

Appendix of “Lattice CNNs for Matching Based Chinese Question Answering”

Yuxuan Lai¹, Yansong Feng^{1,*}, Xiaohan Yu¹
Zheng Wang², Kun Xu³, Dongyan Zhao¹

¹Institute of Computer Science and Technology, Peking University, China

²School of Computing and Communications, Lancaster University, UK ³Tencent AI Lab

¹{erutan, fengyansong, yuxiaohan, zhaodongyan}@pku.edu.cn

²z.wang@lancaster.ac.uk ³syxu828@gmail.com

Appendix

Formula Derivation

$$F_w = g\{f(\mathbf{W}_c(\mathbf{v}_{w_1} : \dots : \mathbf{v}_{w_n}) + \mathbf{b}_c^T) | \forall i, w_i \in V, (w_i, w_{i+1}) \in E, w_{\lceil \frac{n+1}{2} \rceil} = w\} \quad (1)$$

Lattice based CNN layer with average pooling, derived from Eq.1 (the same as Eq.3 in the paper). Omit the condition $w_i \in V$, set $n = 3$ and mark \mathbf{W}_c as $(\mathbf{W}_{c1} : \mathbf{W}_{c2} : \mathbf{W}_{c3})$, where the size of \mathbf{W}_{ci} is $(m', m) \forall i \in \{1, 2, 3\}$, and suppose the pooling mode is average, and ignore the activation function.

$$\begin{aligned} F_w &= \text{Pool}\{f(\mathbf{W}_c(\mathbf{v}_{w_1} : \dots : \mathbf{v}_{w_n}) + \mathbf{b}_c^T) | \\ &\quad \forall i, w_i \in V, (w_i, w_{i+1}) \in E, w_{\lceil \frac{n+1}{2} \rceil} = w\} \\ &= \text{ave}\{(\mathbf{W}_{c1} : \mathbf{W}_{c2} : \mathbf{W}_{c3})(\mathbf{v}_{w_1} : \mathbf{v}_{w_2} : \mathbf{v}_{w_3}) + \mathbf{b}_c^T | \\ &\quad \forall i, (w_i, w_{i+1}) \in E, w_2 = w\} \\ &= \text{ave}\{w_{c1}v_{w_1} + w_{c2}v_w + w_{c3}v_{w_3} + \mathbf{b}_c^T | \\ &\quad (w_1, w), (w, w_3) \in E\} \\ &= w_{c2}v_w + \mathbf{b}_c^T + \frac{\sum_{w_1} w_{c1}v_{w_1}}{t_1} + \frac{\sum_{w_3} w_{c3}v_{w_3}}{t_2} \end{aligned} \quad (2)$$

, where t_1 is the number of previous words of w and t_2 is the number of next words.

The same as Eq.2, but have maximum pooling.

$$\begin{aligned} F_w &= \text{Pool}\{f(\mathbf{W}_c(\mathbf{v}_{w_1} : \dots : \mathbf{v}_{w_n}) + \mathbf{b}_c^T) | \\ &\quad \forall i, w_i \in V, (w_i, w_{i+1}) \in E, w_{\lceil \frac{n+1}{2} \rceil} = w\} \\ &= \max\{(\mathbf{W}_{c1} : \mathbf{W}_{c2} : \mathbf{W}_{c3})(\mathbf{v}_{w_1} : \mathbf{v}_{w_2} : \mathbf{v}_{w_3}) + \mathbf{b}_c^T | \\ &\quad \forall i, (w_i, w_{i+1}) \in E, w_2 = w\} \\ &= \max\{w_{c1}v_{w_1} + w_{c2}v_w + w_{c3}v_{w_3} + \mathbf{b}_c^T | \\ &\quad (w_1, w), (w, w_3) \in E\} \\ &= w_{c2}v_w + \mathbf{b}_c^T + \max\{w_{c1}v_{w_1} | \forall w_1\} + \max\{w_{c3}v_{w_3} | \forall w_3\} \end{aligned} \quad (3)$$

Directed graph convolutional networks (DGCs), derived from Eq.2 in (Vashishth et al. 2018). We also suppose the pooling mode is average(not sum, which is in original paper and the first line of our derivation) and ignore the activation function.

$$\begin{aligned} h_v^{k+1} &= f(\sum_{u \in \mathcal{N}} (\mathbf{W}_{l(u,v)}^k h_u^k + \mathbf{b}_{l(u,v)}^k)) \\ &= \text{ave}\{(\mathbf{W}_{l(u,v)}^k h_u^k + \mathbf{b}_{l(u,v)}^k) | u \in \mathcal{N}\} \\ &= \text{ave}\{W_{c1}Vw_1 + \mathbf{b}_1^T, W_{c3}Vw_3 + \mathbf{b}_3^T, W_{c2}Vw + \mathbf{b}_2^T | \\ &\quad (w_1, w), (w, w_3) \in E\} \\ &= \frac{1}{1 + t_1 + t_2} (w_{c2}v_w + (\mathbf{b}_2^T + \mathbf{b}_1^T * t_1 + \mathbf{b}_3^T * t_2) \\ &\quad + \sum_{w_1} w_{c1}v_{w_1} + \sum_{w_3} w_{c3}v_{w_3}) \end{aligned} \quad (4)$$

, where W_{c1}, \mathbf{b}_1^T are parameters for forward connections, W_{c3}, \mathbf{b}_3^T are parameters for backward connections, and W_{c2}, \mathbf{b}_2^T are parameters for self-loops.

The same as Eq.4, but have maximum pooling.

$$\begin{aligned} h_v^{k+1} &= f(\sum_{u \in \mathcal{N}} (\mathbf{W}_{l(u,v)}^k h_u^k + \mathbf{b}_{l(u,v)}^k)) \\ &= \max\{(\mathbf{W}_{l(u,v)}^k h_u^k + \mathbf{b}_{l(u,v)}^k) | u \in \mathcal{N}\} \\ &= \max\{W_{c1}Vw_1 + \mathbf{b}_1^T, W_{c3}Vw_3 + \mathbf{b}_3^T, W_{c2}Vw + \mathbf{b}_2^T | \\ &\quad (w_1, w), (w, w_3) \in E\} \\ &= \max\{\max\{W_{c1}Vw_1 + \mathbf{b}_1^T | \forall w_1\}, \\ &\quad W_{c2}Vw + \mathbf{b}_2^T, \max\{W_{c3}Vw_3 + \mathbf{b}_3^T | \forall w_3\}\} \end{aligned} \quad (5)$$

References

Vashishth, S.; Dasgupta, S. S.; Ray, S. N.; and Talukdar, P. 2018. Dating documents using graph convolution networks. In *ACL 2018*, 1605–1615.