

A Surprisingly Effective Fix for Deep Latent Variable Modeling of Text

Bohan Li*¹, Junxian He*¹, Graham Neubig¹, Taylor Berg-Kirkpatrick², Yiming Yang¹

Core idea:

- Pretraining encoder using AE objective
- Initializing VAE with pertained AE and training with FB objective

Results:

	Yahoo				
LSTM-LM	60.75	-	-	-	-
VAE	61.52	329.10	0	0.00	329.10
+ anneal	61.21	328.80	0	0.00	328.80
+ cyclic	66.93	333.80	4	2.83	336.63
+ aggressive	59.77	322.70	15	5.70	328.40
+ FBP ($\lambda = 9$)	62.59	322.91	6	9.08	331.99
+ FBP ($\lambda = 7$)	62.76	324.66	5	7.03	331.69
+ FBP ($\lambda = 5$)	62.78	326.26	3	5.07	331.32
+ FBP ($\lambda = 3$)	62.88	328.13	2	3.06	331.19
Ours ($\lambda = 6$)	59.23	317.39	32	12.09	329.48
Ours ($\lambda = 8$)	59.51	315.31	32	15.02	330.33
Ours ($\lambda = 9$)	59.60	315.09	32	15.49	330.58

Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang[†], Chunyuan Li[§], Jianfeng Gao[§], Wen Dong[†], Changyou Chen[†]

Core ideas:

- Attributing posterior collapse to the restrictive Gaussian assumption, and advocate more flexible sample-based posterior representation.
- Proposed iVAE:
 - Instead of assuming posterior as Gaussian, use sampling mechanism for $q_\phi(z|x)$. $z_{x,i} = G_\phi(x, \epsilon_i)$, $\epsilon_i \sim q(\epsilon)$
 - How to evaluate KL on this? Use dual form of $KL(q_\phi(z|x) \parallel p(z))$:

$$\begin{aligned} & KL(q_\phi(z|x) \parallel p(z)) \\ &= \max_{\nu} \mathbb{E}_{z \sim q_\phi(z|x)} \nu_\psi(x, z) - \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(x, z)), \end{aligned} \quad (7)$$

- New iVAE objective:

$$\begin{aligned} \mathcal{L}_{iVAE} = & \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x|z) \\ & - \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{z \sim q_\phi(z|x)} \nu_\psi(x, z) \\ & + \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(x, z)), \end{aligned} \quad (8)$$

- Here, ν_ψ is a MLP taking in (x, z) .

Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang[†], Chunyuan Li[§], Jianfeng Gao[§], Wen Dong[†], Changyou Chen[†]

Core ideas:

- Training Scheme:

- Sample a mini-batch of $\mathbf{x}_i \sim \mathcal{D}$, $\epsilon_i \sim q(\epsilon)$, and generate $\mathbf{z}_{\mathbf{x}_i, \epsilon_i} = G(\mathbf{x}_i, \epsilon_i; \phi)$; Sample a mini-batch of $\mathbf{z}_i \sim p(\mathbf{z})$.

- Update ψ in $\nu_\psi(\mathbf{x}, \mathbf{z})$ to maximize

$$\sum_i \nu_\psi(\mathbf{x}_i, \mathbf{z}_{\mathbf{x}_i, \epsilon_i}) - \sum_i \exp(\nu_\psi(\mathbf{x}_i, \mathbf{z}_i)) \quad (9)$$

- Update parameters $\{\phi, \theta\}$ to maximize

$$\sum_i \log p_\theta(\mathbf{x}_i | \mathbf{z}_{\mathbf{x}_i, \epsilon_i}) - \sum_i \nu_\psi(\mathbf{x}_i, \mathbf{z}_{\mathbf{x}_i, \epsilon_i}) \quad (10)$$

Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang[†], Chunyuan Li[§], Jianfeng Gao[§], Wen Dong[†], Changyou Chen[†]

Core ideas:

- Proposed Mutual Information Regularized iVAE
 - Replace $-\text{KL}(q_\phi(z|x) \parallel p(z))$ with $-\text{KL}(q_\phi(z) \parallel p(z))$, where $q_\phi(z) = \int q(x)q_\phi(z|x)dx$, estimated by ancestral sampling in practice.
 - This objective also maximize the mutual information $I(x,z)$, as they claimed.
 - The new objective:

$$\begin{aligned} & \text{KL}(q_\phi(z) \parallel p(z)) \\ &= \max_{\nu} \mathbb{E}_{z \sim q_\phi(z)} \nu_\psi(z) - \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(z)). \end{aligned} \quad (13)$$

$$\begin{aligned} \mathcal{L}_{\text{iVAE}_{\text{MI}}} &= \mathbb{E}_{x \sim D} \mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x|z) \\ &\quad - \mathbb{E}_{z \sim q_\phi(z)} \nu_\psi(z) + \mathbb{E}_{z \sim p(z)} \exp(\nu_\psi(z)), \end{aligned} \quad (14)$$

- The only difference with iVAE is that ν_ψ is a MLP taking in only z .

Implicit Deep Latent Variable Models for Text Generation (EMNLP 2019)

Le Fang[†], Chunyuan Li[§], Jianfeng Gao[§], Wen Dong[†], Changyou Chen[†]

Results:

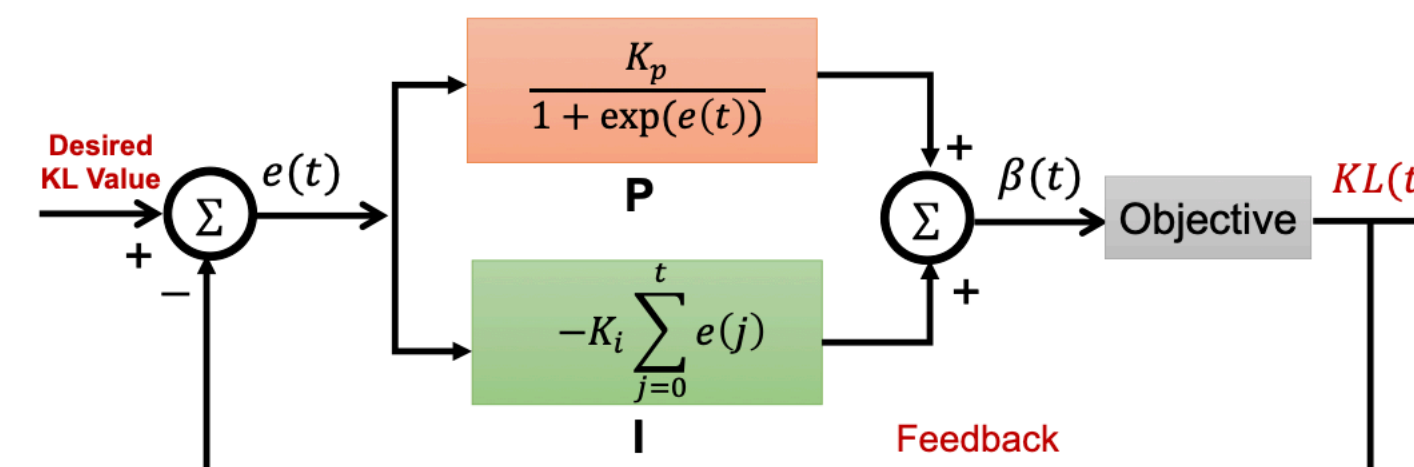
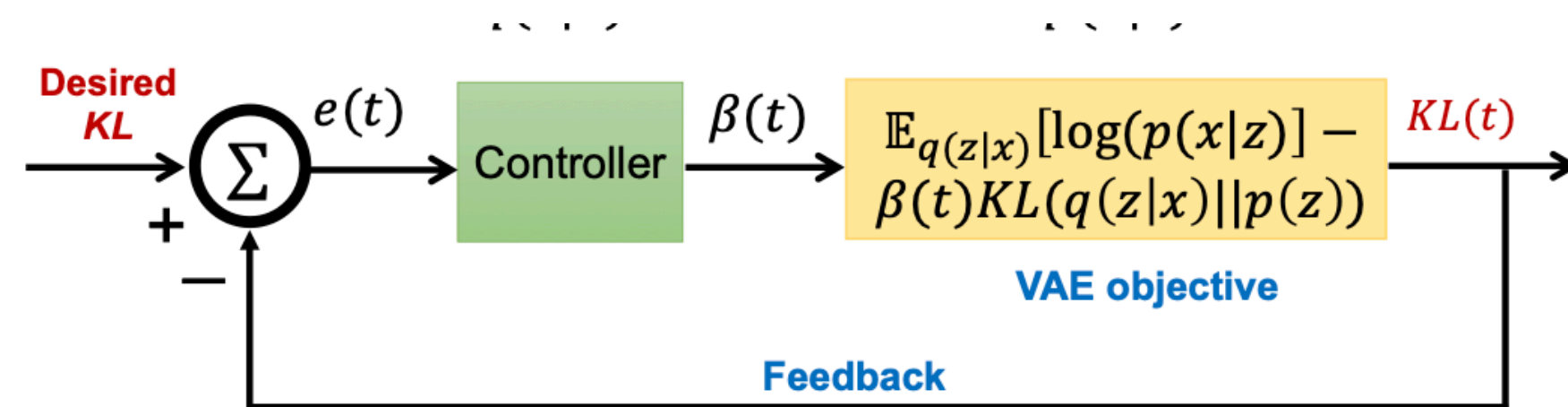
Methods	-ELBO↓	PPL↓	KL↑	MI↑	AU↑
Dataset: PTB					
VAE	102.6	108.26	1.08	0.8	2
$\beta(0.5)$ -VAE	104.5	117.92	7.50	3.1	5
SA-VAE	102.6	107.71	1.23	0.7	2
Cyc-VAE	103.1	110.50	3.48	1.8	5
iVAE	87.6	54.46	6.32	3.5	32
iVAE _{MI}	87.2	53.44	12.51	12.2	32
Dataset: Yahoo					
VAE	328.6	61.21	0.0	0.0	0
$\beta(0.4)$ -VAE	328.7	61.29	6.3	2.8	8
SA-VAE	327.2	60.15	5.2	2.7	10
Lag-VAE	326.7	59.77	5.7	2.9	15
iVAE	309.5	48.22	8.0	4.4	32
iVAE _{MI}	309.1	47.93	11.4	10.7	32
Dataset: Yelp					
VAE	357.9	40.56	0.0	0.0	0
$\beta(0.4)$ -VAE	358.2	40.69	4.2	2.0	4
SA-VAE	355.9	39.73	2.8	1.7	8
Lag-VAE	355.9	39.73	3.8	2.4	11
iVAE	348.2	36.70	7.6	4.6	32
iVAE _{MI}	348.7	36.88	11.6	11.0	32

ControlVAE: Controllable Variational Autoencoder (ICML 2020)

Huajie Shao et al

Core ideas:

- Enhance β -VAE by controlling β with a new non-linear PI controller, $\beta(t)$. The controller is not learnable, with K_p , K_i , Desired KL being hyperparameters.



- Intuition: Comparing error between Desired KL and Sampled KL (from inference output); if Sampled KL is too small, $\beta(t)$ will be small so KL-divergence is encouraged to grow, and vice versa.

ControlVAE: Controllable Variational Autoencoder (ICML 2020)

Huajie Shao et al

Core ideas:

- PI Algorithm

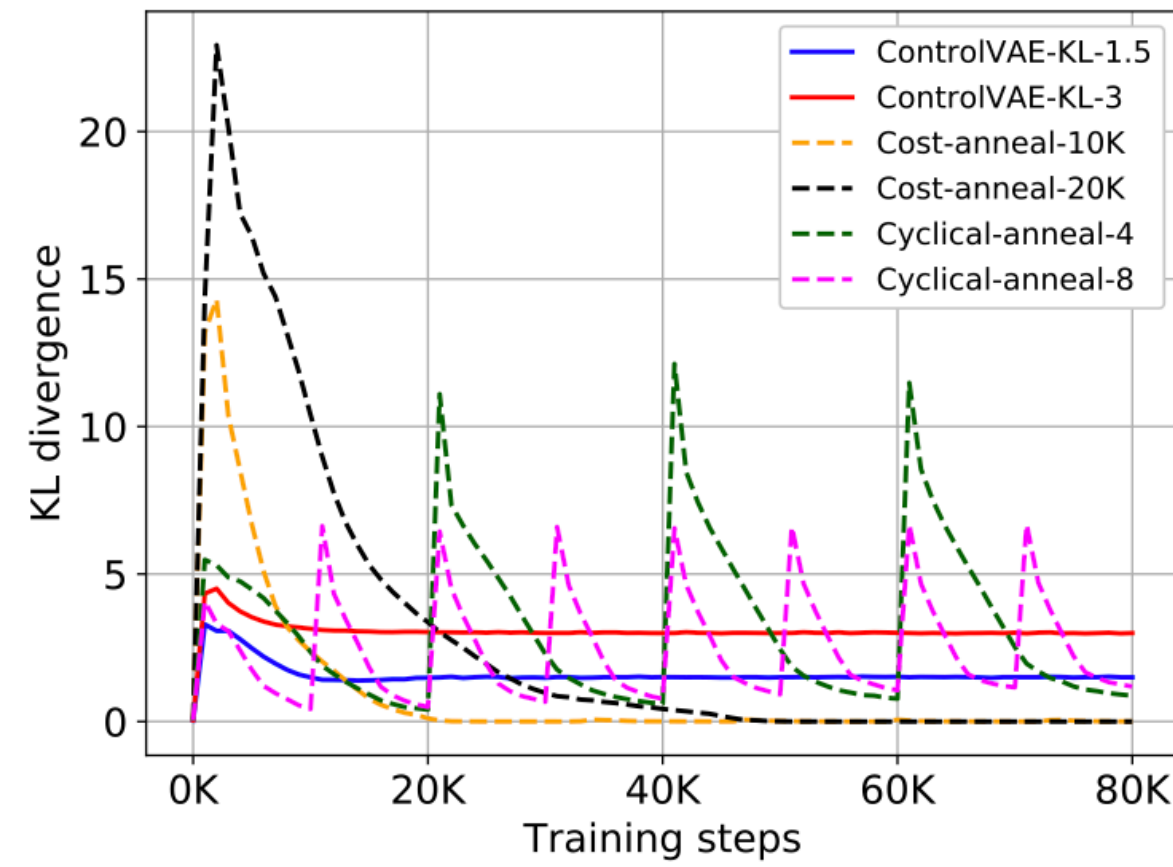
Algorithm 1 PI algorithm.

```
1: Input: desired KL  $v_{kl}$ , coefficients  $K_p$ ,  $K_i$ , max/min value  
    $\beta_{max}$ ,  $\beta_{min}$ , iterations  $N$   
2: Output: hyperparameter  $\beta(t)$  at training step  $t$   
3: Initialization:  $I(0) = 0$ ,  $\beta(0) = 0$   
4: for  $t = 1$  to  $N$  do  
5:   Sample KL-divergence,  $\hat{v}_{kl}(t)$   
6:    $e(t) \leftarrow v_{kl} - \hat{v}_{kl}(t)$   
7:    $P(t) \leftarrow \frac{K_p}{1 + \exp(e(t))}$   
8:   if  $\beta_{min} \leq \beta(t-1) \leq \beta_{max}$  then  
9:      $I(t) \leftarrow I(t-1) - K_i e(t)$   
10:  else  
11:     $I(t) = I(t-1)$  // Anti-windup  
12:  end if  
13:   $\beta(t) = P(t) + I(t) + \beta_{min}$   
14:  if  $\beta(t) > \beta_{max}$  then  
15:     $\beta(t) = \beta_{max}$   
16:  end if  
17:  if  $\beta(t) < \beta_{min}$  then  
18:     $\beta(t) = \beta_{min}$   
19:  end if  
20:  Return  $\beta(t)$   
21: end for
```

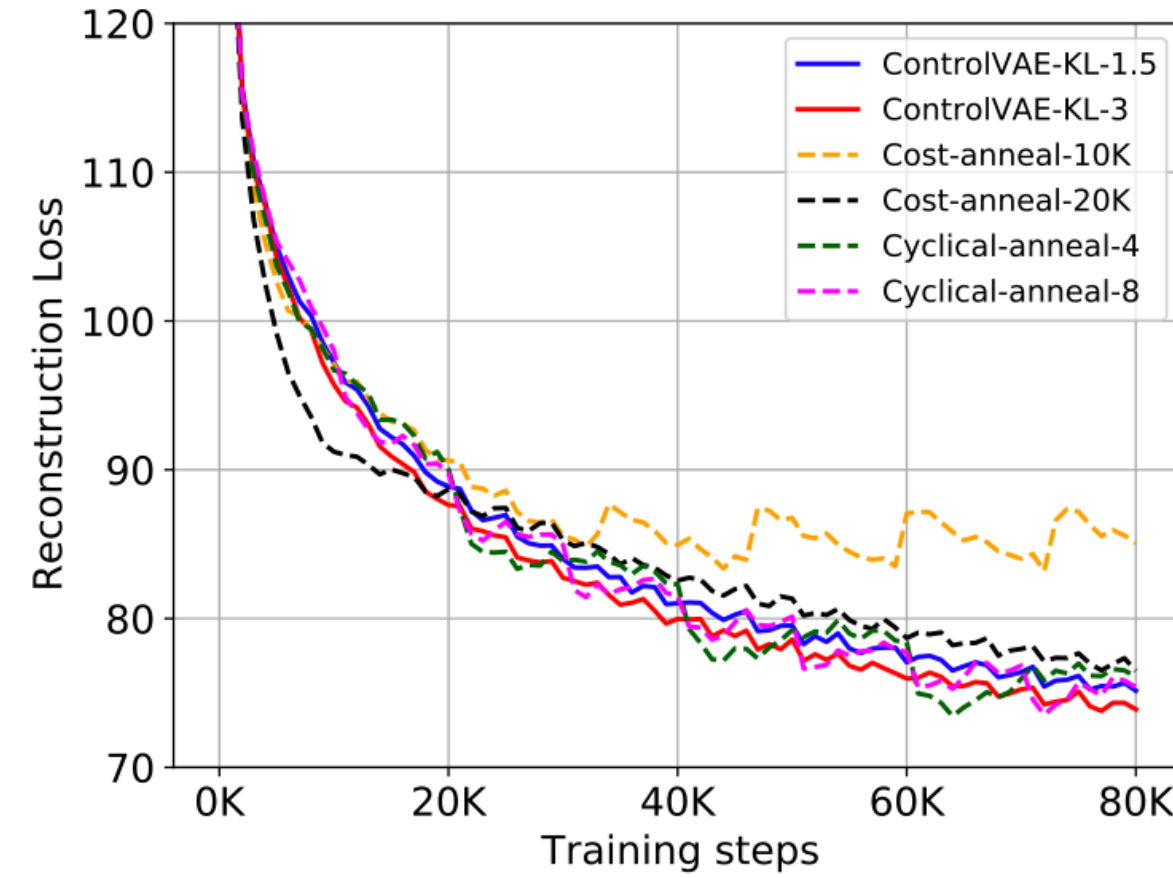
ControlVAE: Controllable Variational Autoencoder (ICML 2020)

Huajie Shao et al

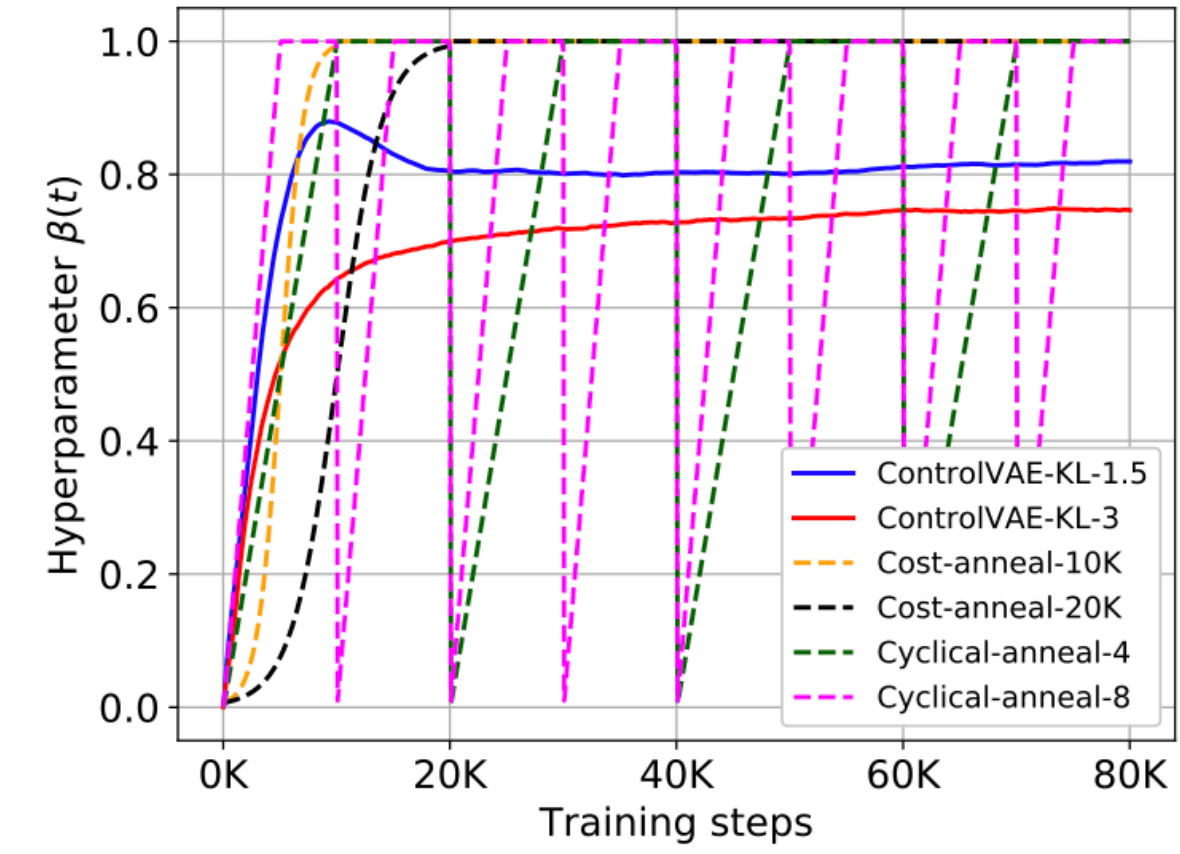
Results:



(a) KL divergence



(b) Reconstruction loss



(c) $\beta(t)$

Methods/metric	Dis-1	Dis-2	self-BLEU-2	self-BLEU-3	PPL
ControlVAE-KL-35	6.27K \pm 41	95.86K \pm 1.02K	0.663 \pm 0.012	0.447 \pm 0.013	8.81 \pm 0.05
ControlVAE-KL-25	6.10K \pm 60	83.15K \pm 4.00K	0.698 \pm 0.006	0.495 \pm 0.014	12.47 \pm 0.07
Cost anneal-KL-17	5.71K \pm 87	69.60K \pm 1.53K	0.721 \pm 0.010	0.536 \pm 0.008	16.82 \pm 0.11
Cyclical (KL = 21.5)	5.79K \pm 81	71.63K \pm 2.04K	0.710 \pm 0.007	0.524 \pm 0.008	17.81 \pm 0.33

FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

Core ideas:

- Using normalizing flows (NF) for the prior distribution. In this paper, using *real-valued non-volume preserving* (real NVP) transformations (Dinh et al., 2016)
 - $z_L = f_L \circ f_{L-1} \circ \dots \circ f_1(z_0)$, $z_0 \sim N(0, I)$, z_L is the sentence latent variable.

$$z_0 = f_1^{-1} \circ \dots \circ f_{L-1}^{-1} \circ f_L^{-1}(z_L) \quad (5)$$

$$\log p_\psi(z_L) = \log p_0(z_0) - \sum_{l=1}^L \log \left| \det \left(\frac{\partial f_l(z_{l-1})}{\partial z_{l-1}} \right) \right| \quad (6)$$

- Inspired by Real-NVP, the choice of f is an affine transformation, which is $f: z_{i-1} \rightarrow z_i$

$$\begin{aligned} z_i^{(1:d)} &= z_{i-1}^{(1:d)} \\ z_i^{(d+1:D)} &= z_{i-1}^{(d+1:D)} \odot \exp(s(z_{i-1}^{(1:d)})) + t(z_{i-1}^{(1:d)}) \end{aligned}$$

- 1) easily invertible
- 2) its Jacobian determinant is easy to compute. (Do not require inverse/Jacobian of s and t)
- So we can model s and t using NNs, denoted by $p_\psi(z)$

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix}$$

FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

Core ideas:

- New Objective: Using importance weighting + Monte Carlo for KL

$$\begin{aligned}\mathcal{L}(\theta, \phi, \psi; x) = & \log \frac{1}{N} \sum_{i=1}^N \frac{p_{\theta}(x|z^{(i)})p_{\psi}(z^{(i)})}{q_{\phi}(z^{(i)}|x)} \\ & + \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x|z^{(i)}) - \text{KL}_{\phi, \psi}(x, \{z^{(i)}\}_{i=1}^N) \\ & \text{s.t. } z^{(i)} \sim q_{\phi}(z|x) \quad (7)\end{aligned}$$

FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

Core ideas:

- Training Scheme

1. Draw N samples $z_L^{(1)}, z_L^{(2)}, \dots, z_L^{(N)}$ from the inference network using the reparameterization trick.
2. Perform the inverse transformation to get the image of each point under the base distribution: $z_0^{(1)}, z_0^{(2)}, \dots, z_0^{(N)}$.
3. Compute the exact log likelihood of the sample prior with change of variable theorem (Eq. 6).
4. Compute and backpropagate the loss (Eq. 7).

FlowPrior: Learning Expressive Priors for Latent Variable Sentence Models (NAACL 2021)

Xiaoan Ding, Kevin Gimpel

Results:

Model	PPL(↓)	Recon(↓)	KL	AU(↑)	MI(↑)
VAE	101.40	101.28	0.00	0	0.00
Cyc-VAE	107.73	101.17	2.01	5	1.24
Lag-VAE	100.25	100.41	1.04	3	0.79
VAE + FB	101.56	99.84	4.46	32	0.90
Pre-VAE + FB	96.35	94.52	8.15	32	6.30
MoG-VAE	98.22	100.54	0.00	0	0.00
MoG-VAE + FB	97.50	99.44	2.35	32	0.68
Vamp-VAE	98.27	100.56	0.00	0	0.00
Vamp-VAE + FB	97.83	99.53	2.31	32	0.72
FlowPrior	94.72	98.46	3.28	2	2.25
FlowPrior + FB	93.58	99.20	7.21	31	2.83

Table 1: Language modeling results on PTB dataset.

Adversarial Poincaré Variational Autoencoder (APo- VAE) (NAACL 2021)

Shuyang Dai·Zhe Gan·Yu Cheng·Chenyang Tao·Lawrence Carin·Jingjing Liu

Core ideas:

- Natural languages have latent hierarchy, but prior distribution from Euclidean space couldn't capture it, so resorting to Riemannian Geometry.
- Proposing a prior based on Poincaré Ball Model in Hyperbolic space, with the following advantages:

$$\mathbb{B}_c^n := \{z \in \mathbb{R}^n \mid c\|z\|^2 < 1\}$$

- It implies significant representation ability.
- It's a well defined compact metric space, with well defined algebraic operations that allows back-propagation.

$$z \oplus_c z' := \frac{(1 + 2c\langle z, z' \rangle + c\|z'\|^2)z + (1 - c\|z\|^2)z'}{1 + 2c\langle z, z' \rangle + c^2\|z\|^2\|z'\|^2}. \quad (4)$$

$$\begin{aligned} \exp_\mu^c(u) &:= \mu \oplus_c \left(\tanh\left(\sqrt{c}\frac{\lambda_\mu^c\|u\|}{2}\right) \frac{u}{\sqrt{c}\|u\|} \right), \\ \log_\mu^c(y) &:= \frac{2}{\sqrt{c}\lambda_\mu^c} \tanh^{-1}\left(\sqrt{c}\|\kappa_{\mu,y}\|\right) \frac{\kappa_{\mu,y}}{\|\kappa_{\mu,y}\|}, \end{aligned} \quad (5)$$

$$P_{0 \rightarrow \mu}^c(v) = \log_\mu^c(\mu \oplus_c \exp_0^c(v)) = \frac{\lambda_0^c}{\lambda_\mu^c} v. \quad (6)$$

where $\kappa_{\mu,y} := (-\mu) \oplus_c y$.

Adversarial Poincaré Variational Autoencoder (APo- VAE) (NAACL 2021)

Shuyang Dai·Zhe Gan·Yu Cheng·Chenyang Tao·Lawrence Carin·Jingjing Liu

Core ideas:

- How does the defined compact metric space contribute to the formulation of a prior?
 - Poincaré ball based normal prior $z \sim N_{B^n_c}(\mu, \Sigma)$:

$$z = \exp_{\mu}^c \left(\frac{\lambda_0^c}{\lambda_{\mu}^c} v \right), v \sim \mathcal{N}(\mathbf{0}, \Sigma). \quad (7)$$

- Inspired by iVAE, we can enhance $v := G(\mathbf{x}, \xi; \phi_1)$ and $\mu := F(\mathbf{x}; \phi_2)$, as outputs of encoder (ϕ). Then we get z from the above formula.
- Same as iVAE, they use dual form to evaluate $KL(q_{\phi}(z|\mathbf{x}) \parallel p(z))$:

$$\mathbb{D}_{KL}(q_{\phi}(z|\mathbf{x}) \parallel p(z)) = \max_{\psi} \quad (9)$$
$$\left\{ \mathbb{E}_{z \sim q_{\phi}(z|\mathbf{x})} \nu_{\psi}(\mathbf{x}, z) - \mathbb{E}_{z \sim p(z)} \exp \nu_{\psi}(\mathbf{x}, z) \right\},$$

- Same as iVAE, ν_{ψ} is optimized by the following:

$$\mathcal{L}_1 = \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \left[\mathbb{E}_{z \sim q_{\phi}(z|\mathbf{x})} \nu_{\psi}(\mathbf{x}, z) - \mathbb{E}_{z \sim p(z)} \exp \nu_{\psi}(\mathbf{x}, z) \right], \quad (10)$$

Adversarial Poincaré Variational Autoencoder (APo- VAE) (NAACL 2021)

Shuyang Dai · Zhe Gan · Yu Cheng · Chenyang Tao · Lawrence Carin · Jingjing Liu

Core ideas:

- Different from iVAE, here the real prior is not assumed Gaussian. It's estimated by sampling scheme used in VampPrior (Tomczak and Welling, 2018)

$$p_{\delta}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^K q_{\phi}(\mathbf{z} | \mathbf{s}_k), \quad (12)$$

- $\delta := \{\mathbf{s}_k\}_{k=1}^K$ is now learnable pseudo inputs. Replacing $p(\mathbf{z})$ with $p_{\delta}(\mathbf{z})$ seeks to match the aggregated posterior $q(\mathbf{z}) = \sum q_{\phi}(\mathbf{z} | \mathbf{x}_i) / N$.

Adversarial Poincaré Variational Autoencoder (APo- VAE) (NAACL 2021)

Shuyang Dai·Zhe Gan·Yu Cheng·Chenyang Tao·Lawrence Carin·Jingjing Liu

Core ideas:

- For the geometry-aware decoder, they use a deterministic (and learnable) hyperbolic linear function f to extract feature, then pass to LSTM decoder. (a, b , together with the LSTM, are trainable parameters of decoder θ)

$$f_{a,b}^c(z) = \text{sign}(\langle \mathbf{a}, \log_b^c(z) \rangle_b) \|\mathbf{a}\|_b d_c^{\mathbb{B}}(z, H_{a,b}^c), \quad (8)$$

where $H_{a,b}^c = \{z \in \mathbb{B}_c^n | \langle \mathbf{a}, \log_b^c(z) \rangle_b = 0\}$,

$$d_c^{\mathbb{B}}(z, H_{a,b}^c) = \frac{1}{\sqrt{c}} \sinh^{-1} \left(\frac{2\sqrt{c} |\langle \kappa_{b,z}, \mathbf{a} \rangle|}{(1 - c \|\kappa_{b,z}\|^2) \|\mathbf{a}\|} \right)$$

- θ and ϕ are optimized by the following objective

$$\mathcal{L}_2 = \mathbb{E}_{\mathbf{x} \sim \mathbf{X}} \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}) - \nu_\psi(\mathbf{x}, \mathbf{z})]. \quad (11)$$

Adversarial Poincaré Variational Autoencoder (APo- VAE) (NAACL 2021)

Shuyang Dai·Zhe Gan·Yu Cheng·Chenyang Tao·Lawrence Carin·Jingjing Liu

Core ideas:

- Training Procedure

Algorithm 1 Training procedure of APo-VAE.

- 1: **Input:** Data samples $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$, Poincaré curvature c , and number of pseudo-input K .
- 2: Initialize θ , ϕ , ψ , and δ .
- 3: **for** $iter$ from 1 to max_iter **do**
- 4: Sample a mini-batch $\{\mathbf{x}_m\}_{m=1}^M$ from \mathbf{X} of size M .
- 5: *# Sampling in the Hyperbolic Space.*
- 6: Obtain μ_m and \mathbf{v}_m from $\text{EncNet}_\phi(\mathbf{x}_m)$.
- 7: Move \mathbf{v}_m to $\mathbf{u}_m = P_{\mathbf{0} \rightarrow \mu_m}^c(\mathbf{v}_m)$ by (6).
- 8: Map \mathbf{u}_m to $\mathbf{z}_m = \exp_{\mu_m}^c(\mathbf{u}_m)$ by (5).
- 9: *# Update the dual function and the pseudo-input.*
- 10: Sample $\tilde{\mathbf{z}}_m$ by (12).
- 11: Update ψ and δ by gradient ascent on (10).
- 12: *# Update the encoder and decoder networks.*
- 13: Update θ and ϕ by gradient ascent on (11).
- 14: **end for**

Adversarial Poincaré Variational Autoencoder (APo- VAE) (NAACL 2021)

Shuyang Dai·Zhe Gan·Yu Cheng·Chenyang Tao·Lawrence Carin·Jingjing Liu

Results:

Model	-ELBO	PPL	KL	MI	AU
	PTB				
VAE	102.6	108.26	1.1	0.8	2
β -VAE	104.5	117.92	7.5	3.1	5
SA-VAE	102.6	107.71	1.2	0.7	2
vMF-VAE	95.8	93.70	2.9	3.2	21
\mathcal{P} -VAE	91.4	76.13	4.5	2.9	23
iVAE	87.2	53.44	12.5	12.2	32
APo-VAE	87.2	53.32	8.4	4.8	32
APo-VAE+VP	87.0	53.02	8.9	4.5	32
Yahoo					
VAE	328.6	61.21	0.0	0.0	0
β -VAE	328.7	61.29	6.3	2.8	8
SA-VAE	327.2	60.15	5.2	2.9	10
LAG-VAE	326.7	59.77	5.7	2.9	15
vMF-VAE	318.5	53.92	6.3	3.7	23
\mathcal{P} -VAE	313.4	50.57	7.2	3.3	27
iVAE	309.1	47.93	11.4	10.7	32
APo-VAE	286.2	47.00	6.9	4.1	32
APo-VAE+VP	285.6	46.61	8.1	4.9	32
Yelp					
VAE	357.9	40.56	0.0	0.0	0
β -VAE	358.2	40.69	4.2	2.0	4
SA-VAE	357.8	40.51	2.8	1.7	8
LAG-VAE	355.9	39.73	3.8	2.4	11
vMF-VAE	356.2	51.03	4.1	3.9	13
\mathcal{P} -VAE	355.4	50.64	4.3	4.8	19
iVAE	348.7	36.88	11.6	11.0	32
APo-VAE	319.7	34.10	12.1	7.5	32
APo-VAE+VP	316.4	32.91	12.7	6.2	32