

Supplementary material for the paper Dynamic Neural Language Models

Edouard Delasalles, Sylvain Lamprier, and Ludovic Denoyer

Sorbonne Université, LIP6, F-75005, Paris, France
`firstname.lastname@lip6.fr`

A Full ELBO derivation

The generative model of DRLM writes as follow:

$$p_{\theta, \psi}(\mathcal{D}, \mathbf{Z}) = \prod_{i=1}^N p_{\theta}(\mathbf{d}^{(i)} | \mathbf{z}_{t(i)}) \prod_{t=0}^{T-1} p_{\psi}(\mathbf{z}_{t+1} | \mathbf{z}_t).$$

We use Variational Inference (VI) to learn this model, and consider a variational distribution $q_{\phi}(\mathbf{Z})$ that factorizes across all timesteps:

$$q_{\phi}(\mathbf{Z}) = \prod_{t=1}^T q_{\phi}^t(\mathbf{z}_t),$$

where q_{ϕ}^t are independent Gaussian distributions $\mathcal{N}(\mu_t, \sigma_t^2)$ with diagonal covariance matrices σ_t^2 , and ϕ is the total set of variational parameters.

A particularity of our approach is that we have several documents published at the same timestep. So, to obtain an Evidence Lower Bound (ELBO) we adapt

the derivation in [3] as follows:

$$\begin{aligned}
\log p_{\theta, \phi}(\mathcal{D}) &= \log \int_{\mathbf{Z}} p_{\psi}(\mathbf{Z}) \prod_{t=1}^T p_{\theta}(\mathcal{D}_t | \mathbf{z}_t) d\mathbf{Z} \\
&= \log \int_{\mathbf{Z}} q_{\phi}(\mathbf{Z}) p_{\psi}(\mathbf{Z}) \frac{\prod_{t=1}^T p_{\theta}(\mathcal{D}_t | \mathbf{z}_t)}{q_{\phi}(\mathbf{Z})} d\mathbf{Z} \\
&\geq \int_{\mathbf{Z}} q_{\phi}(\mathbf{Z}) \log \left(p_{\psi}(\mathbf{Z}) \frac{\prod_{t=1}^T p_{\theta}(\mathcal{D}_t | \mathbf{z}_t)}{q_{\phi}(\mathbf{Z})} \right) d\mathbf{Z} \\
&= \sum_{t=1}^T \int_{\mathbf{z}_t} q_{\phi}^t(\mathbf{z}_t) \log p_{\theta}(\mathcal{D}_t | \mathbf{z}_t) d\mathbf{z}_t \\
&\quad + \sum_{t=1}^T \int_{\mathbf{z}_{t-1}} q_{\phi}^{t-1}(\mathbf{z}_{t-1}) \int_{\mathbf{z}_t} q_{\phi}^t(\mathbf{z}_t) \log \frac{p_{\psi}(\mathbf{z}_t | \mathbf{z}_{t-1})}{q_{\phi}^t(\mathbf{z}_t)} d\mathbf{z}_{t-1} d\mathbf{z}_t \\
&= \sum_{t=1}^T \mathbb{E}_{q_{\phi}^t(\mathbf{z}_t)} [\log p_{\theta}(\mathcal{D}_t | \mathbf{z}_t)] - \sum_{t=1}^T \mathbb{E}_{q_{\phi}^{t-1}(\mathbf{z}_{t-1})} [D_{\text{KL}}(q_{\phi}^t(\mathbf{z}_t) \| p_{\psi}(\mathbf{z}_t | \mathbf{z}_{t-1}))] \\
&= \mathcal{L}(\theta, \psi, \phi), \tag{1}
\end{aligned}$$

where \mathcal{D}_t is the set of all documents published at timestep t , and the inequality is obtained thanks to the Jensen theorem on concave functions.

The KL between two Gaussians owns an analytically closed form. This allows us to rewrite our log-likelihood lower-bound, noted $\mathcal{L}(\theta, \psi, \phi)$, as follows:

$$\begin{aligned}
\mathcal{L}(\theta, \psi, \phi) &= \sum_{t=1}^T \mathbb{E}_{q_{\phi}^t(\mathbf{z}_t)} [p_{\theta}(\mathcal{D}_t | \mathbf{z}_t)] + \frac{Td}{2} \\
&\quad - \frac{1}{2} \left(T \sum_{i=0}^{d_z-1} \log \sigma_i^2 - \sum_{t=1}^T \sum_{i=0}^{d_z-1} \log \eta_{t,i}^2 + \sum_{t=1}^T \sum_{i=0}^{d_z-1} \frac{\eta_{t,i}^2}{\sigma_i^2} \right. \\
&\quad \left. + \sum_{t=1}^T \mathbb{E}_{q_{\phi}^{t-1}(\mathbf{z}_{t-1})} [(g(\mathbf{z}_{t-1}; \mathbf{w}) - \mu_t)' (\sigma^2)^{-1} (g(\mathbf{z}_{t-1}; \mathbf{w}) - \mu_t)] \right)
\end{aligned}$$

where we note A' the matrix transpose of a matrix A and where σ_i^2 and $\eta_{t,i}^2$ stand for the i -th component of diagonals σ^2 and η_t^2 respectively. This re-writing improve learning stability w.r.t. a version in which sampling would be done on the KL components too. Since the observation model p_{θ} is an RNN, the model is non-conjugate, and the ELBO in Equation 1 cannot be computed in closed form. We thus use the re-parametrization trick [2,4] to learn the model. It allows us to learn all parameters of the model jointly.

B Quantitative results for prediction

We report here the quantitative results for the prediction configuration on the language modeling task on 1.

Table 1: Prediction perplexity

Corpus	S2		NYT		Reddit	
Perplexity	<i>micro</i>	<i>macro</i>	<i>micro</i>	<i>macro</i>	<i>micro</i>	<i>macro</i>
LSTM	84.7	82.7	128.5	128.4	125.8	126.1
DT	92.0	89.6	137.1	137.0	151.1	151.6
DWE	87.0	84.8	140.1	140.0	136.5	139.9
DRLM-Id	81.2	79.2	123.7	123.6	124.7	125.0
DRLM	79.7	77.8	123.3	123.1	123.9	124.3

C Deriving Temporal Word Embedding Methods for Recurrent Language Modeling

We detail here how we adapt temporal word embeddings baselines to recurrent language modeling. The baselines are Dynamic Word Embeddings (DWE) [1], and DiffTime [5]. For both methods, we get rid of the context embeddings and only keep word embeddings \mathbf{U} .

C.1 Dynamic Word Embeddings

In DWE [1], Gaussian word embeddings are learned at each timestep with a temporal diffusion prior:

$$\mathbf{U}_{t+1}|\mathbf{U}_t \sim \mathcal{N}\left(\frac{U_t}{1 + \sigma_t^2/\sigma_0^2}, \frac{1}{\sigma_t^{-2} + \sigma_0^{-2}}I\right),$$

where σ_0^2 and σ_t^2 are hyperparameters of the model.

We derive their skip-gram algorithm for our setting by maximizing the following approximate ELBO:

$$\begin{aligned} \mathcal{L}_{\mathcal{DWE}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = & \sum_{t=1}^T \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{U}_t)} [\log p_{\boldsymbol{\theta}}(\mathbf{X}^t|\mathbf{U}^t)] + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{U}_t)} [\log \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{U}_{t-1})} [p(\mathbf{U}_t|\mathbf{U}_{t-1})]] \\ & - \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{U}_t)} [\log q_{\boldsymbol{\phi}}(\mathbf{U}_t)], \end{aligned} \tag{2}$$

where p_{θ} is parametrized by an LSTM. q_{ϕ} is a variational Gaussian distribution that factorizes as:

$$q_{\phi}(\mathbf{U}) = \prod_{t=1}^T q_{\phi}(\mathbf{U}_t),$$

and ϕ are its parameters.

To learn this model, we sample a mini-batches \mathbf{M} that contains text coming from different training timesteps. We must hence rescale the ELBO in Equation 2. We do so by estimating the probability that a given word appears in a particular mini-batch:

$$\begin{aligned} \mathcal{L}_{minibatch}(\theta, \phi) = & \frac{|\mathbf{X}|}{|\mathbf{M}|} \mathbb{E}_{q_{\phi}(\mathbf{U}^{\mathbf{M}})} \left[\sum_{\mathbf{x} \in \mathbf{M}} \log p_{\theta}(\mathbf{x} | \mathbf{U}^{\mathbf{M}}) \right] \\ & + \sum_{\mathbf{u} \in \mathbf{U}^{\mathbf{M}}} \frac{1}{(1 - (1 - \nu_{\mathbf{u}})^{|\mathbf{M}|})} \sum_{t=1}^T \mathbb{E}_{q_{\phi}(\mathbf{u})} [\log \mathbb{E}_{q_{\phi}(\mathbf{u}_{t-1})} [p(\mathbf{u}_t | \mathbf{u}_{t-1})]] \\ & - \mathbb{E}_{q_{\phi}(\mathbf{u}_t)} [\log q_{\phi}(\mathbf{u}_t)], \end{aligned}$$

where $\mathbf{U}^{\mathbf{M}}$ are the embeddings of words in \mathbf{M} , $\nu_{\mathbf{u}}$ is the apparition frequency of term whose embedding is \mathbf{u} in \mathbf{X} , and $|\mathbf{X}|$ (respectively $|\mathbf{M}|$) is the number of words in \mathbf{X} (\mathbf{M}). In this formulation, gradient computation does not require any approximation, while allowing it to flow through all timesteps.

C.2 DiffTime

The adaptation of the DiffTime baseline [5] is straightforward. It learns a non-linear function d that outputs temporal word embeddings:

$$\mathbf{u}_t = d(\mathbf{u}, t; \phi)$$

where \mathbf{u} is a learned word embedding, t is a scalar timestep, and ϕ are the function's parameters. We refer the reader to the complete paper for more details on the implementation of d .

For recurrent language modeling adaptation, we simply learn jointly the word embeddings \mathbf{U} , the parameters ϕ of d and the parameters θ of an LSTM by maximizing the following likelihood:

$$\mathcal{L}_{\mathcal{DT}}(\theta, \phi, \mathbf{U}) = \prod_{t=1}^T \prod_{\mathbf{x} \in \mathbf{X}^t} \prod_{k=1}^{|\mathbf{x}|-1} p_{\theta}(\mathbf{x}_{k+1} | \mathbf{u}_{1:k}^t).$$

References

1. Bamler, R., Mandt, S.: Dynamic word embeddings. In: ICML (2017)
2. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: ICLR (2014)

3. Krishnan, R.G., Shalit, U., Sontag, D.: Structured inference networks for nonlinear state space models. In: AAAI (2017)
4. Rezende, D.J., Mohamed, S., Wierstra, D.: Stochastic backpropagation and approximate inference in deep generative models. In: ICML (2014)
5. Rosenfeld, A., Erk, K.: Deep neural models of semantic shift. In: NAACL (2018)