

CS224N: Assignment #1

Kangwei Ling

February 27, 2017

1 Tensorflow Softmax

(a) `q1_softmax.py`

(b) `q1_softmax.py`

- (c)
- placeholder variables are nodes in the computation graph that need the data from outside to give its value. We use placeholder variables for input data and labels.
 - feed dictionaries is used to map placeholder variables to data provided when we are about to do the computation.

(d) `q1_classifier.py`

(e) `q1_classifier.py`

2 Neural Transition-Based Dependency Parsing

(a)

stack	buffer	new dependency	transition
...
[ROOT, parsed]	[this, sentence, correctly]	parsed → I	LEFT-ARC
[ROOT, parsed, this]	[sentence, correctly]		SHIFT
[ROOT, parsed, this, sentence]	[correctly]		SHIFT
[ROOT, parsed, sentence]	[correctly]	sentence → this	LEFT-ARC
[ROOT, parsed]	[correctly]	parsed → sentence	RIGHT-ARC
[ROOT, parsed, correctly]	[]		SHIFT
[ROOT, parsed]	[]	parsed → correctly	RIGHT-ARC
[ROOT]	[]	ROOT → parsed	RIGHT-ARC

- (b) Every word will be pushed onto the stack exactly 1 time. And each word will be associated with exactly 1 arc which itself is the pointee of the arc, so there are n -ARC step. Therefore, a sentence containing n words will be parsed in $2n$ steps.

(c) `q2_parser_transitions.py`

(d) `q2_parser_transitions.py`

(e) `q2_initialization.py`

(f) $\gamma = \frac{1}{p_{drop}}$. (\mathbf{d} is independent of \mathbf{h})

$$\mathbb{E}[\mathbf{h}_{drop}]_i = \gamma \mathbb{E}[\mathbf{d}]_i \cdot \mathbb{E}[\mathbf{h}]_i = \gamma p_{drop} \mathbb{E}[\mathbf{h}]_i$$

(g) (i) Using \mathbf{m} , the update of each time is a combination of current gradient and previous updates. This method helps damp the oscillation (reversing direction updates will be combined and damped), while their momentum towards the optima will be compounded, greatly boosting the converging speed.

(ii) The model parameters received smaller and infrequent gradients will get larger updates, while those with larger and frequent gradients will get lower learning rate. In this way, the learning rate is auto-adjusted and tuned for the training.

(h) UAS on dev set: 88.56

UAS on train set: 89.16

3 Recurrent Neural Networks: Language Modeling

(a) Suppose $y_c^t = 1$ (since y is one-hot), then $J^t(\theta) = -\log \hat{y}_c^{(t)}$. Thus,

$$\begin{aligned} PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) &= \frac{1}{\hat{y}_j^{(t)}} \\ &= 2^J \end{aligned}$$

From the relation between perplexity and the cross-entropy loss, it's obvious that minimizing the (arithmetic) mean cross-entropy loss will also minimize the (geometric) mean perplexity.

For a totally random model, the perplexity would be $|V|$, the cross-entropy would be $J = -\log \frac{1}{|V|} = |V|$.

(b) Let $\mathbf{z}^{(t)} = \mathbf{h}^{(t-1)}\mathbf{H} + \mathbf{e}^{(t)}\mathbf{I} + \mathbf{b}_1$, $\boldsymbol{\theta} = \mathbf{h}^{(t)}\mathbf{U} + \mathbf{b}_2$

$$\begin{aligned} \frac{\partial J^{(t)}}{\partial \mathbf{b}_2} &= \hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)} \\ \boldsymbol{\delta} &= \frac{\partial J^{(t)}}{\partial \boldsymbol{\theta}} \frac{\partial \boldsymbol{\theta}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} = (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)})\mathbf{U}^T \circ \sigma'(\mathbf{z}^{(t)}) \\ \left. \frac{\partial J^{(t)}}{\partial \mathbf{I}} \right|_{(t)} &= \boldsymbol{\delta} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{I}} = (\mathbf{e}^{(t)})^T \boldsymbol{\delta} \\ \left. \frac{\partial J^{(t)}}{\partial \mathbf{H}} \right|_{(t)} &= \boldsymbol{\delta} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{H}} = (\mathbf{h}^{(t-1)})^T \boldsymbol{\delta} \\ \frac{\partial J^{(t)}}{\partial \mathbf{L}_{\mathbf{x}^{(t)}}} &= \frac{\partial J^{(t)}}{\partial \mathbf{e}^{(t)}} = \boldsymbol{\delta} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{e}^{(t)}} = \boldsymbol{\delta} \mathbf{I}^T \\ \frac{\partial J^{(t)}}{\partial \mathbf{h}^{(t-1)}} &= \boldsymbol{\delta} \mathbf{H}^T \end{aligned}$$

(c) Use the notation in (c).

$$\begin{aligned}\frac{\partial J^{(t)}}{\partial \mathbf{H}} \Big|_{(t-1)} &= (\mathbf{h}^{(t-2)})^T (\boldsymbol{\delta}^{(t-1)} \circ \sigma'(\mathbf{z}^{(t-1)})) \\ \frac{\partial J^{(t)}}{\partial \mathbf{I}} \Big|_{(t-1)} &= (\mathbf{e}^{(t-1)})^T (\boldsymbol{\delta}^{(t-1)} \circ \sigma'(\mathbf{z}^{(t-1)})) \\ \frac{\partial J^{(t)}}{\partial \mathbf{L}_{\mathbf{x}^{(t-1)}}} &= (\boldsymbol{\delta}^{(t-1)} \circ \sigma'(\mathbf{z}^{(t-1)})) \mathbf{I}^T\end{aligned}$$

(d) forward: $O(D_h^2 + D_h(d + |V|))$
backward: $\tau \cdot O(D_h^2 + D_h(d + |V|))$

(e) bonus: hierarchical softmax?