

Supplementary Materials: On Tree-Based Neural Sentence Modeling

Haoyue Shi^{†,‡} Hao Zhou[‡] Jiaze Chen[‡] Lei Li[‡]

[†]: School of EECS, Peking University, Beijing, China

hyshi@pku.edu.cn

[‡]: ByteDance AI Lab, Beijing, China

{zhouhao.nlp, chenjiaze, lileilab}@bytedance.com

1 Implementation Details

Our codebase is built on PyTorch 0.3.0.¹ All the sentences was tokenized with SpaCy.²

1.1 Sentence Encoding

We use LSTM based sentence encodings as the extracted features of sentences for downstream classification or generation tasks. We use typical long short term memory (LSTM; Hochreiter and Schmidhuber, 1997) units for linear structures, which can be summarized as:

$$\begin{aligned}f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\\tilde{c}_t &= \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \\o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\c_t &= f_t c_{t-1} + i_t \tilde{c}_t \\h_t &= o_t \tanh(c_t)\end{aligned}$$

where t indicates the time step of a state; h_t is the hidden state and x_t is the input vector. We apply binary tree LSTM units adapted from Zhu et al. (2015) for binary tree LSTMs, which can be summarized as:

$$\begin{aligned}f_l &= \sigma(W_l \cdot [h_l, h_r] + b_l) \\f_r &= \sigma(W_r \cdot [h_l, h_r] + b_r) \\i_t &= \sigma(W_i \cdot [h_l, h_r] + b_i) \\\tilde{c}_t &= \tanh(W_c \cdot [h_l, h_r] + b_c) \\o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\c_t &= f_l c_l + f_r c_r + i_t \tilde{c}_t \\h_t &= o_t \tanh(c_t)\end{aligned}$$

where the subscript t denotes the current state, and l, r denote the left and right child states respectively. We also apply LSTM (Hochreiter and

Schmidhuber, 1997) as leaf-node RNN when necessary.

It is worth noting that left-branching tree LSTM without leaf-node RNN is structurally equivalent to unidirectional LSTM. The only difference between them, which may cause the slight difference on performance, comes from the implementation of LSTM units.

The candidate set of dropout ratio we explore for the task of word-level semantic relation (WSR) is $\{0, 0.1, 0.15, 0.2, 0.3, 0.5\}$.

1.2 Sentence Relation Classification

In the task of sentence relation classification, the feature vector consists of the concatenation of two sentence vectors, their difference, and their element-wise product (Mou et al., 2016):

$$z = \begin{pmatrix} s_1 \\ s_2 \\ s_1 - s_2 \\ s_1 \odot s_2 \end{pmatrix}$$

1.3 Pooling Mechanism

Following (Socher et al., 2011), we apply pooling mechanism to all leaf states (of tree LSTMs) and hidden states. The detailed pooling methods are described as follows.

Max Pooling. Max pooling takes the max value for each dimension

$$\begin{aligned}H &= (h_1, h_2, \dots, h_m) \\s_i &= \max_{j=1}^m h_{j,i} \quad i = 1, 2, \dots, d \\s &= (s_1, s_2, \dots, s_d)^T\end{aligned}$$

where h_i denotes a leaf state in tree LSTMs or a hidden state; $m = 2n - 1$ for tree LSTMs and $m = n$ for linear LSTMs; s denotes the final sentence encoding.

¹<https://pytorch.org/docs/0.3.0>

²<https://spacy.io>

Mean Pooling. Mean pooling (average pooling) takes the average of all hidden states as the sentence representation, which can be summarized as:

$$\mathbf{H} = (\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_m)$$

$$\mathbf{s} = \frac{1}{m} \sum_{i=1}^m \mathbf{h}_i$$

Self-Attention. We follow [Conneau et al. \(2017\)](#) and [Lin et al. \(2017\)](#) to build a self-attentive mechanism, which can be summarized as:

$$\mathbf{H} = (\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_m)$$

$$\mathbf{a} = \text{softmax}(\mathbf{w}_\beta^T \tanh(\mathbf{W}_\alpha \mathbf{H}))$$

$$\mathbf{s} = \mathbf{H} \mathbf{a}^T$$

where \mathbf{a} denotes attention weights computed by learned parameters \mathbf{W}_α and \mathbf{w}_β . In all experiments, \mathbf{w}_β is a 128-d vector.

References

- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. Supervised Learning of Universal Sentence Representations from Natural Language Inference Data. In *Proc. of EMNLP*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation*.
- Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. A Structured Self-Attentive Sentence Embedding. In *Proc. of ICLR*.
- Lili Mou, Rui Men, Ge Li, Yan Xu, Lu Zhang, Rui Yan, and Zhi Jin. 2016. Natural language inference by tree-based convolution and heuristic matching. In *Proc. of ACL*.
- Richard Socher, Eric H. Huang, Jeffrey Pennington, Christopher D. Manning, and Andrew Y. Ng. 2011. Dynamic Pooling and Unfolding Recursive Autoencoders for Paraphrase Detection. In *Proc. of NIPS*.
- Xiaodan Zhu, Parinaz Sobihani, and Hongyu Guo. 2015. Long Short-Term Memory Over Recursive Structures. In *Proc. of ICML*.