
CS 224N: Assignment 2

PETER888@STANFORD.EDU

SUNDAY 4TH FEBRUARY, 2018

Problem 1: Tensorflow Softmax (25 points)

1.1 (a) Implement Softmax use Tensorflow (5 points, coding)

Answer:

See code: `~/q1_softmax.py`.

1.2 (b) Implement Cross-Entropy use Tensorflow (5 points, coding)

Answer:

See code: `~/q1_softmax.py`.

1.3 (c) Tensorflow Placeholder and Feed Dictionary (5 points, coding/written)

Answer:

See code: `~/q1_classifier.py`.

Explanation:

Placeholder variables and feed dictionary make it possible to feed external data (such as training data for network) into Tensorflow's computational graph.

1.4 (d) Implement Classifier (5 points, coding)

Answer:

See code: `~/q1_classifier.py`.

1.5 (e) Implement Model (5 points, coding/written)

Answer:

See code: `~/q1_classifier.py`.

Explanation:

When the model's `train_op` is called. Computation going forward through the computation graph with Placeholder's input and the init/start values in Variables. The Tensorflow automatic differentiation will start the backward propagation, it will automatically calculate the gradients and differentiations, and update the Variables back through the computation graph.

Problem 2: Neural Transition-Based Dependency Parsing (50 points + 2 bonus points)

2.1 (a) Dependency Parsing (6 points, written)

Answer:

stack	buffer	new dependency	transition
[ROOT]	[I, parsed, this, sentence, correctly]		Initial Configuration
[ROOT, I]	[parsed, this, sentence, correctly]		SHIFT
[ROOT, I, parsed]	[this, sentence, correctly]		SHIFT
[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

2.2 (b) How many steps (2 points, written)

Answer:

*2n parse steps. Because each word take exactly one shift transition from buffer to stack, take exactly one *-ARC (either LEFT-ARC or RIGHT-ARC) transition move out from stack, and every transition either add or remove word in the stack.*

2.3 (c) Parser Step (6 points, coding)

See code: ~/q2_parser_transitions.py.

2.4 (d) Parser Transitions (6 points, coding)

See code: ~/q2_parser_transitions.py.

2.5 (e) Xavier Initialization (4 points, coding)

See code: ~/q2_initialization.py.

2.6 (f) Dropout (2 points, written)

Answer:

$$\gamma = \frac{1}{1 - p_{\text{drop}}}$$

The mask vector \mathbf{d} set entries in \mathbf{h} to zero at probability p_{drop} , let's say \mathbf{h} has full expectation value $\mathbb{E}[\mathbf{h}] = 1$, so $\mathbb{E}[\mathbf{d} \circ \mathbf{h}] = (1 - p_{\text{drop}})$, this because p_{drop} of \mathbf{h} values are become zero.

$$\text{So, } 1 = \frac{\mathbb{E}[\mathbf{d} \circ \mathbf{h}]}{(1 - p_{\text{drop}})} \Rightarrow 1 = \mathbb{E}[\mathbf{h}_{\text{drop}}] = \frac{1}{(1 - p_{\text{drop}})} \mathbb{E}[\mathbf{d} \circ \mathbf{h}] \Rightarrow \gamma = \frac{1}{1 - p_{\text{drop}}}$$

2.7 (g) Adam Optimizer (4 points, written)

2.7.1 i) Momentum

Answer:

The rolling average, m , can reduce the variance in updates, by down-weighting the impact of the current minibatch's gradient. For instance, if β is 0.5, when updating θ , it will use half of past rolling average m , and half of current minibatch's gradient. Consider if the current minibatch's gradient varies too larger or too smaller than the last rolling average m , the momentum will reduce the larger or smaller part of the gradient to just half to get new m , and use the reduced new m to do update.

2.7.2 ii) Adaptive Learning Rates

Answer:

The parameters with smaller gradients (especial rolling average smaller) will get larger updates. For parameters been changing flatly, the adaptive learning rate will help to these parameter get larger updates, and helping these parameter move off plateaus. This can increase learning speed.

2.8 (h) Parser Model (20 points, coding/written)

Answer:

See code: `~/q2_parser_model.py`.

Report the best UAS

The best dev UAS: 88.11, and it's test UAS: 88.31

List of predicted labels, see file `q2_test.predicted.pkl`

2.9 (i) Bonus (2 points, coding/written)

Answer:

Implemented L2 regularization. The best dev UAS: 87.49, and it's test UAS: 87.86. There is no significant improvement by using L2 regularization. This could be the overfit not due to parameters values, need to tune other hyper-parameters such as hidden layer size, or dropout rate.

Problem 3: Recurrent Neural Networks (25 points + 1 bonus point)

3.1 (a) Perplexity (4 points, written)

3.1.1 i) Derive Perplexity (2 points)

Answer:

As $\mathbf{y}^{(t)}$ is an one-hot vector, assume the k -th element $y_k^{(t)}$ is 1

so,

$$\mathbf{J}^{(t)}(\boldsymbol{\theta}) = -\log(\hat{y}_k^{(t)}) = \log\left(\frac{1}{\hat{y}_k^{(t)}}\right) \quad (3.1)$$

and we also have

$$PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = \frac{1}{\hat{y}_k^{(t)}} \quad (3.2)$$

thus, we have

$$\mathbf{J}^{(t)}(\boldsymbol{\theta}) = \log(PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})) \quad (3.3)$$

3.1.2 ii) Equivalent (1 point)

Answer:

We know $\min\{\log(f(x))\} = \min\{f(x)\}, f(x) > 0$

then we have,

$$\min\{\mathbf{J}^{(t)}(\boldsymbol{\theta}) = \log(PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})) : \boldsymbol{\theta} > 0\} = \min\{PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})\} \forall (t \in [1..T])$$

from convex theory not hard to know minimizing the geometric mean equivalent to minimizing the arithmetic mean if the related function have same minimizing equivalent.

$$\min\left\{\left(\prod_{j=1}^T f_j(x)\right)^{\frac{1}{T}}\right\} = \min\left\{\frac{1}{T} \sum_{i=1}^T (g_i(x))\right\}, \text{ when } \min\{\mathbf{f}(x)\} = \min\{\mathbf{g}(x)\}, \text{ for } \mathbf{f}, \mathbf{g} \text{ are positive functions}$$

Finally, we can get the minimizing geometric mean perplexity equivalent to minimizing the arithmetic mean cross-entropy loss

$$\min\left\{\left(\prod_{t=1}^T PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})\right)^{\frac{1}{T}}\right\} = \min\left\{\frac{1}{T} \sum_{t=1}^T CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)})\right\}$$

3.1.3 iii) Perplexity for a single word (1 point)

Answer:

for given word ω_j ,

$$\bar{P}(x^{(t+1)} = \omega_j | x^{(t)}, \dots, x^{(1)}) = \frac{1}{|V|}$$

so the perplexity for that single word ω_j , is $|V|$

$$PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = 1/\frac{1}{|V|} = |V|$$

because, $\mathbf{J}^{(t)}(\boldsymbol{\theta}) = \log(PP^{(t)}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}))$, when $|V|=10000$

$$\mathbf{J}^{(t)}(\boldsymbol{\theta}) = \log(|V|) \approx 9.213$$

3.2 (b) Gradients on Single Point (7 points, written)

Answer:

$$\delta_1^{(t)} = \frac{\partial \mathbf{J}}{\partial \boldsymbol{\theta}^{(t)}} = -\mathbf{y} + \hat{\mathbf{y}} \quad (3.4)$$

Write sigmoid(x) as $\sigma(x)$, and it's derivative as $\sigma'(x)$

$$\begin{aligned} \delta_2^{(t)} &= \frac{\partial \mathbf{J}}{\partial \mathbf{z}^{(t)}} = \delta_1^{(t)} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} \\ &= \mathbf{U}^T \cdot \delta_1^{(t)} \circ \sigma'(\mathbf{z}^{(t)}) \\ &= \mathbf{U}^T \cdot \delta_1^{(t)} \circ \mathbf{h}^{(t)} \circ (1 - \mathbf{h}^{(t)}) \end{aligned} \quad (3.5)$$

$$\frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{U}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{U}} = \delta_1^{(t)} (\mathbf{h}^{(t)})^T \quad (3.6)$$

$$\frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{e}^{(t)}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{e}^{(t)}} = (\mathbf{W}_e|_{(t)})^T \delta_2^{(t)} \quad (3.7)$$

$$\left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)} = \left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)} = \delta_2^{(t)} (\mathbf{e}^{(t)})^T \quad (3.8)$$

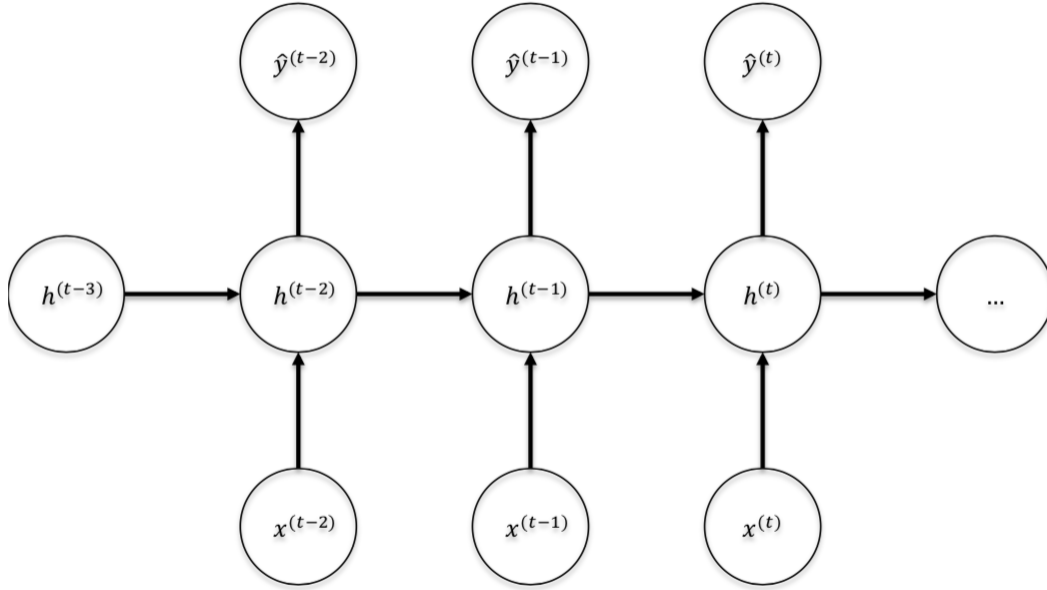
$$\left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)} = \left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)} = \delta_2^{(t)} (\mathbf{h}^{(t-1)})^T \quad (3.9)$$

$$\frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = (\mathbf{W}_h|_{(t)})^T \delta_2^{(t)} \quad (3.10)$$

3.3 (c) Gradients (7 points, written)

Answer:

Figure 3.1: Unrolled RNN



We know:

$$\mathbf{z}^{(t-1)} = \mathbf{W}_h|_{(t-1)} \mathbf{h}^{(t-2)} + \mathbf{W}_e|_{(t-1)} \mathbf{e}^{(t-1)} + \mathbf{b}_1|_{(t-1)} \quad (3.11)$$

$$\mathbf{h}^{(t-1)} = \text{sigmoid}(\mathbf{z}^{(t-1)}) = \sigma(\mathbf{z}^{(t-1)}) \quad (3.12)$$

$$\gamma^{(t-1)} = \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \quad (3.13)$$

Then,

$$\begin{aligned} \delta_2^{(t-1)} &= \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t-1)}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{h}^{(t-1)}} \frac{\partial \mathbf{h}^{(t-1)}}{\partial \mathbf{z}^{(t-1)}} \\ &= \gamma^{(t-1)} \circ \sigma'(\mathbf{z}^{(t-1)}) \\ &= \gamma^{(t-1)} \circ \mathbf{h}^{(t-1)} \circ (1 - \mathbf{h}^{(t-1)}) \end{aligned} \quad (3.14)$$

$$\frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{e}^{(t-1)}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{z}^{(t-1)}}{\partial \mathbf{e}^{(t-1)}} = (\mathbf{W}_e|_{(t-1)})^T \delta_2^{(t-1)} \quad (3.15)$$

$$\left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{W}_e} \right|_{(t-1)} = \left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{z}^{(t-1)}}{\partial \mathbf{W}_e} \right|_{(t-1)} = \delta_2^{(t-1)} (\mathbf{e}^{(t-1)})^T \quad (3.16)$$

$$\left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{W}_h} \right|_{(t-1)} = \left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t-1)}} \frac{\partial \mathbf{z}^{(t-1)}}{\partial \mathbf{W}_h} \right|_{(t-1)} = \delta_2^{(t-1)} (\mathbf{h}^{(t-2)})^T \quad (3.17)$$

3.4 (d) How Many Operations for Single Timestep (3 points, written)

Answer:

Go through each equations in part (b), those backpropagation's operations are $O(|V| \times D_h)$, $O(d \times D_h)$ and $O(D_h \times D_h)$, adding up all those operations, for a single step the operation count is $O(|V| \times D_h + d \times D_h + D_h \times D_h)$

3.5 (e) How Many Operations for Entire Sequence (3 points, written)

Answer:

Go through each backpropagation equations operation in part (c), the operation count also $O(|V| \times D_h + d \times D_h + D_h \times D_h)$. So for T words, the total operation count is $O(T \times (|V| \times D_h + d \times D_h + D_h \times D_h))$

3.6 (f) Which largest? Term RNN? (1 point, written)

Answer:

$O(|V| \times D_h)$ term is the largest. Word dictionary $|V|$ is normally large, and $D_h \ll |V|$. It happens when backpropagate from output layer to hidden layer.

3.7 (g) Bonus (1 point, written)

Answer:

If we can change the operation count from $O(T \times (|V| \times D_h + d \times D_h + D_h \times D_h))$ to $O(T \times (C \times D_h + d \times D_h + D_h \times D_h))$, by mapping words to classes, and C is the number of classes, because $C < |V|$, this will speed up the calculation. And this approach suggested by the paper "Extensions of recurrent neural network language model" (Mikolov, et al. 2013)