

Towards Decoding as Continuous Optimisation in Neural Machine Translation

Cong Duy Vu Hoang

University of Melbourne

Melbourne, VIC, Australia

`vhoang2@student.unimelb.edu.au`

Gholamreza Haffari

Monash University

Clayton, VIC, Australia

`gholamreza.haffari@monash.edu`

Trevor Cohn

University of Melbourne

Melbourne, VIC, Australia

`t.cohn@unimelb.edu.au`

Neural Machine Translation

$$\begin{aligned} P_{\Theta}(\mathbf{y}|\mathbf{x}) &= \sum_{i=1}^{|\mathbf{y}|} \log P_{\Theta}(y_i|\mathbf{y}_{<i}, \mathbf{x}) \\ y_i|\mathbf{y}_{<i}, \mathbf{x} &\sim \text{softmax}(\mathbf{f}(\Theta, \mathbf{y}_{<i}, \mathbf{x})) \end{aligned}$$

In this paper:

$$\begin{aligned} \mathbf{f}(\Theta, \mathbf{y}_{<i}, \mathbf{x}) &= \mathbf{W}_o \cdot \text{MLP}(\mathbf{c}_i, \mathbf{E}_T^{y_{i-1}}, \mathbf{g}_i) + \mathbf{b}_o \\ \mathbf{g}_i &= \text{RNN}_{dec}^{\phi}(\mathbf{c}_i, \mathbf{E}_T^{y_{i-1}}, \mathbf{g}_{i-1}) \end{aligned}$$

Training objective:

$$\Theta^* := \operatorname{argmax}_{\Theta} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log P_{\Theta}(\mathbf{y} | \mathbf{x}).$$

Discrete decoding

- Decoding objective:

$$\begin{aligned} \mathbf{y}^* = \arg \min_{y_1, \dots, y_\ell} & \sum_{i=1}^{\ell} -\log P_{\Theta}(y_i \mid \mathbf{y}_{<i}, \mathbf{x}) \\ \text{s.t.} \quad & \forall i \in \{1 \dots \ell\} : y_i \in V_T. \end{aligned}$$

- Re-write:

$$\begin{aligned} \arg \min_{\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_\ell} & - \sum_{i=1}^{\ell} \tilde{\mathbf{y}}_i \cdot \log \text{softmax}(\mathbf{f}(\Theta, \tilde{\mathbf{y}}_{<i}, \mathbf{x})) \\ \text{s.t.} \quad & \forall i \in \{1 \dots \ell\} : \tilde{\mathbf{y}}_i \in \mathbb{I}^{|V_T|} \end{aligned} \quad (5)$$

where $\tilde{\mathbf{y}}_i$ are vectors using the one-hot representation of the target words $\mathbb{I}^{|V_T|}$.

discrete optimization problem

Continuous decoding

$$\begin{aligned} \arg \min_{\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_\ell} & - \sum_{i=1}^{\ell} \hat{\mathbf{y}}_i \cdot \log \text{softmax}(\mathbf{f}(\Theta, \hat{\mathbf{y}}_{<i}, \mathbf{x})) \\ \text{s.t.} \quad & \forall i \in \{1 \dots \ell\} : \hat{\mathbf{y}}_i \in \Delta_{|V_T|} \end{aligned}$$

where $\Delta_{|V_T|}$ is the $|V_T|$ -dimensional probability simplex, i.e., $\{\hat{\mathbf{y}}_i \in [0, 1]^{|V_T|} : \|\hat{\mathbf{y}}_i\|_1 = 1\}$. Intuitively, this amounts to replacing $\mathbf{E}_T^{y_i}$ with the *expected* embedding of target language words $\mathbb{E}_{\hat{\mathbf{y}}_i(w)}[\mathbf{E}_T^w]$ under the distribution $\hat{\mathbf{y}}_i$ in the NMT model.

\mathbf{y}_i is a distribution over words

Generation words: choose the highest probability

Optimization method

- Exponentiated Gradient (EG)

Problem:

$$\begin{aligned} & \arg \min_{\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_\ell} Q(\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_\ell) \\ & \text{s.t. } \forall i \in \{1 \dots \ell\} : \hat{\mathbf{y}}_i \in \Delta_{|V_T|} \end{aligned}$$

where $Q(\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_\ell)$ is defined as

$$- \sum_{i=1}^{\ell} \hat{\mathbf{y}}_i \cdot \log \text{softmax}(\mathbf{f}(\Theta, \hat{\mathbf{y}}_{<i}, \mathbf{x})).$$

Algorithm:

$$\forall w \in V_T : \hat{\mathbf{y}}_i^t(w) = \frac{1}{Z_i^t} \hat{\mathbf{y}}_i^{t-1}(w) \exp(-\eta \nabla_{i,w}^{t-1})$$

where η is the step size, $\nabla_{i,w}^{t-1} = \frac{\partial Q(\hat{\mathbf{y}}_1^{t-1}, \dots, \hat{\mathbf{y}}_\ell^{t-1})}{\partial \hat{\mathbf{y}}_i(w)}$
and Z_i^t is the normalisation constant

$$Z_i^t = \sum_{w \in V_T} \hat{\mathbf{y}}_i^{t-1}(w) \exp(-\eta \nabla_{i,w}^{t-1}).$$

Algorithm 1 The EG Algorithm for Decoding by Optimisation

- 1: For all i initialise $\hat{\mathbf{y}}_i^0 \in \Delta_{|V_T|}$
 - 2: **for** $t = 1, \dots, \text{MaxIter}$ **do**
 - 3: For all i, w : calculate $\nabla_{i,w}^{t-1} = \frac{\partial Q(\hat{\mathbf{y}}_1^{t-1}, \dots, \hat{\mathbf{y}}_\ell^{t-1})}{\partial \hat{\mathbf{y}}_i(w)}$
 - ▷ $Q(\cdot)$ is defined as eqn (6)
 - ▷ using backpropagation
 - 4: For all i, w : update $\hat{\mathbf{y}}_i^t(w) \propto \hat{\mathbf{y}}_i^{t-1}(w) \cdot \exp(-\eta \nabla_{i,w}^{t-1})$
 - ▷ η is the step size
 - 5: **return** $\arg \min_t Q(\hat{\mathbf{y}}_1^t, \dots, \hat{\mathbf{y}}_\ell^t)$
-

Optimization method

- Stochastic Gradient Descent (SGD)

Make sure the simplex constraints:

$$\hat{\mathbf{y}}_i = \text{softmax}(\hat{\mathbf{r}}_i).$$

$$\arg \min_{\hat{\mathbf{r}}_1, \dots, \hat{\mathbf{r}}_\ell} - \sum_{i=1}^{\ell} \text{softmax}(\hat{\mathbf{r}}_i) \cdot \log \text{softmax}(\mathbf{f}(\mathbf{\Theta}, \hat{\mathbf{y}}_{<i}, \mathbf{x}))$$

Algorithm 2 The SGD Algorithm for Decoding by Optimisation

- 1: For all i initialise $\hat{\mathbf{r}}_i^0$
 - 2: **for** $t = 1, \dots, \text{MaxIter}$ **do**
 - $\triangleright Q(.)$ is defined in eqn (6) and $\hat{\mathbf{y}}_i = \text{softmax}(\hat{\mathbf{r}}_i)$
 - 3: For all i, w : calculate $\nabla_{i,w}^{t-1} = \sum_{w' \in V_T} \frac{\partial Q(\hat{\mathbf{y}}_1^{t-1}, \dots, \hat{\mathbf{y}}_\ell^{t-1})}{\partial \hat{\mathbf{y}}_i(w')} \frac{\partial \hat{\mathbf{y}}_i(w')}{\partial \hat{\mathbf{r}}_i(w)}$ \triangleright using backpropagation
 - 4: For all i, w : update $\hat{\mathbf{r}}_i^t(w) = \hat{\mathbf{r}}_i^{t-1}(w) - \eta \nabla_{i,w}^{t-1}$ $\triangleright \eta$ is the step size
 - 5: **return** $\arg \min_t Q(\text{softmax}(\hat{\mathbf{r}}_1^t), \dots, \text{softmax}(\hat{\mathbf{r}}_\ell^t))$
-

Decoding in Extended NMT

- Allowing decoding for richer global models, for which there is no effective means of greedy decoding or beam search.
- Bidirectional Ensemble:

$$\mathcal{C}_{+\text{bidir}} := -\alpha \log P_{\Theta_{\leftarrow}}(\mathbf{y} \mid \mathbf{x}) \\ - (1 - \alpha) \log P_{\Theta_{\rightarrow}}(\mathbf{y} \mid \mathbf{x});$$

- Bilingual Ensemble:

$$\mathcal{C}_{+\text{biling}} := -\alpha \log P_{\Theta_{s \rightarrow t}}(\mathbf{y} \mid \mathbf{x}) \\ - (1 - \alpha) \log P_{\Theta_{s \leftarrow t}}(\mathbf{x} \mid \mathbf{y});$$

Experiment

- Main results:

	BTEC	TEDTalks	WMT
	zh→en	de→en	de→en
$\text{gdec}_{\text{left-to-right}}$	35.98	23.16	24.41
$\text{gdec}_{\text{right-to-left}}$	35.86	21.95	23.59
$\text{EGdec}_{\text{greedy init}}$	36.34	23.28	24.63
+bidirectional	36.67	23.91	25.37[†]
+bilingual	36.88[†]	24.01[†]	25.21
$\text{bdec}_{\text{left-to-right}}$	38.02	23.95	26.69
$\text{bdec}_{\text{right-to-left}}$	37.38	23.13	26.11
$\text{EGdec}_{\text{beam init}}$	38.38	24.02	26.66
+bidirectional	39.13[†]	24.72[†]	27.34[†]
+bilingual	38.25	24.60	26.82

Table 3: The BLEU evaluation results across evaluation datasets for EG algorithm variants against the baselines; **bold**: statistically significantly better than the best greedy or beam baseline, [†]: best performance on dataset.