

Latent Constituents via Self-Attention

August 2, 2018

1 Introduction

2 Problem

We would like to learn a generative model over source sentences $\mathbf{x} = \{x_0, x_1, \dots\}$, using a distribution over latent trees \mathbf{z} . (Shen et al., 2018) Yin et al. (2018)

3 PRPN (Shen et al., 2018)

The Parsing-Reading-Predict Network (Shen et al., 2018) is an approach to unsupervised tree induction that uses a soft approximation to a latent variable to a degree of success. The paper is inspired by the following hypothesis: given a binary tree, a token x_t only requires information up to a token x_{l_t} that satisfies either of the following conditions:

- (a) its leftmost sibling
- (b) or, if the token x_t is a leftmost child, its parent’s left sibling’s leftmost child.

Although the hypothesis itself is not tested and seems to only serve as inspiration, the model does see empirical success in the task of language modeling. The model is realized through the following insight: given a ranking of tokens $\mathbf{x} = \{x_0, \dots, x_T\}$, recursively splitting \mathbf{x} using the following procedure induces a binary tree where the first token to the left of token x_t that has higher rank, denoted x_{l_t} , also satisfies condition (a) or (b): given token i with the next highest rank, create a subtree $(x_{<i}, (x_i, x_{>i}))$ and recursively perform this procedure until only terminal nodes remain.

[Example on board]

3.1 The Model

Let x_t be the current token, $\mathbf{x}_{0:t-1}$ all previous tokens, z_t the prediction for the current score, and $\tilde{\mathbf{z}}_{0:t-1} \in \mathbb{R}_+^{t-1}$ be all previous scores. The model takes the form a language model, where the distribution over the next token

$$p(x_t \mid \mathbf{x}_{0:t-1}, \tilde{\mathbf{z}}_{0:t-1}, z_t) = f([h_{l_t:t-1}, h_t])$$

is parameterized by either a linear projection or additional residual blocks applied to a concatenation of the output of modified LSTMN (LSTM Memory-Network) h_t and a convex combination of the previous LSTMN outputs (attention over $h_{l_t:t-1}$). f is referred to as the Predict Network, and we will cover the attention mechanism shortly. The LSTMN, referred to as the Reading Network, is as follows: at each time step, the hidden and cell input are given by

$$\{\tilde{h}_t, \tilde{c}_t\} = \sum_{i=1}^{t-1} s_i^t \{h_i, c_i\},$$

a convex combination of the previous hidden and cell outputs. Then

$$h_t, c_t = \text{LSTM}(\tilde{h}_t, \tilde{c}_t).$$

The coefficients s_i^t are computed by first predicting a key

$$k = W_h h_{t-1} + W_x x_t$$

then querying all previous hidden states in order to compute the intermediate attention

$$\tilde{s}_i^t = \text{softmax}\left(\frac{h_i k_t^T}{\sqrt{\delta_k}}\right),$$

where δ_k is the dimension of the hidden state. This same attention mechanism is used in the Predict Network f . We then incorporate a form of multiplicative gating g_i^t that is a function of the induced ranking of the tokens, which we described in the previous section. Finally,

$$s_i^t = \frac{g_i^t}{\sum_j g_j^t} \tilde{s}_i^t,$$

so at every timestep t the self-attention is renormalized using the gates $g_{l_t:t-1}^t$.

3.2 Gating Self-attention

Recall that the gates g_i^t used modulate self-attention also indirectly induce tree structure (due to the previously stated insight that modulating self-attention using a ranking induces a binary tree). At every timestep for token x_t , Shen et al. (2018) define a latent variable l_t which corresponds to the index of the token satisfying either condition (a) or (b). We then would have

$$g_i^t = \begin{cases} 1, & l_t \leq i < t \\ 0, & \text{otherwise.} \end{cases}$$

where we allow the model at timestep t to only attend up to the token at l_t . We would then renormalize the attention accordingly. In order to calculate the gates g_i^t we must model $p(l_t \mid \mathbf{x}_{0:t})$ (Shen et al. (2018) use the posterior distribution in their notation). They propose to model

$$p(l_t = i \mid \mathbf{x}_{0:t-1}) = (1 - \alpha_i^t) \prod_{j=i+1}^{t-1} \alpha_j^t$$

using a stick-breaking process. Rather than resorting to approximate inference, they instead use the expected value of g_i^t

$$\mathbb{E}[g_i^t] = F_{l_t}(l_t \leq i \mid \mathbf{x}_{0:t-1}) = \prod_{j=i+1}^{t-1} \alpha_j^t.$$

Shen et al. (2018) derive the expectation in the appendix of their paper. The $\alpha_j^t = \sigma(d_t - d_j)$, where σ is the hard sigmoid function. The $d_j \in \mathbb{R}_+$ are given by a CNN over the token embeddings.

The network used to compute the gates is called the Parsing Network. When computing $p(x_t \mid \dots)$, we cannot use the posterior score d_t , so an estimate is used based on the previous k words, where k is the kernel width of the convolution.

4 Comparison to CCM (Klein and Manning, 2002a)

(Klein and Manning, 2002b) (Golland et al., 2012; Huang et al., 2012)

5 Alternative Latent Ranking Model

5.1 Generative Model

PRPN (Shen et al., 2018) choose to view l_t as a latent variable. There are at least a few other choices:

- Model the scores $d_t \sim \mathcal{N}(\mu_t, \sigma_t)$ or $d_t \sim \text{Gamma}$ and use the Plackett-Luce ranking distribution.
- Model the comparisons $p(d_t < d_i)$ as order statistics of Gammas.
- Model the permutation matrix $Z \sim \mathcal{B}_n$.
- Model $l_t \sim \text{Cat}$.

Parameterize with $d_t \sim \mathcal{N}(\mu_t, \sigma_t)$. Reparameterize comparisons with gumbel softmax? Leave all self attentions as is?

- $p(z)$
- $p(x|z)$

6 Training and Inference

References

- Dave Golland, John DeNero, and Jakob Uszkoreit. A feature-rich constituent context model for grammar induction. In *The 50th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference, July 8-14, 2012, Jeju Island, Korea - Volume 2: Short Papers*, pages 17–22. The Association for Computer Linguistics, 2012. ISBN 978-1-937284-25-1. URL <http://www.aclweb.org/anthology/P12-2004>.
- Yun Huang, Min Zhang, and Chew Lim Tan. Improved constituent context model with features. In *PACLIC*, pages 564–573. PACLIC 26 Organizing Committee and PACLIC Steering Committee / ACL / Faculty of Computer Science, Universitas Indonesia, 2012.
- Dan Klein and Christopher D. Manning. A generative constituent-context model for improved grammar induction. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA.*, pages 128–135. ACL, 2002a. URL <http://www.aclweb.org/anthology/P02-1017.pdf>.
- Dan Klein and Christopher D. Manning. A generative constituent-context model for improved grammar induction. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, pages 128–135, Stroudsburg, PA, USA, 2002b. Association for Computational Linguistics. doi: 10.3115/1073083.1073106. URL <https://doi.org/10.3115/1073083.1073106>.
- Yikang Shen, Zhouhan Lin, Chin wei Huang, and Aaron Courville. Neural language modeling by jointly learning syntax and lexicon. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rkgOLb-0W>.
- Pengcheng Yin, Chunting Zhou, Junxian He, and Graham Neubig. StructVAE: Tree-structured latent variable models for semi-supervised semantic parsing. In *The 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, Melbourne, Australia, July 2018. URL <https://arxiv.org/abs/1806.07832v1>.