
CS 224N: Assignment 2

PETER888@STANFORD.EDU

SATURDAY 3RD 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:

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:

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
[<i>ROOT</i>]	[<i>I, parsed, this, sentence, correctly</i>]		Initial Configuration
[<i>ROOT, I</i>]	[<i>parsed, this, sentence, correctly</i>]		SHIFT
[<i>ROOT, I, parsed</i>]	[<i>this, sentence, correctly</i>]		SHIFT
[<i>ROOT, parsed</i>]	[<i>this, sentence, correctly</i>]	parsed → I	LEFT-ARC
[<i>ROOT, parsed, this</i>]	[<i>sentence, correctly</i>]		SHIFT
[<i>ROOT, parsed, this, sentence</i>]	[<i>correctly</i>]		SHIFT
[<i>ROOT, parsed, sentence</i>]	[<i>correctly</i>]	sentence → this	LEFT-ARC
[<i>ROOT, parsed</i>]	[<i>correctly</i>]	parsed → sentence	RIGHT-ARC
[<i>ROOT, parsed, correctly</i>]	[]		SHIFT
[<i>ROOT, parsed</i>]	[]	parsed → correctly	RIGHT-ARC
[<i>ROOT</i>]	[]	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:

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

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 $\mathbb{E}[\mathbf{d} \odot \mathbf{h}] = (1 - p_{drop})$, this because p_{drop} of \mathbf{h} values are become zero.

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

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

2.7.1 i) Momentum

Answer:

2.7.2 ii) Adaptive Learning Rates

Answer:

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

Answer:

See code: `~/q2_parser_model.py`. Report the best UAS

The dev UAS: 87.83, the test UAS: 88.09

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

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

Answer: Implemented L2 regularization. And attempt

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, suppose the k -th element $y_k^{(t)}$ is 1

so,

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

and we also have

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

thus, we have

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

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.

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

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

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

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

Answer:

for given word ω_j ,

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

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:

$$\delta_1^{(t)} = \frac{\partial \mathbf{J}}{\partial \boldsymbol{\theta}^{(t)}} = -\mathbf{y} + \hat{\mathbf{y}} \quad (3.1)$$

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

$$\begin{aligned} \delta_2^{(t)} &= \frac{\partial \mathbf{J}}{\partial \mathbf{z}^{(t)}} = \delta_1^{(t)} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{h}^{(t)}} \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{z}^{(t)}} \\ &= \mathbf{U}^T \cdot \delta_1^{(t)} \circ \sigma'(z) \end{aligned} \quad (3.2)$$

$$\frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{U}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \boldsymbol{\theta}^{(t)}} \frac{\partial \boldsymbol{\theta}^{(t)}}{\partial \mathbf{U}} = \delta_1^{(t)} (\mathbf{h}^{(t)})^T \quad (3.3)$$

$$\frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{e}^{(t)}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{e}^{(t)}} = (\mathbf{W}_e|_{(t)})^T \delta_2^{(t)} \quad (3.4)$$

$$\left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)} = \left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_e} \right|_{(t)} = \delta_2^{(t)} (\mathbf{e}^{(t)})^T \quad (3.5)$$

$$\left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)} = \left. \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{W}_h} \right|_{(t)} = \delta_2^{(t)} (\mathbf{h}^{(t-1)})^T \quad (3.6)$$

$$\frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = \frac{\partial \mathbf{J}^{(t)}}{\partial \mathbf{z}^{(t)}} \frac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{h}^{(t-1)}} = (\mathbf{W}_h|_{(t)})^T \delta_2^{(t)} \quad (3.7)$$

3.3 (c) Gradients (7 points, written)

Answer:

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

Answer:

Derivatives for the skip-gram model

$$\frac{J_{\text{skip-gram}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{v}_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c}$$

$$\frac{J_{\text{skip-gram}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{v}_j} = 0, \forall j \neq c$$

$$\frac{J_{\text{skip-gram}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{U}} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{U}}$$

Derivatives for the CBOW model

$$\frac{J_{\text{CBOW}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{v}_j} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \hat{\mathbf{v}}}, \forall (j \neq c) \in \{c-m\dots c+m\}$$

$$\frac{J_{\text{CBOW}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{v}_j} = 0, \forall (j \neq c) \notin \{c-m\dots c+m\}$$

$$\frac{J_{\text{CBOW}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{U}} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{U}}$$

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

Answer: See code: `~/q3_word2vec.py`.

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

Answer: See code: `~/q3_sgd.py`.

3.7 (g) Bonus (1 point, written)

Answer:

Explain: In the Word Vectors image, words clustered at similarity, such as the emotion words "amazing", "wonderful" and "great" are very close to each other. The word "well" a little further but still close to "amazing", the connection characters and words "the" "a" ", " etc are spread around alone.