

Neural Machine Translation with V-AE

Baosong Yang

2018.03.12

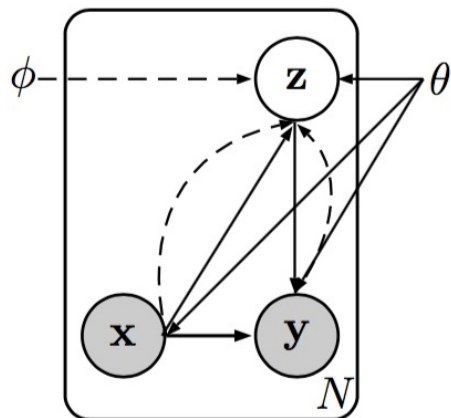
Contents

- ▶ Variational Neural Machine Translation
 - ▶ Motivations
 - ▶ Model
 - ▶ Experiments
- ▶ Variational Recurrent Neural Machine Translation
 - ▶ Motivations
 - ▶ Model
 - ▶ Experiments
- ▶ Conclusions

Variational Neural Machine Translation

Motivations

- ▶ The semantic representations are learned in an implicit way.
- ▶ Far from being sufficient for capturing all semantic details and dependencies.
- ▶ **Solution:** Incorporates the latent random variable z into NMT



$$p(\mathbf{y}|\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{y}, \mathbf{z}|\mathbf{x}) d_{\mathbf{z}} = \int_{\mathbf{z}} p(\mathbf{y}|\mathbf{z}, \mathbf{x}) p(\mathbf{z}|\mathbf{x}) d_{\mathbf{z}}$$

- ▶ Solid lines:

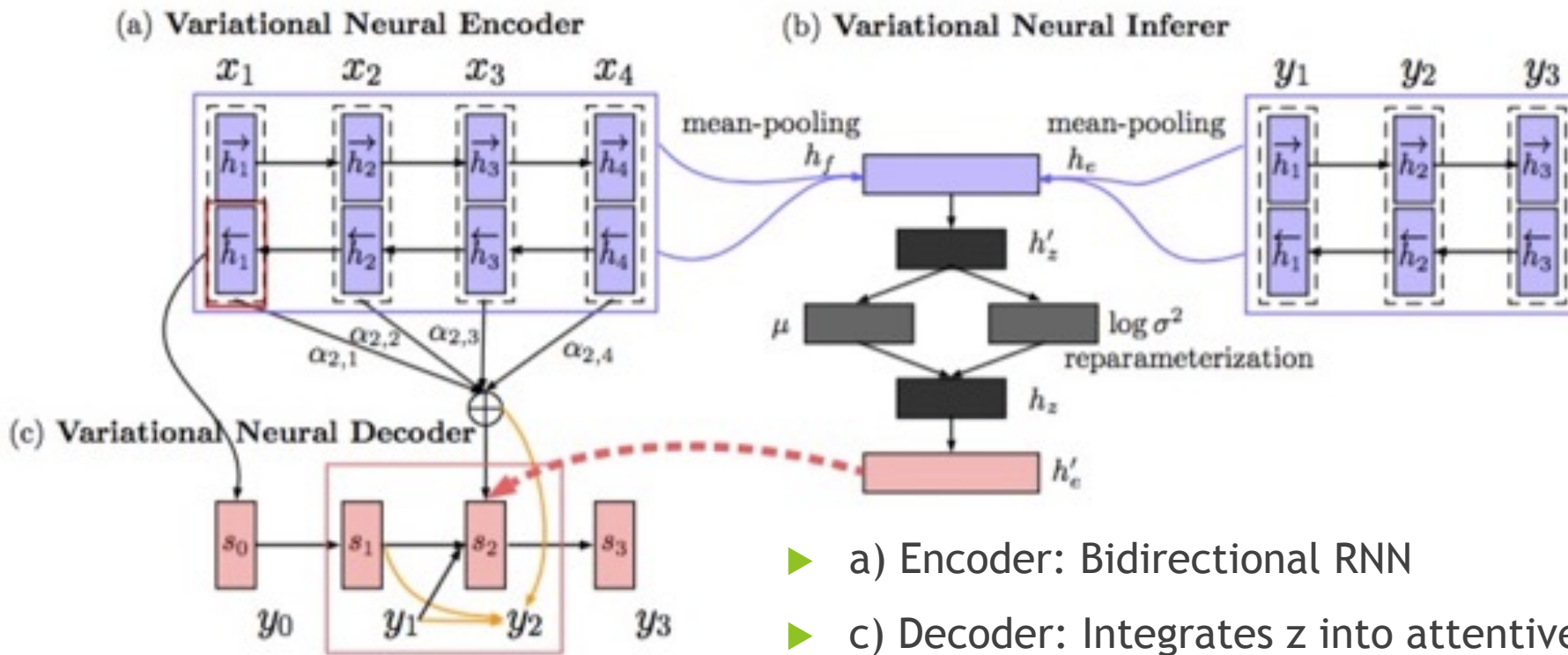
- ▶ Generative model $p_{\theta}(\mathbf{z}|\mathbf{x}) p_{\theta}(\mathbf{y}|\mathbf{z}, \mathbf{x})$

- ▶ Dashed lines:

- ▶ Variational approximation $q_{\phi}(\mathbf{z}|\mathbf{x})$
- ▶ Intractable posterior $p(\mathbf{z}|\mathbf{x}, \mathbf{y})$

Variational Neural Machine Translation

Model



Variational Neural Machine Translation

Variational Neural Inferer

- Infers the representation of \mathbf{z} according to

- 1) the learned source representations; $p_{\theta}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu'(\mathbf{x}), \sigma'(\mathbf{x})^2 \mathbf{I})$

- 2) target ones.

$$q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y}) = \mathcal{N}(\mathbf{z}; \mu(\mathbf{x}, \mathbf{y}), \sigma(\mathbf{x}, \mathbf{y})^2 \mathbf{I})$$

- Obtain \mathbf{z} :

- **Reparameterization** reparameterizes \mathbf{z} as a function of μ and σ , rather than using the standard sampling method.

$$\mathbf{h}_z = \mu + \sigma \odot \epsilon, \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

- Variational Lower Bound

-> Training: Monte Carlo ($L = 1$)

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

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}, \mathbf{y})} [\cdot] \simeq \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{y}|\mathbf{x}, \mathbf{h}_z^{(l)})$$

Variational Neural Machine Translation

Experiments

► 2.9M LDC Zh-En

System	MT05	MT02	MT03	MT04	MT06	MT08	AVG
<i>Moses</i>	33.68	34.19	34.39	35.34	29.20	22.94	31.21
<i>GroundHog</i>	31.38	33.32	32.59	35.05	29.80	22.82	30.72
<i>VNMT w/o KL</i>	31.40	33.50	32.92	34.95	28.74	22.07	30.44
<i>VNMT</i>	32.25	34.50 ⁺⁺	33.78 ⁺⁺	36.72 ⁺⁺⁺	30.92 ⁺⁺⁺	24.41 ⁺⁺⁺	32.07

► 4.5M WMT14 En-De (Evaluated on newstest14)

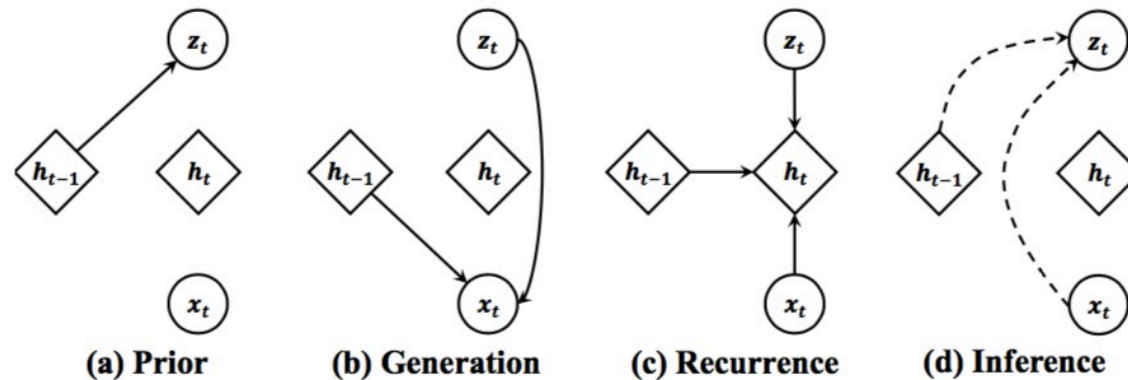
System	Architecture	BLEU
<i>Existing end-to-end NMT systems</i>		
Jean et al. (2015)	RNNSearch	16.46
Jean et al. (2015)	RNNSearch + unk replace	18.97
Jean et al. (2015)	RNNsearch + unk replace + large vocab	19.40
Luong et al. (2015a)	LSTM with 4 layers + dropout + local att. + unk replace	20.90
<i>Our end-to-end NMT systems</i>		
<i>this work</i>	RNNSearch	16.40
	VNMT	17.13 ⁺⁺
	VNMT + unk replace	19.58 ⁺⁺

Variational Recurrent Neural Machine Translation

Motivations

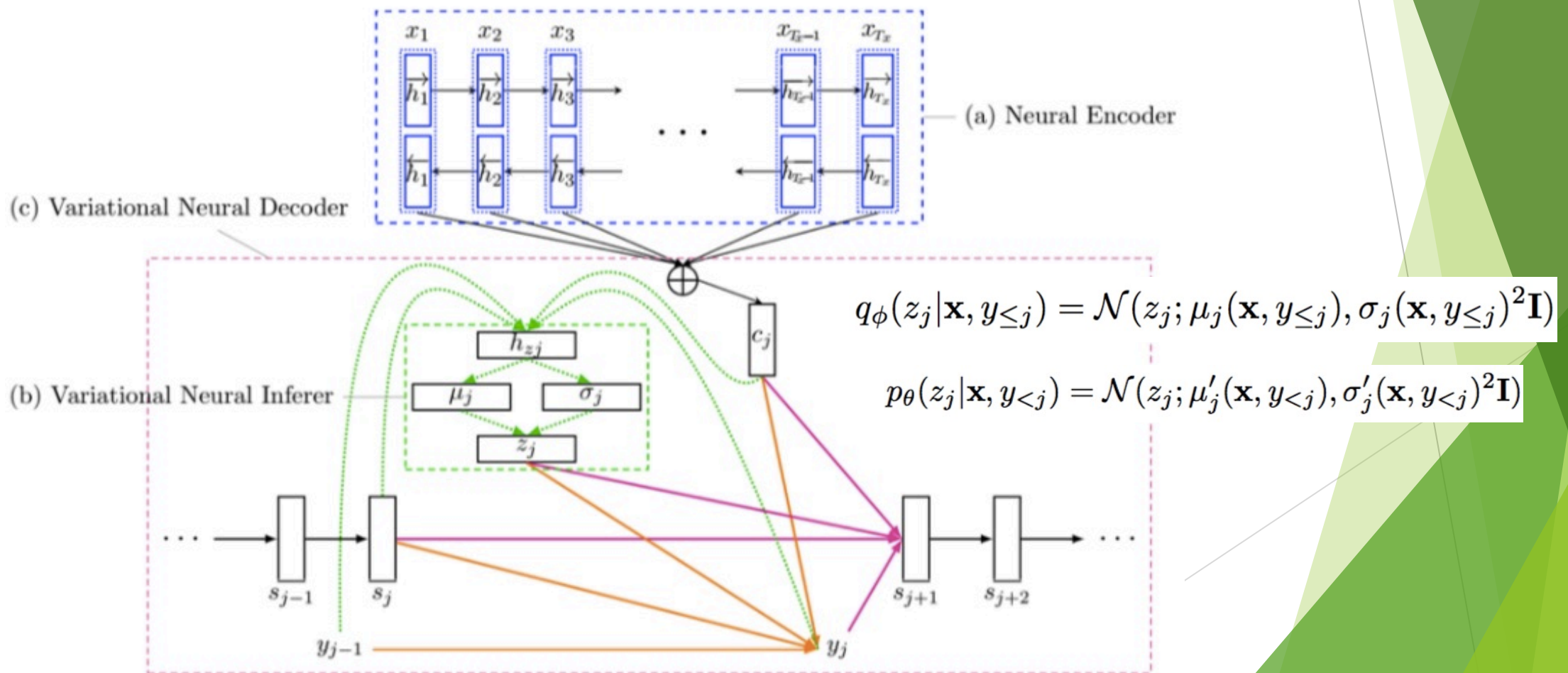
- ▶ The internal transition structure of VNMT is entirely deterministic, and hence, this implementation may not be an effective way to model high variability
- ▶ **Solution:** VRNMT is adopted continuous latent random variable sequence.

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \prod_{j=1}^{T_y} p(y_j|\mathbf{x}, y_{<j}) = \prod_{j=1}^{T_y} \int_{z_j} p(y_j, z_j|\mathbf{x}, y_{<j}) dz_j \\ &= \prod_{j=1}^{T_y} \int_{z_j} p(y_j|\mathbf{x}, y_{<j}, z_j) p(z_j|\mathbf{x}, y_{<j}) dz_j \end{aligned}$$



Variational Recurrent Neural Machine Translation

Model



Variational Recurrent Neural Machine Translation

Experiments

► 1.25M LDC Zh-En

System	MT03	MT04	MT05	MT06	Ave.
COVERAGE	34.49	38.34	34.91	34.25	35.50
MemDec	35.09	37.73	35.53	34.32	35.67
DeepLAU	36.16	39.81	35.91	35.98	36.97
DMAAen	38.33	40.11	36.71	35.29	37.61
Moses	32.93	34.76	31.31	31.05	32.51
DL4MT	36.59	39.57	35.56	35.29	36.75
VNMT	37.23	40.32	36.28	35.73	37.39
VRNMT(-TD)	36.97	40.07	36.13	35.49	37.17
VRNMT	38.08* ₊₊	41.07** ₊₊	36.82** ₊₊	36.72* ₊₊	38.17

► 4.46M WMT14 En-De (Evaluated on newstest15)

System	BLEU
BPEChar	23.9
RecAtten	25.0
ConvEncoder	24.2
Moses	20.54
DL4MT	24.88
VNMT	25.49
VRNMT(-TD)	25.34
VRNMT	25.93* ₊₊

Conclusions

► Discussions:

- The papers improved NMT models with VAE to capture the underlying semantics of sentence pairs.
- How to better exploit latent variables may be the future direction under this category.
- An alternative solution for the post-edit NMT (using latent variable to model the edit vector)

► Shortages:

- Approach 1 can not model the variants in each time step.
- Approach 2 can not model the future context of the target.