

Semi-supervised Multitask Learning for Sequence Labeling

Marek Rei

The ALTA Institute

Computer Laboratory

University of Cambridge

United Kingdom

Motivation

- Due to the Sparseness of labels, many available data can't be utilized effectively.
- Due to the limitation of the single task, many features can't be learnt effectively.
- Due to the limitation of distribution of data of the single task, it can occur over-fitting easily.

Solution-Multitask

- Multitask can incorporate new supervise signals so that learn the data representation better.
- Multitask can learn the feature representation not easily learnt in task A by incorporating task B.
- Multitask can reduce the over-fitting occurrence by incorporating multi-tasks learning.

Model

Neural Sequence Labeling

Bi-LSTM encoding

$$\vec{h}_t = LSTM(x_t, \vec{h}_{t-1})$$

$$\overleftarrow{h}_t = LSTM(x_t, \overleftarrow{h}_{t+1})$$

$$h_t = [\vec{h}_t; \overleftarrow{h}_t]$$

Non-Linear Transformation

$$d_t = \tanh(W_d h_t)$$

Soft-Max

$$\begin{aligned} P(y_t|d_t) &= \textit{softmax}(W_o d_t) \\ &= \frac{e^{W_{o,k} d_t}}{\sum_{\tilde{k} \in K} e^{W_{o,\tilde{k}} d_t}} \end{aligned}$$

Optimization Object

$$E = - \sum_{t=1}^T \log(P(y_t|d_t))$$

Language Modeling Objective

Bi-Prediction

$$\vec{m}_t = \tanh(\vec{W}_m \vec{h}_t)$$

$$\overleftarrow{m}_t = \tanh(\overleftarrow{W}_m \overleftarrow{h}_t)$$

$$P(w_{t+1} | \vec{m}_t) = \text{softmax}(\vec{W}_q \vec{m}_t)$$

$$P(w_{t-1} | \overleftarrow{m}_t) = \text{softmax}(\overleftarrow{W}_q \overleftarrow{m}_t)$$

Optimization Object

$$\overrightarrow{E} = - \sum_{t=1}^{T-1} \log(P(w_{t+1} | \overrightarrow{m}_t))$$

$$\overleftarrow{E} = - \sum_{t=2}^T \log(P(w_{t-1} | \overleftarrow{m}_t))$$

The Over-Whole Option Object

$$\tilde{E} = E + \gamma(\overrightarrow{E} + \overleftarrow{E})$$