

A APPENDIX

A.1 A PROOF OF MLM ON CODE-SWITCHED DATASET

Proposition Training model with standard masked language modeling on source-language, target-language and source-target code-switched sentences is approximately optimizing:

$$\mathbf{c}_{x_i}^{src} \sim \mathbf{e}_{x_i} \sim \mathbf{e}_{y_i} \sim \mathbf{c}_{y_i}^{tgt}$$

where $\mathbf{c}_{x_i}^{src}$ represents the contextualized embedding of token x_i in source sentence and $\mathbf{c}_{y_i}^{tgt}$ represents the contextualized embedding of token y_i in target sentence. And \mathbf{e}_{x_i} and \mathbf{e}_{y_i} are word embeddings of token x_i and y_i in vocabulary.

Note Under the existing conditions, we can not derive a strict bound but an approximate conclusion. Given $\mathbf{a} \in \mathbb{R}^n, \mathbf{b} \in \mathbb{R}^n, \mathbf{c} \in \mathbb{R}^n$ we assume $\mathbf{a} \sim \mathbf{b}$, if the projection components of \mathbf{a} and \mathbf{b} onto another vector \mathbf{c} are the same: $\mathbf{a} \cdot \mathbf{c} = \mathbf{b} \cdot \mathbf{c}$. We think this is approximately reasonable for word alignment task, because for word alignments method two words are aligned as long as their similarity is higher than other words in two parallel sentences and doesn't need to exceed a fixed number. And in section 3.3, we give a corollary. Subsequent experiments prove that our assumption is reasonable.

Proof We denote the embeddings of the corresponding original tokens as $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_L$. The MLM objective $\mathcal{L}_{MLM}(\mathbf{x})$ can be formulated as:

$$-\frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \log \frac{\exp(\mathbf{m}_i \cdot \mathbf{e}_i)}{\sum_{k=1}^{|\mathcal{V}|} \exp(\mathbf{m}_i \cdot \mathbf{e}_k)} = -\frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \log \sum_{k=1}^{|\mathcal{V}|} \exp(\mathbf{m}_i \cdot \mathbf{e}_k - \mathbf{m}_i \cdot \mathbf{e}_i) \quad (1)$$

where \mathcal{M} denotes the set of masked tokens and $|\mathcal{V}|$ is the size of vocabulary \mathcal{V} . \mathbf{m}_i is hidden state of the last layer at the masked position, and can be regarded as a fusion of contextualized representations of surrounding tokens. Given two sentences: one source-language sentence $\mathbf{x} = \langle x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n \rangle$ of length n and its code-switched sentence $\mathbf{x}' = \langle x_1, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n \rangle$, where $\langle x_i, y_i \rangle$ is aligned pair. If we only mask x_i in the \mathbf{x} and y_i in the \mathbf{x}' , then $\mathbf{x}_{mask} = \langle x_1, \dots, x_{i-1}, \langle mask \rangle, x_{i+1}, \dots, x_n \rangle = \langle x_1, \dots, x_{i-1}, \langle mask \rangle, x_{i+1}, \dots, x_n \rangle = \mathbf{x}'_{mask}$, the loss function can be written as

$$\mathcal{L}_{MLM} = \mathcal{L}_{\mathbf{x}} + \mathcal{L}_{\mathbf{x}'} = -\frac{1}{2} (\log \sum_{k=1}^{|\mathcal{V}|} \exp(\mathbf{m} \cdot \mathbf{e}_k - \mathbf{m} \cdot \mathbf{e}_{x_i}) + \log \sum_{k=1}^{|\mathcal{V}|} \exp(\mathbf{m} \cdot \mathbf{e}_k - \mathbf{m} \cdot \mathbf{e}_{y_i})) \quad (2)$$

This inequality below is easily proved.

$$\max \{x_1, \dots, x_n\} \leq \log \sum_{i=0}^n e^{x_i} \leq \max \{x_1, \dots, x_n\} + \log n \quad (3)$$

So for

$$\mathcal{L}_{\mathbf{x}} = -\log \sum_{k=1}^{|\mathcal{V}|} \exp(\mathbf{m} \cdot \mathbf{e}_k - \mathbf{m} \cdot \mathbf{e}_{x_i}) \quad (4)$$

We have:

$$\max \begin{pmatrix} \mathbf{m} \cdot \mathbf{e}_0 - \mathbf{m} \cdot \mathbf{e}_{x_i} \\ \vdots \\ \mathbf{m} \cdot \mathbf{e}_{x_{i-1}} - \mathbf{m} \cdot \mathbf{e}_{x_i} \\ 0 \\ \mathbf{m} \cdot \mathbf{e}_{x_{i+1}} - \mathbf{m} \cdot \mathbf{e}_{x_i} \\ \vdots \\ \mathbf{m} \cdot \mathbf{e}_{|\mathcal{V}|} - \mathbf{m} \cdot \mathbf{e}_{x_i} \end{pmatrix} \leq \log \sum_{k=1}^{|\mathcal{V}|} e^{\mathbf{m} \cdot \mathbf{e}_k - \mathbf{m} \cdot \mathbf{e}_{x_i}} \leq \max \begin{pmatrix} \mathbf{m} \cdot \mathbf{e}_0 - \mathbf{m} \cdot \mathbf{e}_{x_i} \\ \vdots \\ \mathbf{m} \cdot \mathbf{e}_{x_{i-1}} - \mathbf{m} \cdot \mathbf{e}_{x_i} \\ 0 \\ \mathbf{m} \cdot \mathbf{e}_{x_{i+1}} - \mathbf{m} \cdot \mathbf{e}_{x_i} \\ \vdots \\ \mathbf{m} \cdot \mathbf{e}_{|\mathcal{V}|} - \mathbf{m} \cdot \mathbf{e}_{x_i} \end{pmatrix} + \log n \quad (5)$$

In Ineq.5, 0 is fixed value. So when training model with this loss function, model is optimized to learn $\mathbf{m} \cdot \mathbf{e}_k - \mathbf{m} \cdot \mathbf{e}_{x_i} \leq 0, \forall k \in |\mathcal{V}|$. In other words, $\mathbf{m} \cdot \mathbf{e}_k \leq \mathbf{m} \cdot \mathbf{e}_{x_i}, \forall k \in |\mathcal{V}|$. When $k = y_i$,

we have $\mathbf{m} \cdot e_{y_i} \leq \mathbf{m} \cdot e_{x_i}$. Similarly, for $\mathcal{L}_{\mathbf{x}'}$, we have $\mathbf{m} \cdot e_{x_i} \leq \mathbf{m} \cdot e_{y_i}$. So when training model with loss function $\mathcal{L}_{\text{MLM}} = \mathcal{L}_{\mathbf{x}} + \mathcal{L}_{\mathbf{x}'}$, model will be optimized to learn $\mathbf{m} \cdot e_{x_i} = \mathbf{m} \cdot e_{y_i}$. This equation can't ensure $e_{x_i} = e_{y_i}$ but $e_{x_i} \sim e_{y_i}$ to some extent. For standard masked language modeling, there is a probability that the original token will not be masked and we use $c_{x_i}^{src}$ to represent the hidden state of the last layer, which is the contextualized embedding of token x_i . So we have $c_{x_i}^{src} \cdot e_k \leq c_{x_i}^{src} \cdot e_{x_i}$, $\forall k \in |\mathcal{V}|$. Obviously, $e_{x_i} \sim c_{x_i}^{src}$. Similarly, if we consider target language sentence $\mathbf{y} = \langle y_1, \dots, y_{i-1}, y_i, y_{i+1}, \dots, y_n \rangle$, we have $e_{y_i} \sim c_{y_i}^{tgt}$. So training model with masked language modeling on source-language, target-language and source-target code-switched sentences is approximately optimizing:

$$c_{x_i}^{src} \sim e_{x_i} \sim e_{y_i} \sim c_{y_i}^{tgt}$$