

Variational Attention for Sequence-to-Sequence Models

Hareesh Bahuleyan,¹ Lili Mou,¹ Olga Vechtomova, Pascal Poupart

University of Waterloo

March, 2018

Bahuleyan, Mou, Vechtomova, Poupart

Variational Attention for Seq2Seq Models



< ロ > < 同 > < 回 > < 回 > < 回

1 Introduction

- 2 Background & Motivation
- **3** Variational Attention

4 Experiments

5 Conclusion

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models



1 Introduction

- 2 Background & Motivation
- **3** Variational Attention

4 Experiments

5 Conclusion

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seq2Seq Models



Variational Autoencoder (VAE)

Kingma and Welling (2013)

A combination of (neural) autoencoders and variational inference

- Compared with traditional variational inference
 Takes use of neural networks as a powerful density estimator
- Compared with traditional autoencoders
 Imposes a probabilistic distribution on latent representations

(日) (同) (三) (三)

- Learns the distribution of data
 - Instead of the MAP sample
 - Implicitly capture the diversity of data
- Generating samples from scratch
 - Image generation, sentence generation, etc.
 - Controlling the generated samples
 - \Rightarrow Not quite feasible in deterministic AE
- Regularization



Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models



Variational Encoder-Decoder (VED)

- Encoder-decoder frameworks Machine translation Dialog systems Summarization
- VAE \Rightarrow VED (Variational encoder-decoder)

<ロト < 囲 > < 国 > < 国 > < 国 >

- Seq2Seq models as encoders & decoders
- Attention mechanism serving as dynamic alignment

However, we observe the bypassing phenomenon



Image: A matrix of the second seco

U Waterloo

Bahuleyan, Mou, Vechtomova, Poupart

Variational attention mechanism

- Modeling attention probability/vector as random variables
- Imposing some prior and posterior distributions over attention



• Experimental results show that the variational space is more effective when combined with variational attention.

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models

1 Introduction

- 2 Background & Motivation
- **3** Variational Attention

4 Experiments

5 Conclusion

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models



- Z Hidden variables
- Y Observable variables



$$\log p_{\boldsymbol{\theta}}(\boldsymbol{y}^{(n)}) = \mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{y}^{(n)})} \left[\log \left\{ \frac{p_{\boldsymbol{\theta}}(\boldsymbol{y}^{(n)}, \boldsymbol{z})}{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{y}^{(n)})} \right\} \right] \\ + \operatorname{KL}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{y}^{(n)}) \| p(\boldsymbol{z}|\boldsymbol{y})) \\ \geq \mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{y}^{(n)})} \left[\log p_{\boldsymbol{\theta}}(\boldsymbol{y}^{(n)}|\boldsymbol{z}) \right] \\ - \operatorname{KL}\left(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{y}^{(n)}) \| p(\boldsymbol{z}) \right) \stackrel{\Delta}{=} \mathcal{L}^{(n)}(\boldsymbol{\theta}, \boldsymbol{\phi})$$

 $loss = \mathbb{E}[reconstruction \ loss] + KL \ divergence$

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models <ロト </p>

$\mathsf{loss} = \mathbb{E}[\mathsf{reconstruction} \ \mathsf{loss}] + \mathsf{KL} \ \mathsf{divergence}$

Prior

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

Posterior

$$q_{\phi}(\boldsymbol{z}|\boldsymbol{y}) = \mathcal{N}(\boldsymbol{\mu}, \operatorname{diag}\{\boldsymbol{\sigma}\})$$

where $oldsymbol{\mu}, oldsymbol{\sigma} = \mathrm{NN}(oldsymbol{y})$

Bahuleyan, Mou, Vechtomova, Poupart

Variational Attention for Seq2Seq Models



・ロト ・聞 ト ・ ヨト ・ ヨト

- Encoder-Decoder: Transforming X to Y
- Attempt#1: Condition any distribution on X (Zhang et al., 2016; Cao and Clark, 2017; Serban et al., 2017) $\log p_{\theta}(\boldsymbol{y}^{(n)}) \geq \mathbb{E}_{\boldsymbol{z} \sim q_{\phi}(\boldsymbol{z} | \boldsymbol{y}^{(n)})} \left[\log \left\{ \frac{p_{\theta}(\boldsymbol{y}^{(n)}, \boldsymbol{z})}{q_{\phi}(\boldsymbol{z} | \boldsymbol{y}^{(n)})} \right\} \right]$ $= \mathbb{E}_{\boldsymbol{z} \sim q_{\phi}(\boldsymbol{z} | \boldsymbol{y}^{(n)})} \left[\log p_{\theta}(\boldsymbol{y}^{(n)} | \boldsymbol{z}) \right]$ $- \operatorname{KL} \left(q_{\phi}(\boldsymbol{z} | \boldsymbol{y}^{(n)}) \| p(\boldsymbol{z}) \right) \triangleq \mathcal{L}^{(n)}(\boldsymbol{\theta}, \boldsymbol{\phi})$
- Doubts:
 - The posterior contains $oldsymbol{Y}$
 - Cannot have fine-grained (word/character-level) variational modeling

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

U Waterloo

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seq2Seq Models

Image: A matrix of the second seco

Variational Encoder-Decoder

- Encoder-Decoder: Transforming X to Y
- Attempt#2 (Cao and Clark, 2017): Assuming \boldsymbol{Y} is some function of \boldsymbol{X} i.e., $\boldsymbol{Y} = \boldsymbol{Y}(\boldsymbol{X})$, then $q(\boldsymbol{z}|\boldsymbol{y}) = q(\boldsymbol{z}|\boldsymbol{Y}(\boldsymbol{x})) = \widetilde{q}(\boldsymbol{z}|\boldsymbol{x})$



Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models



At each step of decoding

Attention probability

$$\alpha_{ji} = \frac{\exp\{\widetilde{\alpha}_{ji}\}}{\sum_{i'=1}^{|\boldsymbol{x}|} \exp\{\widetilde{\alpha}_{ji'}\}}$$

Attention vector

$$oldsymbol{a}_j = \sum_{i=1}^{|oldsymbol{x}|} lpha_{ji} oldsymbol{h}_i^{(extsf{src})}$$

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models



イロト イ団ト イヨト イヨト

At each step of decoding

Attention probability

$$\alpha_{ji} = \frac{\exp\{\widetilde{\alpha}_{ji}\}}{\sum_{i'=1}^{|\boldsymbol{x}|} \exp\{\widetilde{\alpha}_{ji'}\}}$$

Attention vector

$$oldsymbol{a}_j = \sum_{i=1}^{|oldsymbol{x}|} lpha_{ji} oldsymbol{h}_i^{(extsf{src})}$$



U Waterloo

Bahuleyan, Mou, Vechtomova, Poupart

Bypassing Phenomenon

A deterministic connection bypasses variational space

Hidden state initialization:



Input: the men are playing musical instruments

(a) VAE w/o hidden state init. (Avg entropy: 2.52)

the men are playing musical instruments the men are playing video games the musicians are playing musical instruments the women are playing musical instruments

(b) VAE w/ hidden state init. (Avg entropy: 2.01)

the men are playing musical instruments the men are playing musical instruments the men are playing musical instruments the man is playing musical instruments

Bahuleyan, Mou, Vechtomova, Poupart

Variational Attention for Seq2Seq Models

1 Introduction

- 2 Background & Motivation
- **3** Variational Attention

4 Experiments

5 Conclusion

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models





General idea: Model attention as random variables





U Waterloo

イロト イ団ト イヨト イヨト

Bahuleyan, Mou, Vechtomova, Poupart

$$\begin{split} \mathcal{L}_{j}^{(n)}(\boldsymbol{\theta}, \boldsymbol{\phi}) \\ = & \mathbb{E}_{\boldsymbol{z}, \boldsymbol{a} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}, \boldsymbol{a} | \boldsymbol{x}^{(n)})} \left[\log p_{\boldsymbol{\theta}}(\boldsymbol{y}^{(n)} | \boldsymbol{z}, \boldsymbol{a}) \right] - \mathrm{KL} \left(q_{\boldsymbol{\phi}}(\boldsymbol{z}, \boldsymbol{a} | \boldsymbol{x}^{(n)}) \| p(\boldsymbol{z}, \boldsymbol{a}) \right) \\ = & \mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{\phi}}^{(z)}(\boldsymbol{z} | \boldsymbol{x}^{(n)}), \boldsymbol{a} \sim q_{\boldsymbol{\phi}}^{(a)}(\boldsymbol{a} | \boldsymbol{x}^{(n)})} \left[\log p_{\boldsymbol{\theta}}(\boldsymbol{y}^{(n)} | \boldsymbol{z}, \boldsymbol{a}) \right] \\ & - \mathrm{KL} \left(q_{\boldsymbol{\phi}}^{(z)}(\boldsymbol{z} | \boldsymbol{x}^{(n)}) \| p(\boldsymbol{z}) \right) - \mathrm{KL} \left(q_{\boldsymbol{\phi}}^{(a)}(\boldsymbol{a} | \boldsymbol{x}^{(n)}) \| p(\boldsymbol{a}) \right) \end{split}$$

・ロト ・聞 ト ・ ヨト ・ ヨト

U Waterloo



Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seq2Seq Models

・ロト ・聞 ト ・ ヨト ・ ヨト

U Waterloo

Standard norm:
$$p(\boldsymbol{a}_j) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Norm centered at $ar{h}^{(
m src)}$:

$$p(m{a}_j) = \mathcal{N}(m{ar{h}}^{(extsf{src})}, \mathbf{I})$$
 where $ar{m{h}}^{(extsf{src})} = rac{1}{|m{x}|}\sum_{i=1}^{|m{x}|}m{h}^{(extsf{src})}_i$

Bahuleyan, Mou, Vechtomova, Poupart

Introduction	Background & Motivation	Variational Attention	Experiments	Conclusion	References
Posteri	or				

Attention probability

$$\alpha_{ji} = \frac{\exp\{\widetilde{\alpha}_{ji}\}}{\sum_{i'=1}^{|\boldsymbol{x}|} \exp\{\widetilde{\alpha}_{ji'}\}}$$

Attention vector

$$oldsymbol{a}_j^{(\mathsf{det})} = \sum_{i=1}^{|oldsymbol{x}|} lpha_{ji} oldsymbol{h}_i^{(\mathsf{src})}$$

A B > 4
 B > 4
 B

3 D (3 D

U Waterloo

Posterior:
$$\mathcal{N}(\mu_{a_j}, \sigma_{a_j})$$
, where
$$\mu_{a_j} \equiv a_j^{\mathsf{det}}, \qquad \sigma_{a_j} = \mathrm{NN}(a_j^{\mathsf{det}})$$

Bahuleyan, Mou, Vechtomova, Poupart

$$J^{(n)}(\boldsymbol{\theta}, \boldsymbol{\phi}) = J_{\text{rec}}(\boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{y}^{(n)}) \\ + \lambda_{\text{KL}} \Big[\text{KL} \left(q_{\boldsymbol{\phi}}^{(z)}(\boldsymbol{z}) \| p(\boldsymbol{z}) \right) \\ + \gamma_{a} \sum_{j=1}^{|\boldsymbol{y}|} \text{KL} \left(q_{\boldsymbol{\phi}}^{(a)}(\boldsymbol{a}_{j}) \| p(\boldsymbol{a}_{j}) \right) \Big]$$

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seq2Seq Models U Waterloo

イロト イヨト イヨト イヨン

Geometric Interpretation



Bahuleyan, Mou, Vechtomova, Poupart

Variational Attention for Seq2Seq Models



イロト イヨト イヨト イヨン

1 Introduction

- 2 Background & Motivation
- **3** Variational Attention

4 Experiments

5 Conclusion

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models





Question generation (Du et al., 2017)

- Given some information (a sentence), to generate a related question
- Dataset from the Stanford Question Answering Dataset (Rajpurkar et al., 2016, SQuAD)

< ロ > < 同 > < 回 > < 回 > < 回

- KL-annealing: logistic schema
- Word dropout: 25%

Hyperparaemters tuned on VAE and adopted to all VED models

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seq2Seq Models



Overall Performance

Model	Inference	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Entropy	Dist-1	Dist-2
DED (w/o Attn) Du et al. (2017)	MAP	31.34	13.79	7.36	4.26	-	-	-
DED (w/o Attn) DED+DAttn	MAP MAP	29.31 30.24	12.42 14.33	6.55 8.26	3.61 4.96	-	-	-
VED+DAttn	MAP Sampling	31.02 30.87	14.57 14.71	8.49 8.61	5.02 5.08	- 2.214	- 0.132	- 0.176
VED+DAttn (2-stage training)	MAP Sampling	28.88 29.25	13.02 13.21	7.33 7.45	4.16 4.25	- 2.241	- 0.140	- 0.188
VED+VAttn-0	MAP Sampling	29.70 30.22	14.17 14.22	8.21 8.28	4.92 4.87	2.320	- 0.165	- 0.231
VED+VAttn- \bar{h}	MAP Sampling	30.23 30.47	14.30 14.35	8.28 8.39	4.93 4.96	- 2.316	- 0.162	- 0.228

U Waterloo

Bahuleyan, Mou, Vechtomova, Poupart





< □ > < □ > < □

3

U Waterloo

Bahuleyan, Mou, Vechtomova, Poupart

$$J^{(n)}(\boldsymbol{\theta}, \boldsymbol{\phi}) = J_{\mathsf{rec}}(\boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{y}^{(n)}) \\ + \lambda_{\mathrm{KL}} \Big[\mathrm{KL} \left(q_{\boldsymbol{\phi}}^{(z)}(\boldsymbol{z}) \| p(\boldsymbol{z}) \right) \\ + \gamma_{\boldsymbol{a}} \sum_{j=1}^{|\boldsymbol{y}|} \mathrm{KL} \left(q_{\boldsymbol{\phi}}^{(a)}(\boldsymbol{a}_{j}) \| p(\boldsymbol{a}_{j}) \right) \Big]$$



Bahuleyan, Mou, Vechtomova, Poupart

U Waterloo

ntroduc	спон Баскугоц			Experiments	Conclusion	References
Cas	se Study					
	Source	when the british	forces evacuated at the	e close of the wa	ar in 1783 ,	
	Reference	they transported in what year did				
	how many people evacuated in newfoundland ? VED+DAttn how many people evacuated in newfoundland ? what did the british forces seize in the war ?					
	$VED+Vattn-ar{h}$	how many peop where did the br when did the br	le lived in nova scotia ? ritish forces retreat ? itish forces leave the wa	r ?		
	Source	downstream , m mianyang by jun how many peop	ore than 200,000 people the 1 in anticipation of the le were evacuated down.	e were evacuated he dam bursting stream ?	d from	
	VED+DAttn	how many peop how many peop how many peop	le evacuated from the n le evacuated from the n le evacuated from the n	nianyang basin ? nianyang basin ? nianyang basin ?	2 2 2	
	$VED+VAttn_{ar{h}}$	how many peop how many peop how many peop	le evacuated from the tu le evacuated from the d le were evacuated from	unnel ? am ? fort in the dam	?	T 200
Bahuley	an, Mou, Vechtomo	va, Poupart				_ ્ર્ગ્લ્લ U Waterloo

Bahuleyan, Mou, Vechtomova, Poupart

1 Introduction

- 2 Background & Motivation
- **3** Variational Attention

4 Experiments

5 Conclusion

Bahuleyan, Mou, Vechtomova, Poupart Variational Attention for Seg2Seg Models



Our work

- Address the bypassing phenomenon
- Propose a variational attention

Future work

Probabilistic modeling of attention probability

Lesson learned

Design philosophy of VAE/VED



< ロ > < 同 > < 回 > < 回 > < 回

Kris Cao and Stephen Clark. 2017. Latent variable dialogue models and their diversity. In *EACL*. pages 182–187.

Xinya Du, Junru Shao, and Claire Cardie. 2017. Learning to ask: Neural question generation for reading comprehension. In *ACL*. pages 1342–1352.

Diederik P Kingma and Max Welling. 2013. Auto-encoding variational Bayes. *arXiv* preprint arXiv:1312.6114.

- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *EMNLP*. pages 2383–2392.
- Iulian Vlad Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron C Courville, and Yoshua Bengio. 2017. A hierarchical latent variable encoder-decoder model for generating dialogues. In AAAI. pages 3295–3301.

Biao Zhang, Deyi Xiong, jinsong su, Qun Liu, Rongrong Ji, Hong Duan, and Min Zhang. 2016. Variational neural discourse relation recognizer. In *EMNLP*. pages 382–391.

Bahuleyan, Mou, Vechtomova, Poupart

Introduction Background & Motivation Variational Attention Experiments Conclusion References
Thanks for listening

Question?

Bahuleyan, Mou, Vechtomova, Poupart

