

NPFL116 Compendium of Neural Machine Translation

Sequence-Level Training

April 5, 2017

Jindřich Helcl, Jindřich Libovický



Charles University in Prague
Faculty of Mathematics and Physics
Institute of Formal and Applied
Linguistics



Word-Level Training

- ▶ Likelihood of a sentence in word-level training:

$$p(\mathbf{y}|\mathbf{x}, \theta) = \prod_{t=1}^T p(y_t|\mathbf{y}_{<t}, \mathbf{x}, \theta)$$

- ▶ Log-likelihood as a loss function:

$$\begin{aligned}\mathcal{L}(\theta) &= \sum_{i=1}^N \log p(\mathbf{y}_i|\mathbf{x}_i, \theta) \\ &= \sum_{i=1}^N \sum_{t=1}^{T_i} \log p(y_{it}|\mathbf{y}_{i,<t}, \mathbf{x}_i, \theta)\end{aligned}$$

Log-likelihood

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^N \sum_{t=1}^{T_i} \log p(y_{it} | \mathbf{y}_{i,<t}, \mathbf{x}_i, \boldsymbol{\theta})$$

- ▶ Fast, yields good results
- ▶ *Differentiable!*

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_i} = \sum_{i=1}^N \sum_{t=1}^{T_i} \frac{\partial p(y_{it} | \mathbf{y}_{i,<t}, \mathbf{x}_i, \boldsymbol{\theta}) / \partial \boldsymbol{\theta}_i}{p(y_{it} | \mathbf{y}_{i,<t}, \mathbf{x}_i, \boldsymbol{\theta})}$$

Log-likelihood – problems

- ▶ Training objective different from the evaluation metric
- ▶ Loss function defined on word-level
- ▶ Output word distribution compared with an one-hot distribution
- ▶ Suffers from exposure bias

Sentence-level Losses

- ▶ Score the output sentence as a whole
- ▶ Solve the exposure bias problem
- ▶ Many metrics out there: BLEU, METEOR, ...
- ▶ Good correlation with human judgement

- ▶ Drawback: Although differentiable, the derivatives are locally constant

- ▶ Solution?

Sentence-level Training

- ▶ Problem: locally constant derivatives
- ▶ Cause: selecting the best word during decoding
- ▶ Idea: Score the distribution over possible sentences rather than the model output
- ▶ Can we still somehow use metrics like BLEU?

Sentence-level Training

- ▶ Let the loss be the expected value of the scoring function r over all possible outputs \mathcal{Y} w. r. t. reference translation y^* :

$$\mathbb{E}_{y \in \mathcal{Y}} [r(y, y^*)] = \sum_{y \in \mathcal{Y}} p(y|x, \theta) r(y, y^*)$$

- ▶ good model gives high score to good sentences, low score to bad sentences

Minimum Risk Training (1)

Shen et al., 2016 (<https://arxiv.org/abs/1512.02433>)

- ▶ Minimum risk training uses the expected score to calculate the risk:

$$\mathcal{R}(\theta) = \sum_{i=1}^N \mathbb{E}_{y \in \mathcal{Y}} [r(y, y^*)]$$

- ▶ Nice derivative

Minimum Risk Training (2)

- ▶ Problem: The space of all possible outputs $\mathcal{Y}(x)$ for input sentence x is way too large
- ▶ Approximate the expected value by using only a few samples from the distribution

$$|\mathcal{Y}(x)| \gg |\mathcal{S}(x)| = \text{usually around } 100$$

Resampling

$$\begin{aligned}\bar{R}(\theta) &= \sum_{s=1}^S \mathbb{E} \Delta(\mathbf{y}, \mathbf{y}^*) \\ &\approx \sum_{s=1}^S \sum_{\mathbf{y} \in \mathcal{S}(\mathbf{x})} Q(\mathbf{y}|\mathbf{x}, \theta, \alpha) \Delta(\mathbf{y}, \mathbf{y}^*)\end{aligned}$$

$$Q(\mathbf{y}|\mathbf{s}, \theta, \alpha) = \frac{\mathbf{p}(\mathbf{y}|\mathbf{s}, \theta)^\alpha}{\sum_{\mathbf{y}' \in \mathcal{S}(\mathbf{x})} \mathbf{p}(\mathbf{y}'|\mathbf{x}, \theta)}$$

Sequence-level Training using Reinforcement Learning

Ranzato et al., 2015 (<https://arxiv.org/abs/1511.06732>)

- ▶ In minimum risk training, we approximate the expected score
- ▶ Here: approximate the gradients

$$\text{chain rule: } \frac{\partial \mathcal{L}}{\partial \theta} = \sum_t \frac{\partial \mathcal{L}}{\partial \mathbf{o}_t} \cdot \frac{\partial \mathbf{o}_t}{\partial \theta}$$

- ▶ Use the REINFORCE algorithm (Williams, 1992):

$$\frac{\partial \mathcal{L}}{\partial \mathbf{o}_t} \approx (r(\mathbf{y}, \mathbf{y}^*) - \bar{r}_{t+1}) (\mathbf{p}(\mathbf{y}_{t+1} | \mathbf{y}_{<t}, \mathbf{x}, \theta) - \mathbb{1}(\mathbf{y}_{t+1}))$$