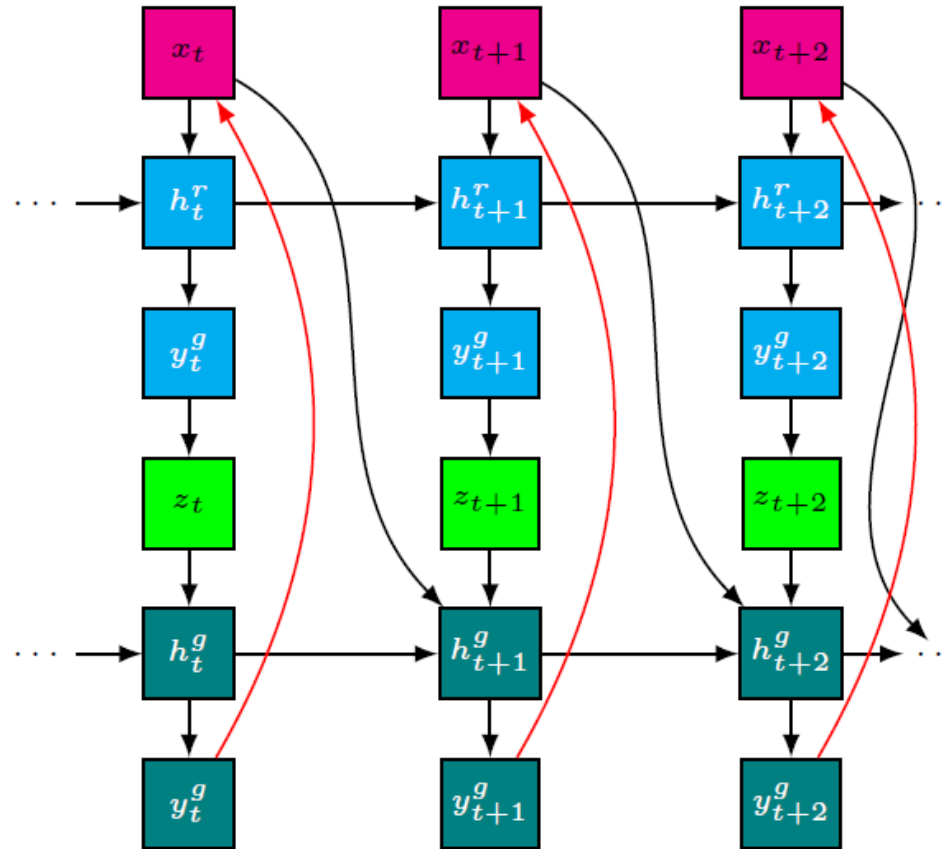


Figure 1: Graphical illustrations of each operation of the VRNN: (a) computing the conditional prior using Eq. (5); (b) generating function using Eq. (6); (c) updating the RNN hidden state using Eq. (7); (d) inference of the approximate posterior using Eq. (9); (e) overall computational paths of the VRNN.

$$p(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}) = \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{z}_{\leq t}, \mathbf{x}_{<t}) p(\mathbf{z}_t \mid \mathbf{x}_{<t}, \mathbf{z}_{<t}). \quad \mathbb{E}_{q(\mathbf{z}_{\leq T} \mid \mathbf{x}_{\leq T})} \left[\sum_{t=1}^T (-\text{KL}(q(\mathbf{z}_t \mid \mathbf{x}_{\leq t}, \mathbf{z}_{<t}) \parallel p(\mathbf{z}_t \mid \mathbf{x}_{<t}, \mathbf{z}_{<t})) + \log p(\mathbf{x}_t \mid \mathbf{z}_{\leq t}, \mathbf{x}_{<t})) \right]$$

- Chung, J., Kastner, K., Dinh, L., Goel, K., Courville, A., and Bengio, Y. (2015). A recurrent latent variable model for sequential data. In NIPS, pages 2962–2970.

Stochastic Recurrent Networks (STORNs)



- Bayer, J. and Osendorfer, C. (2014). Learning stochastic recurrent networks. In NIPS Workshop on Advances in Variational Inference

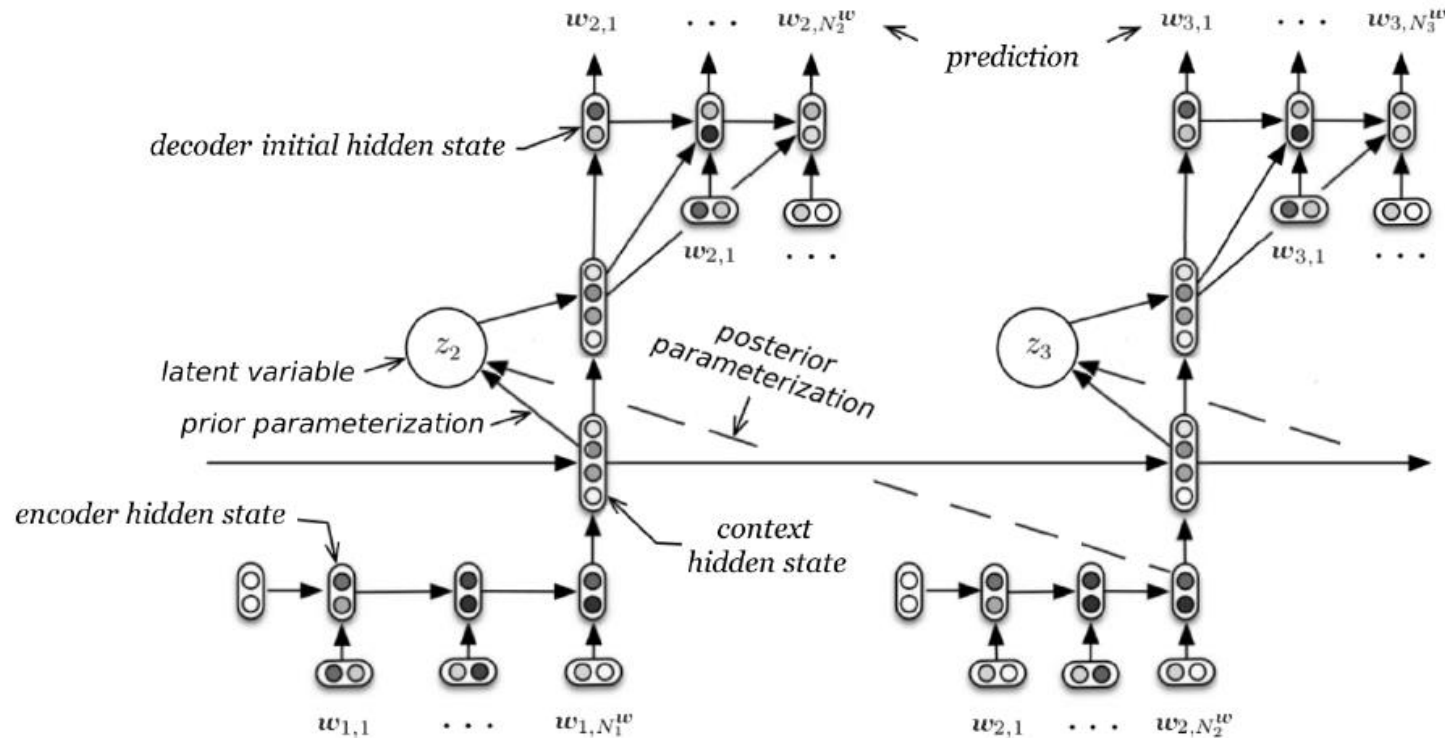
Variational Recurrent AE (VRAE)

$$\begin{aligned}h_{t+1} &= \tanh(W_{enc}^T h_t + W_{in}^T x_{t+1} + b_{enc}) \\ \mu_z &= W_{\mu}^T h_{end} + b_{\mu} \\ \log(\sigma_z) &= W_{\sigma}^T h_{end} + b_{\sigma}\end{aligned}$$

$$\begin{aligned}h_0 &= \tanh(W_z^T z + b_z) \\ h_{t+1} &= \tanh(W_{dec}^T h_t + W_x^T x_t + b_{dec}) \\ x_t &= \text{sigm}(W_{out}^T h_t + b_{out})\end{aligned}$$

- Fabius, O. and van Amersfoort, J. R. (2014). Variational recurrent auto-encoders. arXiv:1412.6581.
- Music generation

Variable Encoder-Decoder RNN



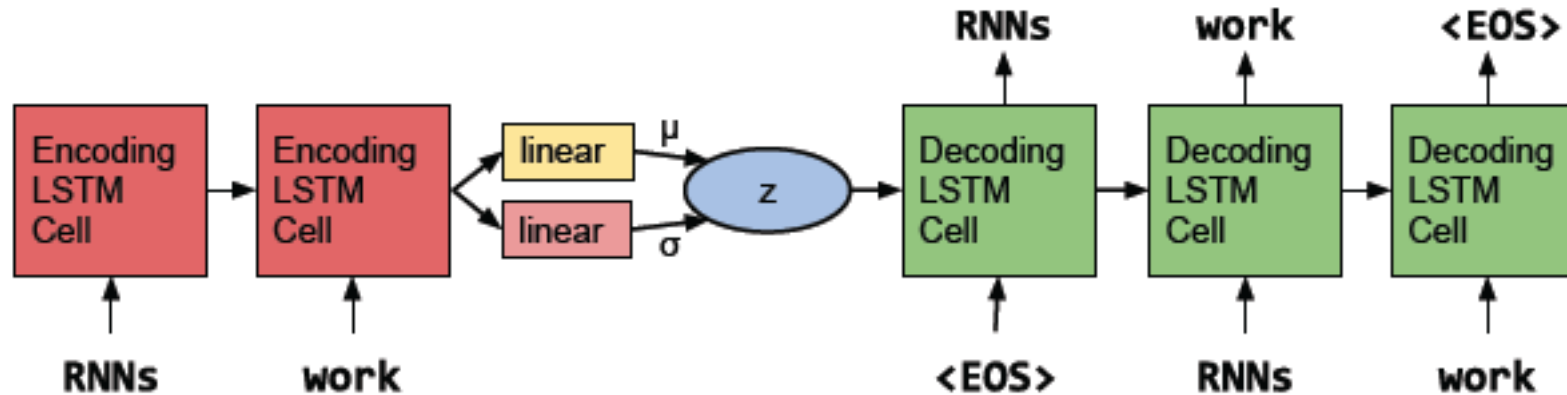
- [Iulian Vlad Serban](#), [Alessandro Sordani](#), [Ryan Lowe](#), [Laurent Charlin](#), [Joelle Pineau](#), [Aaron Courville](#), [Yoshua Bengio](#), **A Hierarchical Latent Variable Encoder-Decoder Model for Generating Dialogues**, 2016/05/20

Variational RNN LM

i went to the store to buy some groceries .
i store to buy some groceries .
i were to buy any groceries .
horses are to buy any groceries .
horses are to buy any animal .
horses the favorite any animal .
horses the favorite favorite animal .
horses are my favorite animal .

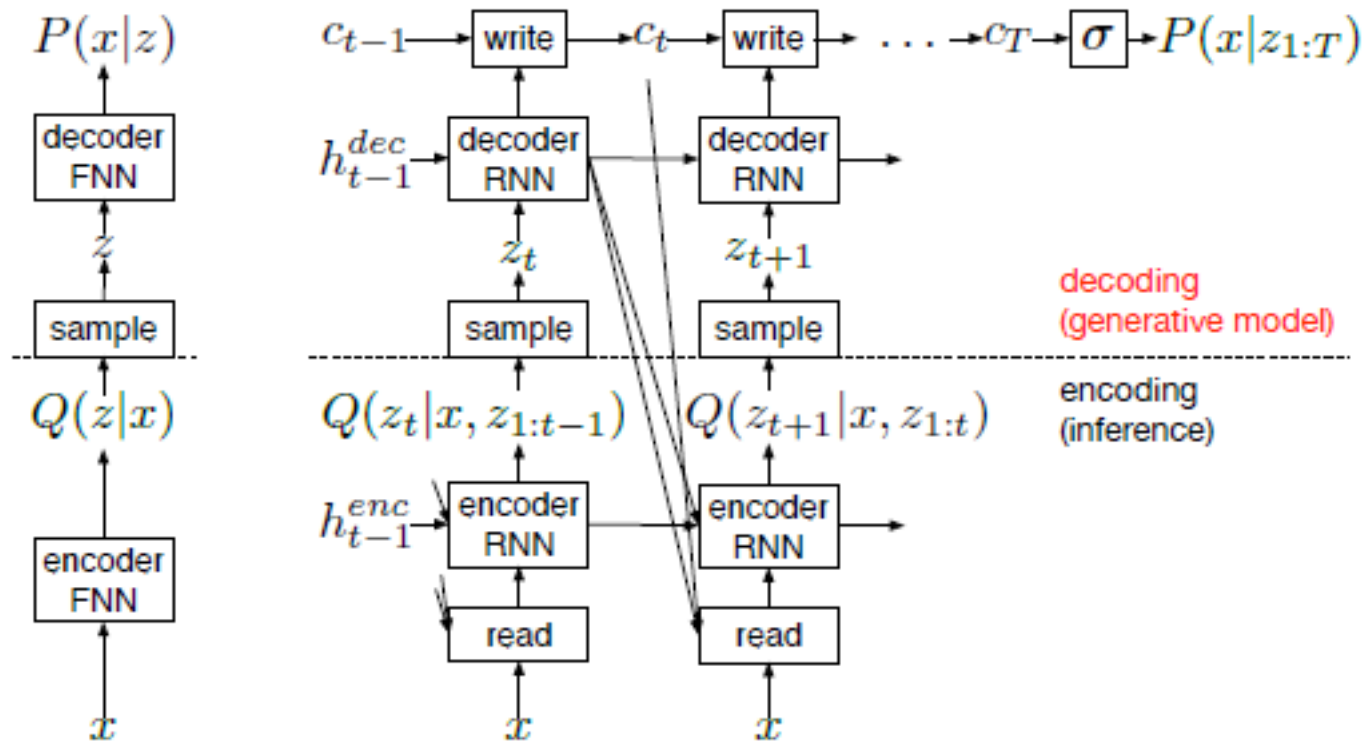
“ i want to talk to you . ”
“i want to be with you . ”
“i do n't want to be with you . ”
i do n't want to be with you .
she did n't want to be with him .

he was silent for a long moment .
he was silent for a moment .
it was quiet for a moment .
it was dark and cold .
there was a pause .
it was my turn .



- Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A. M., Jozefowicz, R., and Bengio, S. (2015). Generating sentences from a continuous space. arXiv:1511.06349.

Variational image generation



$$\mathcal{L}^z = \sum_{t=1}^T KL(Q(Z_t|h_t^{enc})||P(Z_t))$$

$$\mathcal{L}^x = -\log D(x|c_T)$$

$$\mathcal{L} = \langle \mathcal{L}^x + \mathcal{L}^z \rangle_{z \sim Q}$$

- DRAW: A Recurrent Neural Network For Image Generation

Content

- What are they?
- Marriage 1: Variational AE
- Marriage 2: Denoise AE

Denoise AE

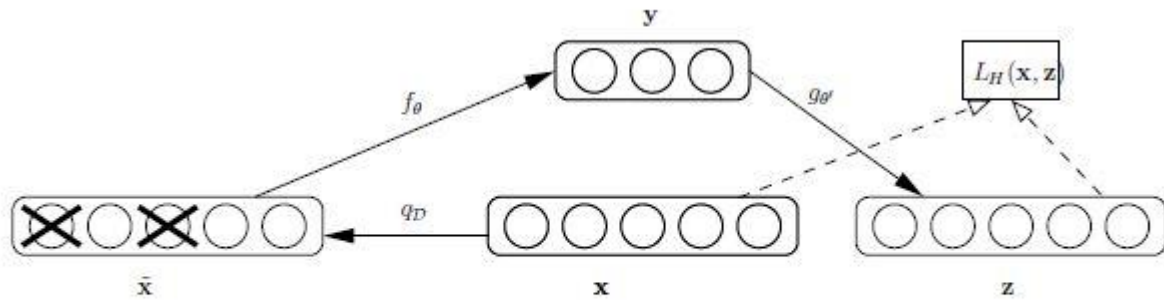
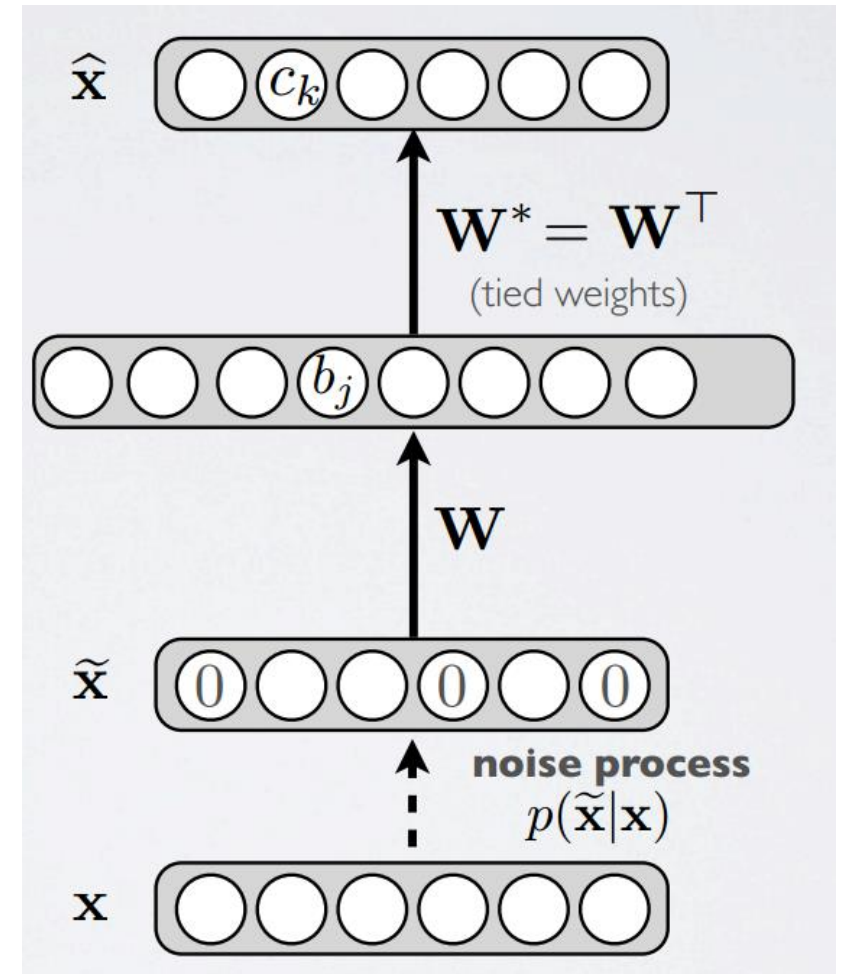
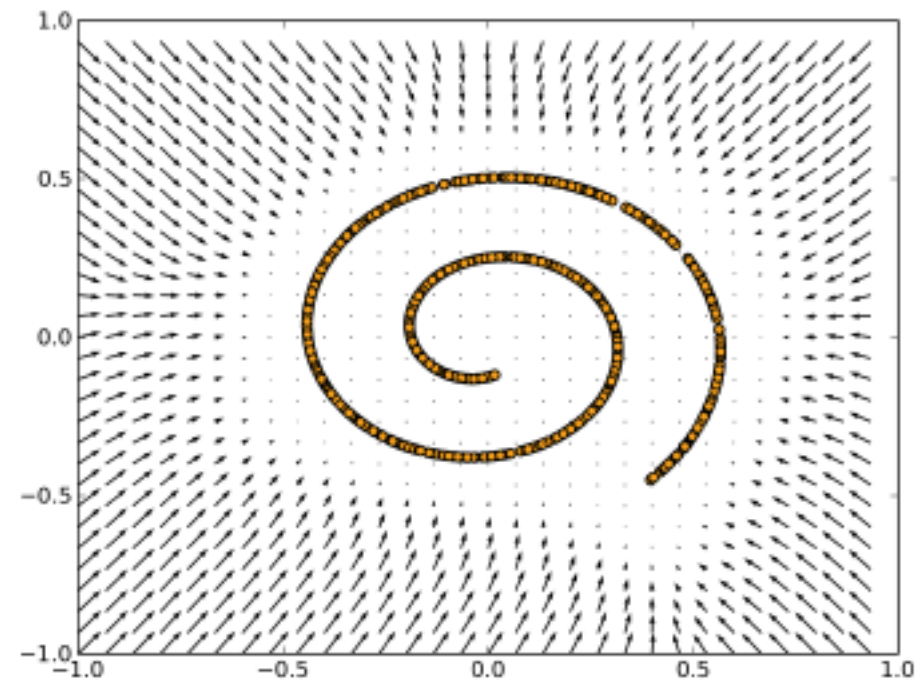
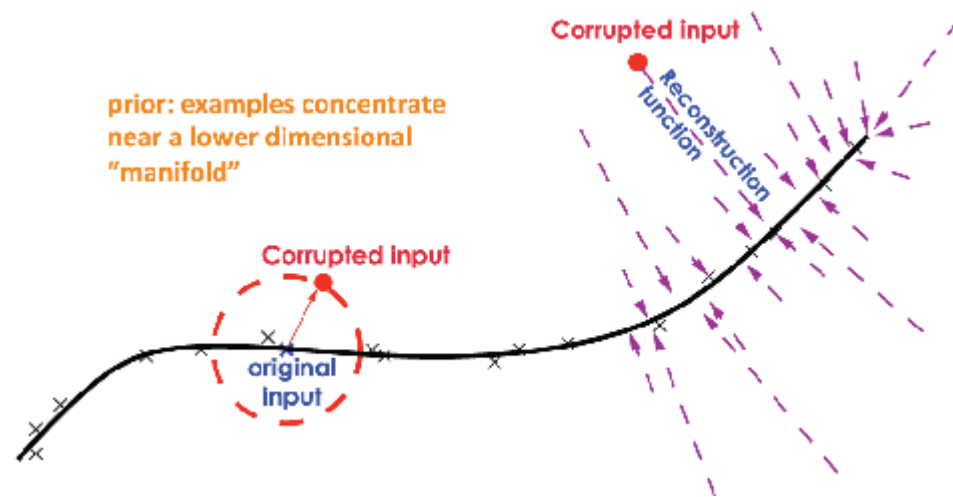


Figure 1. An example \mathbf{x} is corrupted to $\tilde{\mathbf{x}}$. The autoencoder then maps it to \mathbf{y} and attempts to reconstruct \mathbf{x} .

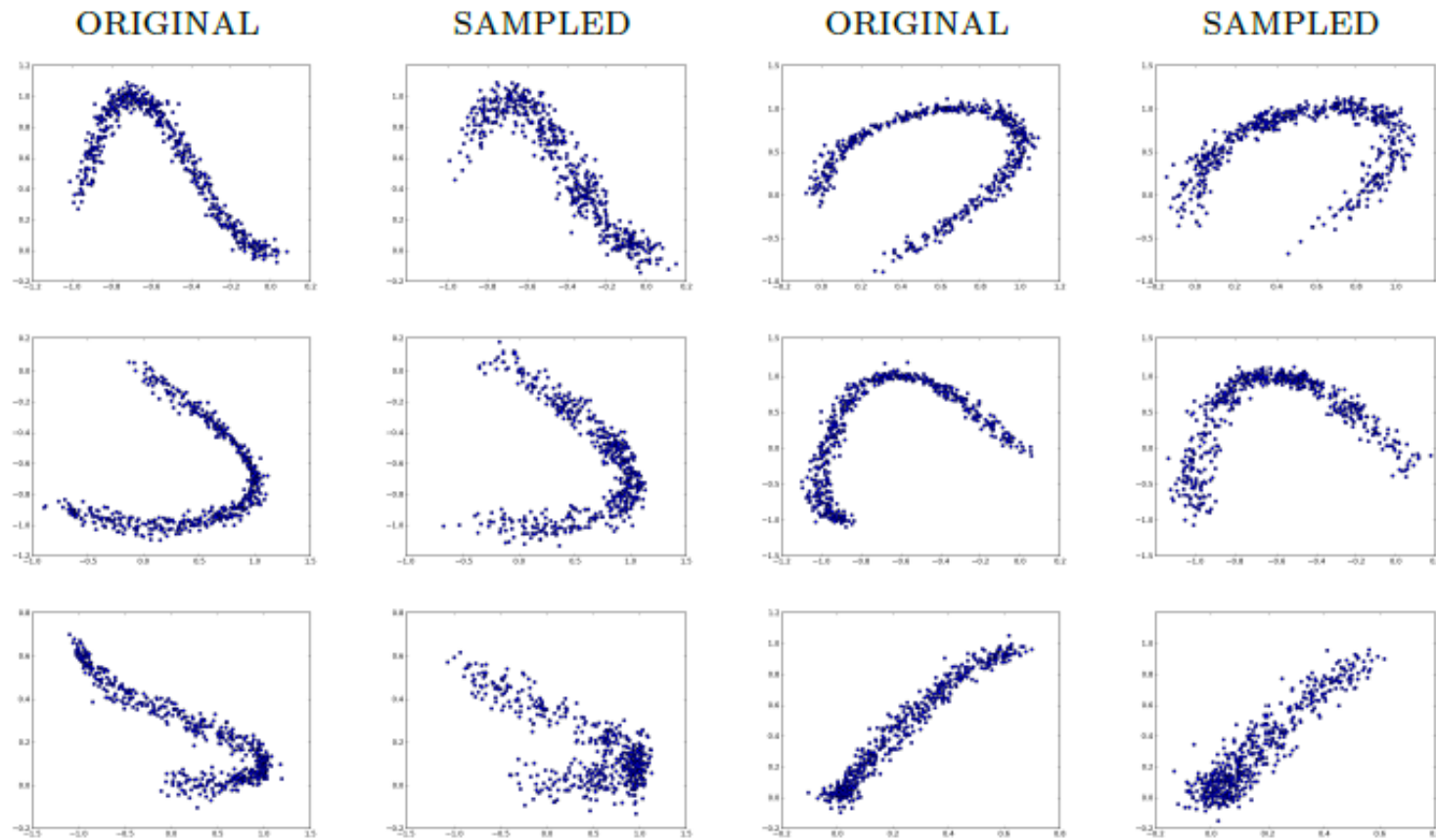


DAE learns scores (gradients)

$$r_{\sigma}^*(x) = x + \sigma^2 \frac{\partial \log p(x)}{\partial x} + o(\sigma^2) \quad \text{as } \sigma \rightarrow 0.$$

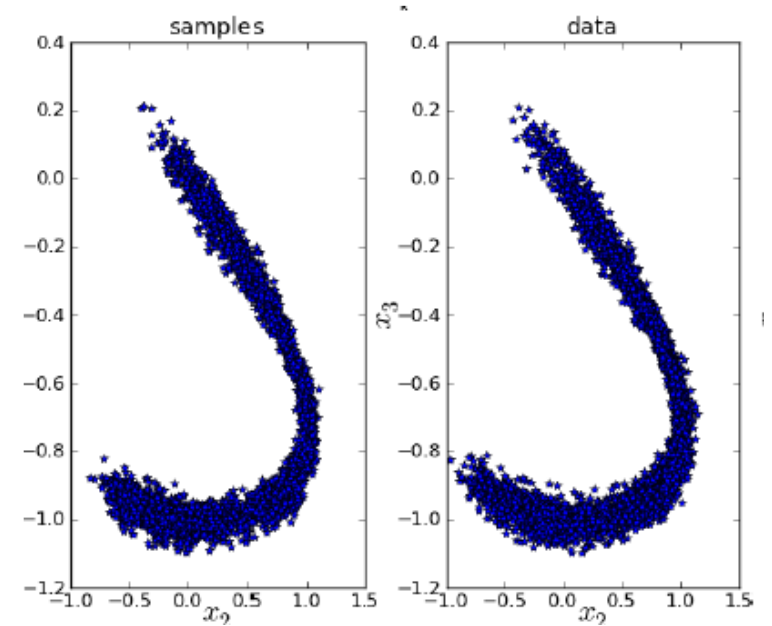
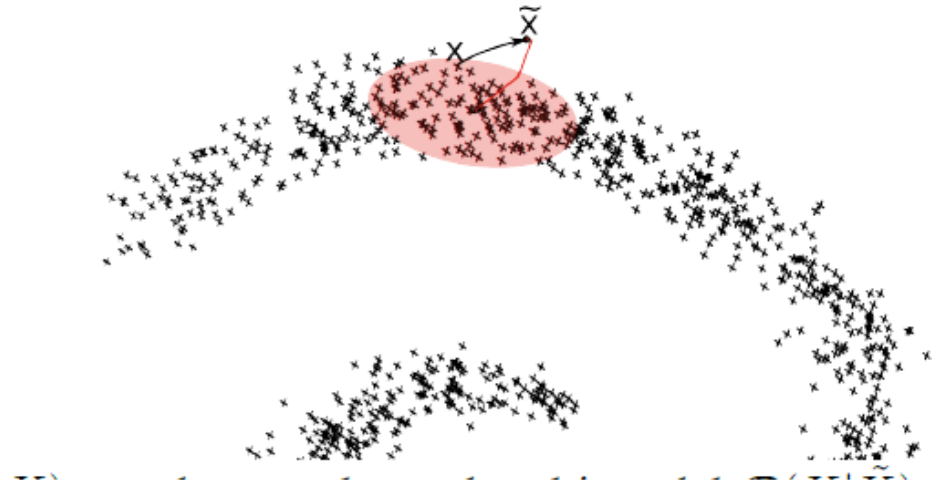


DAE can be used to sampling x



Any corruption + any cost

$$X_t \sim P_\theta(X|\tilde{X}_{t-1})$$
$$\tilde{X}_t \sim \mathcal{C}(\tilde{X}|X_t)$$

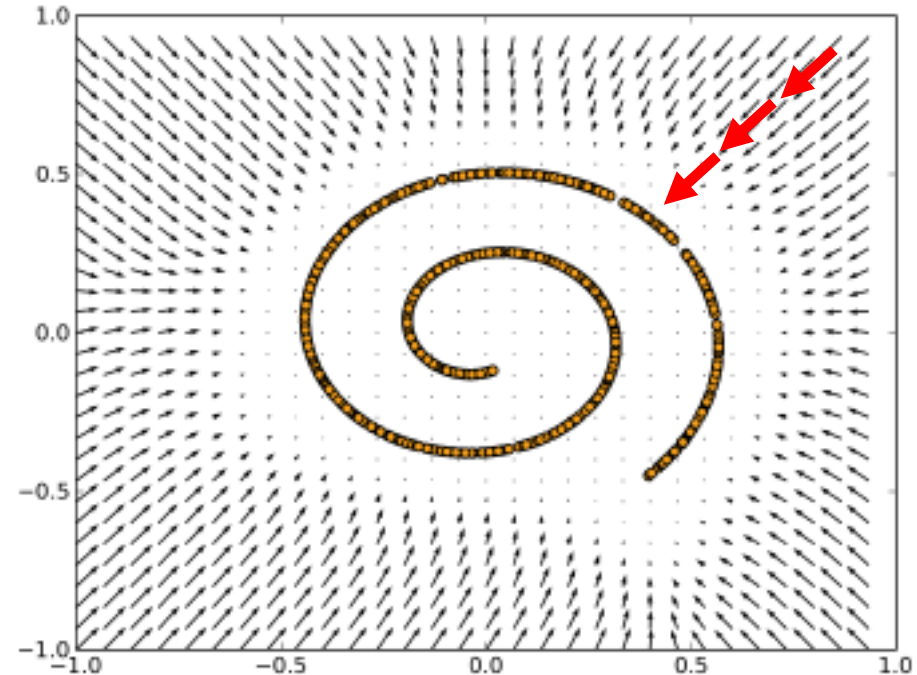


Theorem 1. *If $P_{\theta_n}(X|\tilde{X})$ is a consistent estimator of the true conditional distribution $\mathcal{P}(X|\tilde{X})$ and T_n defines an ergodic Markov chain, then as the number of examples $n \rightarrow \infty$, the asymptotic distribution $\pi_n(X)$ of the generated samples converges to the data-generating distribution $\mathcal{P}(X)$.*

- Yoshua Bengio, Li Yao, Guillaume Alain, and Pascal Vincent, Generalized Denoising Auto-Encoders as Generative Models.

Multi-step generation

- Train DAE with random corruption
- Reconstruct iteratively until converge
- Equals to get stuck to minimum energy, or $\max p(x)$
- It can be proved that with symmetric corruption, the convergence is a stationary point.



Conclusions

- Graphical model and neural model are merging
- Both variational AE and denoise AE seem reasonable to recover data distribution
- If we treat variational AE as a corruption in the encoding phase, then seems it is a special denoise DAE. Is that true?
- How other regulations can be included in both training and decoding, e.g., rythm.