

SUPPLEMENTARY MATERIALS

Reinforcement Learning for Bandit Neural Machine Translation with Simulated Human Feedback

A Neural Machine Translation

Our neural machine translation (NMT) model consists of an encoder and a decoder, each of which is a recurrent neural network (RNN). We closely follow (Luong et al., 2015) for the structure of our model. It directly models the posterior distribution $P_{\theta}(\mathbf{y} \mid \mathbf{x})$ of translating a source sentence $\mathbf{x} = (x_1, \dots, x_n)$ to a target sentence $\mathbf{y} = (y_1, \dots, y_m)$:

$$P_{\theta}(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^m P_{\theta}(y_t \mid \mathbf{y}_{<t}, \mathbf{x}) \quad (1)$$

where $\mathbf{y}_{<t}$ are all tokens in the target sentence prior to y_t .

Each local distribution $P_{\theta}(y_t \mid \mathbf{y}_{<t}, \mathbf{x})$ is modeled as a multinomial distribution over the target language’s vocabulary. We compute this distribution by applying a linear transformation followed by a softmax function on the decoder’s output vector \mathbf{h}_t^{dec} :

$$P_{\theta}(y_t \mid \mathbf{y}_{<t}, \mathbf{x}) = \text{softmax}(\mathbf{W}_s \mathbf{h}_t^{dec}) \quad (2)$$

$$\mathbf{h}_t^{dec} = \tanh(\mathbf{W}_o[\tilde{\mathbf{h}}_t^{dec}; \mathbf{c}_t]) \quad (3)$$

$$\mathbf{c}_t = \text{attend}(\tilde{\mathbf{h}}_{1:n}^{enc}, \tilde{\mathbf{h}}_t^{dec}) \quad (4)$$

where $[\cdot; \cdot]$ is the concatenation of two vectors, $\text{attend}(\cdot, \cdot)$ is an attention mechanism, $\tilde{\mathbf{h}}_{1:n}^{enc}$ are all encoder’s hidden vectors and $\tilde{\mathbf{h}}_t^{dec}$ is the decoder’s hidden vector at time step t . We use the “concat” global attention in (Luong et al., 2015).

During training, the encoder first encodes \mathbf{x} to a continuous vector $\Phi(\mathbf{x})$, which is used as the initial hidden vector for the decoder. In our paper, $\Phi(\mathbf{x})$ simply returns the last hidden vector of the encoder. The decoder performs RNN updates to produce a sequence of hidden vectors:

$$\begin{aligned} \tilde{\mathbf{h}}_0^{dec} &= \Phi(\mathbf{x}) \\ \tilde{\mathbf{h}}_t^{dec} &= f_{\theta} \left(\tilde{\mathbf{h}}_{t-1}^{dec}, \left[\mathbf{h}_{t-1}^{dec}; e(y_t) \right] \right) \end{aligned} \quad (5)$$

where $e(\cdot)$ is a word embedding lookup function and y_t is the ground-truth token at time step t . Feeding the output vector \mathbf{h}_{t-1}^{dec} to the next step is known as “input feeding”.

At prediction time, the ground-truth token y_t in Eq 5 is replaced by the model’s own prediction \hat{y}_t :

$$\hat{y}_t = \arg \max_y P_{\theta}(y \mid \hat{\mathbf{y}}_{<t}, \mathbf{x}) \quad (6)$$

In a supervised learning framework, an NMT model is typically trained under the maximum log-likelihood objective:

$$\max_{\theta} \mathcal{L}_{sup}(\theta) = \max_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim D_{tr}} [\log P_{\theta}(\mathbf{y} \mid \mathbf{x})] \quad (7)$$

where D_{tr} is the training set. However, this learning framework is not applicable to bandit learning since ground-truth translations are not available.

References

Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. In *Empirical Methods in Natural Language Processing (EMNLP)*.