

A Stable Variational Autoencoder for Text Modelling

Ruizhe Li[♣], Xiao Li[♣], Chenghua Lin[♡], Matthew Collinson[♣] and Rui Mao[♣]

[♣]Department of Computing Science, University of Aberdeen, UK

{r02r117, x.li, matthew.collinson, r03rm16}@abdn.ac.uk

[♡]Department of Computer Science, University of Sheffield, UK

c.lin@sheffield.ac.uk

Abstract

Variational Autoencoder (VAE) is a powerful method for learning representations of high-dimensional data. However, VAEs can suffer from an issue known as latent variable collapse (or KL loss vanishing), where the posterior collapses to the prior and the model will ignore the latent codes in generative tasks. Such an issue is particularly prevalent when employing VAE-RNN architectures for text modelling (Bowman et al., 2016). In this paper, we present a simple architecture called holistic regularisation VAE (HR-VAE), which can effectively avoid latent variable collapse. Compared to existing VAE-RNN architectures, we show that our model can achieve much more stable training process and can generate text with significantly better quality.

1 Introduction

Variational Autoencoder (VAE) (Kingma and Welling, 2013) is a powerful method for learning representations of high-dimensional data. However, recent attempts of applying VAEs to text modelling are still far less successful compared to its application to image and speech (Bachman, 2016; Fraccaro et al., 2016; Semeniuta et al., 2017). When applying VAEs for text modelling, recurrent neural networks (RNNs)¹ are commonly used as the architecture for both encoder and decoder (Bowman et al., 2016; Xu and Durrett, 2018; Dieng et al., 2019). While such a VAE-RNN based architecture allows encoding and generating sentences (in the decoding phase) with variable-length effectively, it is also vulnerable to an issue known as latent variable collapse (or KL loss vanishing), where the posterior collapses to the prior and the model will ignore the latent codes in generative tasks.

¹NB: here we refer RNN to any type of recurrent neural architectures including LSTM and GRU.

Various efforts have been made to alleviate the latent variable collapse issue. Bowman et al. (2016) uses KL annealing, where a variable weight is added to the KL term in the cost function at training time. Yang et al. (2017) discovered that there is a trade-off between the contextual capacity of the decoder and effective use of encoding information, and developed a dilated CNN as decoder which can vary the amount of conditioning context. They also introduced a loss clipping strategy in order to make the model more robust. Xu and Durrett (2018) addressed the problem by replacing the standard normal distribution for the prior with the von Mises-Fisher (vMF) distribution. With vMF, the KL loss only depends on the concentration parameter which is fixed during training and testing, and hence results in a constant KL loss. In a more recent work, Dieng et al. (2019) avoided latent variable collapse by including skip connections in the generative model, where the skip connections enforce strong links between the latent variables and the likelihood function.

Although the aforementioned works show effectiveness in addressing the latent variable collapse issue to some extent, they either require carefully engineering to balance the weight between the reconstruction loss and KL loss (Bowman et al., 2016; Sønderby et al., 2016), or resort to designing more sophisticated model structures (Yang et al., 2017; Xu and Durrett, 2018; Dieng et al., 2019).

In this paper, we present a simple architecture called holistic regularisation VAE (HR-VAE), which can effectively avoid latent variable collapse. In contrast to existing VAE-RNN models for text modelling which merely impose a standard normal distribution prior on the last hidden state of the RNN encoder, our HR-VAE model imposes regularisation for all hidden states of the RNN encoder. Another advantage of our model is that it

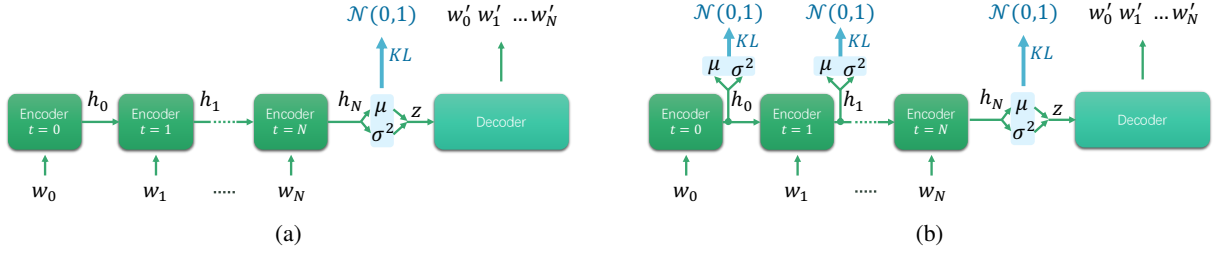


Figure 1: (a) The typical architecture of RNN-based VAE; (b) the proposed HR-VAE architecture.

is generic and can be applied to any existing VAE-RNN-based architectures.

We evaluate our model against several strong baselines which apply VAE for text modelling (Bowman et al., 2016; Yang et al., 2017; Xu and Durrett, 2018). We conducted experiments based on two public benchmark datasets, namely, the Penn Treebank dataset (Marcus and Marcinkiewicz, 1993) and the end-to-end (E2E) text generation dataset (Novikova et al., 2017). Experimental results show that our HR-VAE model not only can effectively mitigate the latent variable collapse issue with a stable training process, but also can give better predictive performance than the baselines, as evidenced by both quantitative (e.g., negative log likelihood and perplexity) and qualitative evaluation. The code for our model is available online².

2 Methodology

2.1 Background of VAE

A variational autoencoder (VAE) is a deep generative model, which combines variational inference with deep learning. The VAE modifies the conventional autoencoder architecture by replacing the deterministic latent representation \mathbf{z} of an input \mathbf{x} with a posterior distribution $P(\mathbf{z}|\mathbf{x})$, and imposing a prior distribution on the posterior, such that the model allows sampling from any point of the latent space and yet able to generate novel and plausible output. The prior is typically chosen to be standard normal distributions, i.e., $P(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{1})$, such that the KL divergence between posterior and prior can be computed in closed form (Kingma and Welling, 2013).

To train a VAE, we need to optimise the marginal likelihood $P_\theta(\mathbf{x}) = \int P(\mathbf{z})P_\theta(\mathbf{x}|\mathbf{z})d\mathbf{z}$,

where the log likelihood can take following form:

$$\log P_\theta(\mathbf{x}) = \mathcal{L}(\theta, \phi; \mathbf{x}) + \text{KL}(Q_\phi(\mathbf{z}|\mathbf{x})\|P_\theta(\mathbf{z}|\mathbf{x})) \quad (1)$$

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{Q_\phi(\mathbf{z}|\mathbf{x})}[\log P_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}(Q_\phi(\mathbf{z}|\mathbf{x})\|P(\mathbf{z})) \quad (2)$$

Here $Q_\phi(\mathbf{z}|\mathbf{x})$ is the variational approximation for the true posterior $P_\theta(\mathbf{z}|\mathbf{x})$. Specifically, $Q_\phi(\mathbf{z}|\mathbf{x})$ can be regarded as an encoder (a.k.a. the recognition model) and $P_\theta(\mathbf{x}|\mathbf{z})$ the decoder (a.k.a. the generative model). Both encoder and decoder are implemented via neural networks. As proved in (Kingma and Welling, 2013), optimising the marginal log likelihood is essentially equivalent to maximising $\mathcal{L}(\theta, \phi; \mathbf{x})$, i.e., the evidence lower bound (ELBO), which consists of two terms. The first term is the expected reconstruction error indicating how well the model can reconstruct data given a latent variable. The the second term is the KL divergence of the approximate posterior from prior, i.e., a regularisation pushing the learned posterior to be as close to the prior as possible.

2.2 Variational Autoencoder with Holistic Regularisation

In this section, we discuss the technical details of the proposed holistic regularisation VAE (HR-VAE) model, a general architecture which can effectively mitigate the KL vanishing phenomenon.

Our model design is motivated by one noticeable defect shared by the VAE-RNN based models in previous works (Bowman et al., 2016; Yang et al., 2017; Xu and Durrett, 2018; Dieng et al., 2019). That is, all these models, as shown in Figure 1a, only impose a standard normal distribution prior on the last hidden state of the RNN encoder, which potentially leads to learning a suboptimal representation of the latent variable and results in model vulnerable to KL loss vanishing. Our hypothesis is that to learn a good representation of

²<https://github.com/ruizheliUOA/HR-VAE>

data and a good generative model, it is crucial to impose the standard normal prior on all the hidden states of the RNN-based encoder (see Figure 1b), which allows a better regularisation of the model learning process.

We implement the HR-VAE model using a two-layer LSTM for both the encoder and decoder. However, one should note that our architecture can be readily applied to other types of RNN such as GRU. For each time stamp t (see Figure 1b), we concatenate the hidden state \mathbf{h}_t and the cell state \mathbf{c}_t of the encoder. The concatenation (i.e., $[\mathbf{h}_t; \mathbf{c}_t]$) is then fed into two linear transformation layers for estimating $\boldsymbol{\mu}_t$ and $\boldsymbol{\sigma}_t^2$, which are parameters of a normal distribution corresponding to the concatenation of \mathbf{h}_t and \mathbf{c}_t . Let $Q_{\phi_t}(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\mathbf{z}_t|\boldsymbol{\mu}_t, \boldsymbol{\sigma}_t^2)$, we wish $Q_{\phi_t}(\mathbf{z}_t|\mathbf{x})$ to be close to a prior $P(\mathbf{z}_t)$, which is a standard Gaussian. Finally, the KL divergence between these two multivariate Gaussian distributions (i.e., Q_{ϕ_t} and $P(\mathbf{z}_t)$) will contribute to the overall KL loss of the ELBO. By taking the average of the KL loss at each time stamp t , the resulting ELBO takes the following form

$$\begin{aligned} \mathcal{L}(\theta, \phi; \mathbf{x}) = & \mathbb{E}_{Q_{\phi}(\mathbf{z}_N|\mathbf{x})}[\log P_{\theta}(\mathbf{x}|\mathbf{z}_N)] \\ & - \frac{1}{N} \sum_{t=0}^N \text{KL}(Q_{\phi_t}(\mathbf{z}_t|\mathbf{x})\|P(\mathbf{z}_t)). \end{aligned} \quad (3)$$

As can be seen in Eq. 3, our solution to the KL collapse issue does not require any engineering for balancing the weight between the reconstruction term and KL loss as commonly the case in existing works (Bowman et al., 2016; Sønderby et al., 2016). The weight between these two terms of our model is simply 1 : 1.

3 Experimental Setup

3.1 Datasets

We evaluate our model on two public datasets, namely, Penn Treebank (PTB) (Marcus and Marcinkiewicz, 1993) and the end-to-end (E2E) text generation corpus (Novikova et al., 2017), which have been used in a number of previous works for text generation (Bowman et al., 2016; Xu and Durrett, 2018; Wiseman et al., 2018; Su et al., 2018). PTB consists of more than 40,000 sentences from Wall Street Journal articles whereas the E2E dataset contains over 50,000 sen-

tences of restaurant reviews. The statistics of these two datasets are summarised in Table 1.

3.2 Implementation Details

For the PTB dataset, we used the train-test split following (Bowman et al., 2016; Xu and Durrett, 2018). For the E2E dataset, we used the train-test split from the original dataset (Novikova et al., 2017) and indexed the words with a frequency higher than 3. We represent input data with 512-dimensional word2vec embeddings (Mikolov et al., 2013). We set the dimension of the hidden layers of both encoder and decoder to 256. The Adam optimiser (Kingma and Ba, 2014) was used for training with an initial learning rate of 0.0001. Each utterance in a mini-batch was padded to the maximum length for that batch, and the maximum batch-size allowed is 128.

3.3 Baselines

We compare our HR-VAE model with three strong baselines using VAE for text modelling:

VAE-LSTM-base³: A variational autoencoder model which uses LSTM for both encoder and decoder. KL annealing is used to tackled the latent variable collapse issue (Bowman et al., 2016);

VAE-CNN⁴: A variational autoencoder model with a LSTM encoder and a dilated CNN decoder (Yang et al., 2017);

vMF-VAE⁵: A variational autoencoder model using LSTM for both encoder and decoder where the prior distribution is the von Mises-Fisher (vMF) distribution rather than a Gaussian distribution (Xu and Durrett, 2018).

4 Experimental Results

We evaluate our HR-VAE model in two experimental settings, following the setup of (Bowman et al., 2016; Xu and Durrett, 2018). In the *standard setting*, the input to the decoder at each time stamp is the concatenation of latent variable \mathbf{z} and the ground truth word of the previous time stamp. Under this setting, the decoder will be more powerful because it uses the ground truth word as input, resulting in little information of the training data captured by latent variable \mathbf{z} . The *inputless setting*, in contrast, does not use the previous ground truth word as input for the decoder. In other words,

³<https://github.com/timbmg/Sentence-VAE>

⁴<https://github.com/kefirski/contiguous-succotash>

⁵https://github.com/jiacheng-xu/vmf_vae_nlp

Dataset	Training	Development	Testing	Avg. sent. length	Vocab.
PTB	42,068	3,370	3,761	21.1	10K
E2E	42,061	4,672	4,693	22.67	2.8K

Table 1: The statistics of the PTB and E2E datasets.

Model	PTB				E2E			
	Standard		Inputless		Standard		Inputless	
	NLL	PPL	NLL	PPL	NLL	PPL	NLL	PPL
VAE-LSTM-base	101 [†] (2 [†])	119 [†]	125 [†] (15 [†])	380 [†]	50 (1.88)	5.77	101 (5.48)	34.70
VAE-CNN	99 (3.1)	113	121 (16.2)	323	41 (3.02)	4.23	82 (5.95)	17.81
vMF-VAE	96 [†] (5.7 [†])	98 [†]	117 [†] (18.6 [†])	262 [†]	34 (7.63)	3.29	61 (19.58)	8.52
HR-VAE (Ours)	79 (10.4)	43	85 (17.32)	54	20 (5.37)	2.02	38 (7.78)	3.74

Table 2: Language modelling results on the PTB and E2E datasets. [†] indicates the results which are reported from the prior publications. KL loss is shown in the parenthesis.

the decoder needs to predict the entire sequence with only the help of the given latent variable \mathbf{z} . In this way, a high-quality representation abstracting the information of the input sentence is much needed for the decoder, and hence enforcing \mathbf{z} to learn the required information.

Overall performance. Table 2 shows the language modelling results of our approach and the baselines. We report negative log likelihood (NLL), KL loss, and perplexity (PPL) on the test set. As expected, all the models have a higher KL loss in the inputless setting than the standard setting, as \mathbf{z} is required to encode more information about the input data for reconstruction. In terms of overall performance, our model outperforms all the baselines in both datasets (i.e., PTB and E2E). For instance, when comparing with the strongest baseline vMF-VAE in the standard setting, our model reduces NLL from 96 to 79 and PPL from 98 to 43 in PTB, respectively. In the inputless setting, our performance gain is even higher, i.e., NLL reduced from 117 to 85 and PPL from 262 to 54. A similar pattern can be observed for the E2E dataset. These observations suggest that our approach can learn a better generative model for data.

Loss analysis. To conduct a more thorough evaluation, we further investigate model behaviours in terms of both reconstruction loss and KL loss, as shown in Figure 2. These plots were obtained based on the E2E training set using the inputless setting.

We can see that the KL loss of VAE-LSTM-base, which uses Sigmoid annealing (Bowman et al., 2016), collapses to zero, leading to a poor generative performance as indicated by the high

reconstruction loss. The KL loss for both VAE-CNN and vMF-VAE are nonzero, where the former mitigates the KL collapse issue with a KL loss clipping strategy and the latter by replacing the standard normal distribution for the prior with the vMF distribution (i.e., with the vMF distribution, the KL loss only depends on a fixed concentration parameter, and hence results in a constant KL loss). Although both VAE-CNN and vMF-VAE outperform VAE-LSTM-base by a large margin in terms of reconstruction loss as shown in Figure 2, one should also notice that these two models actually overfit the training data, as their performance on the test set is much worse (cf. Table 2). In contrast to the baselines which mitigate the KL collapse issue by carefully engineering the weight between the reconstruction loss and KL loss or choosing a different choice of prior, we provide a simple and elegant solution through holistic KL regularisation, which can effectively mitigate the KL collapse issue and achieve a better reconstruction error in both training and testing.

Sentence reconstruction. Lastly, we show some sentence examples reconstructed by vMF-VAE (i.e., the best baseline) and our model in the inputless setting using sentences from the E2E test set as input. As shown in Table 3, the sentences generated by vMF-VAE contain repeated words in quite a few cases, such as ‘city city area’ and ‘blue spice spice’. In addition, vMF-VAE also tends to generate unnecessary or unrelated words at the end of sentences, making the generated sentences ungrammatical. The sentences reconstructed by our model, in contrast, are more grammatical and more similar to the corresponding ground truth sentences than vMF-VAE.

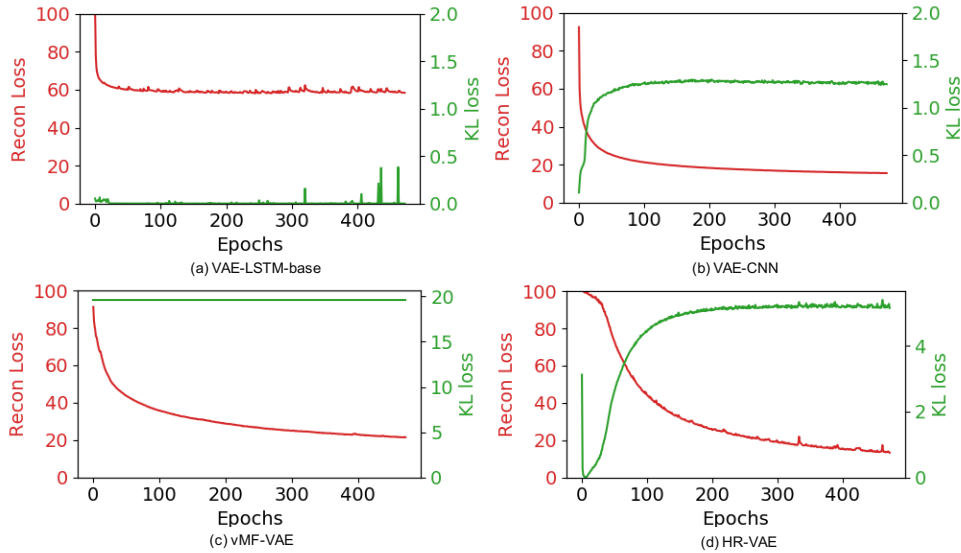


Figure 2: Training curves of reconstruction loss and KL loss of (a) VAE-LSTM-base, (b) VAE-CNN, (c) vMF-VAE, and (d) our model, based on the E2E training set using the inputless setting.

Input	<ol style="list-style-type: none"> 1. blue spice is a coffee shop in city centre . 2. giraffe is a coffee shop found near the bakers . 3. a pub in the city centre area called blue spice 4. pub located near café sicilia called cocum with a high customer rating 5. the cricketers is a one star coffee shop near the ranch that is not family friendly .
vMF-VAE	<ol style="list-style-type: none"> 1. blue spice is a coffee in city centre . it is not , and 2. cotto is a coffee shop located near the bakers . . is 5 out of 3. a coffee in the city city area is blue spice spice . the is is 4. located located near café rouge , cotto has a high customer rating and a customer 5. the cricketers is a low rated coffee shop near the bakers that is a star , is is
Ours	<ol style="list-style-type: none"> 1. blue spice is a coffee shop in city centre . 2. giraffe is a coffee shop located near the bakers . 3. a restaurant in the city centre called blue spice italian 4. located place near café sicilia called punter has a high customer rating 5. the cricketers is a one star coffee shop near ranch ranch that is not family friendly .

Table 3: Example input sentences from the E2E test dataset (top); sentences reconstructed by vMF-VAE (middle); sentences reconstructed by our model (bottom).

5 Conclusion

In this paper, we present a simple and generic architecture called holistic regularisation VAE (HR-VAE), which can effectively avoid latent variable collapse. In contrast to existing VAE-RNN models which merely impose a standard normal distribution prior on the last hidden state of the RNN encoder, our HR-VAE model imposes regularisation on all the hidden states, allowing a better regularisation of the model learning process. Empirical results show that our model can effectively mitigate the latent variable collapse issue while giving a better predictive performance than the baselines.

Acknowledgment

This work is supported by the award made by the UK Engineering and Physical Sciences Research

Council (Grant number: EP/P011829/1).

References

- Philip Bachman. 2016. An architecture for deep, hierarchical generative models. In *Advances in Neural Information Processing Systems*, pages 4826–4834.
- Samuel Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. 2016. Generating sentences from a continuous space. In *Proceedings of the Twentieth Conference on Computational Natural Language Learning (CoNLL)*.
- Adji B Dieng, Yoon Kim, Alexander M Rush, and David M Blei. 2019. Avoiding latent variable collapse with generative skip models. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2397–2405.

- Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. 2016. Sequential neural models with stochastic layers. In *Advances in neural information processing systems*, pages 2199–2207.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Mitchell P Marcus and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of english: The penn treebank. *Computational Linguistics*, 19(2).
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. 2017. The e2e dataset: New challenges for end-to-end generation. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 201–206.
- Stanislau Semeniuta, Aliaksei Severyn, and Erhardt Barth. 2017. A hybrid convolutional variational autoencoder for text generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 627–637.
- Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. 2016. How to train deep variational autoencoders and probabilistic ladder networks. In *33rd International Conference on Machine Learning (ICML 2016)*.
- Shang-Yu Su, Kai-Ling Lo, Yi Ting Yeh, and Yun-Nung Chen. 2018. Natural language generation by hierarchical decoding with linguistic patterns. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 61–66.
- Sam Wiseman, Stuart Shieber, and Alexander Rush. 2018. Learning neural templates for text generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3174–3187.
- Jiacheng Xu and Greg Durrett. 2018. Spherical latent spaces for stable variational autoencoders. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4503–4513.
- Zichao Yang, Zhiting Hu, Ruslan Salakhutdinov, and Taylor Berg-Kirkpatrick. 2017. Improved variational autoencoders for text modeling using dilated convolutions. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3881–3890. JMLR. org.