

---

# CS 224N: Assignment 2

PETER888@STANFORD.EDU

FRIDAY 26<sup>TH</sup> JANUARY, 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:

$$\begin{aligned}\frac{d\sigma(x)}{dx} &= \frac{d}{dx} \left( \frac{1}{1 + e^{-x}} \right) \\ &= -\frac{1}{(1 + e^{-x})^2} \frac{d}{dx} (1 + e^{-x}) = \frac{-1}{(1 + e^{-x})^2} \frac{d}{dx} e^{-x} = \frac{-1}{(1 + e^{-x})^2} \frac{-1}{(e^x)^2} \frac{d}{dx} e^x \\ &= \frac{-1}{(1 + e^{-x})^2} \frac{-1}{(e^x)^2} e^x = \frac{1}{(1 + e^{-x})^2 (e^x)^2} e^x \\ &= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1 + e^{-x} - 1}{(1 + e^{-x})^2} = \frac{1}{1 + e^{-x}} - \frac{1}{(1 + e^{-x})^2} \\ &= \left( \frac{1}{1 + e^{-x}} \right) \left( 1 - \frac{1}{1 + e^{-x}} \right) \\ &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

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

**Answer:**

*With above equations,*

$$\frac{\partial CE(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -\mathbf{y} + \hat{\mathbf{y}}$$

## 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:

## 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  
List of predicted labels

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

Answer:

### **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:

##### **3.1.2 ii) Equivalent (1 point)**

Answer:

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

Answer:

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

Answer:

$$\begin{aligned}
\frac{\partial J}{\partial \mathbf{v}_c} &= \frac{\partial}{\partial \mathbf{v}_c} ((-\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c))) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))) \\
&= -\frac{\partial}{\partial \mathbf{v}_c} (\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c))) - \sum_{k=1}^K \frac{\partial}{\partial \mathbf{v}_c} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \\
&= \frac{-1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \left( \frac{\partial}{\partial \mathbf{v}_c} \sigma(\mathbf{u}_o^T \mathbf{v}_c) \right) - \sum_{k=1}^K \frac{1}{\sigma(-\mathbf{u}_k^T \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sigma(-\mathbf{u}_k^T \mathbf{v}_c) \\
&= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o - \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) (-\mathbf{u}_k) \\
&= (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{u}_o - \sum_{k=1}^K (\sigma(-\mathbf{u}_k^T \mathbf{v}_c) - 1) \mathbf{u}_k
\end{aligned}$$

$$\begin{aligned}
\frac{\partial J}{\partial \mathbf{u}_o} &= \frac{\partial}{\partial \mathbf{u}_o} ((-\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c))) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))) \\
&= \frac{-1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \left( \frac{\partial}{\partial \mathbf{u}_o} \sigma(\mathbf{u}_o^T \mathbf{v}_c) \right) - 0 \\
&= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \frac{\partial}{\partial \mathbf{u}_o} (\mathbf{u}_o^T \mathbf{v}_c) \\
&= (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{u}_o \\
\frac{\partial J}{\partial \mathbf{u}_k} &= \frac{\partial}{\partial \mathbf{u}_k} ((-\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c))) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c))) \\
&= 0 - \frac{\partial}{\partial \mathbf{u}_k} \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \\
&= -\frac{1}{\sigma(-\mathbf{u}_k^T \mathbf{v}_c)} \sigma(-\mathbf{u}_k^T \mathbf{v}_c) (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \frac{\partial}{\partial \mathbf{u}_k} (-\mathbf{u}_k^T \mathbf{v}_c) \\
&= -(\sigma(-\mathbf{u}_k^T \mathbf{v}_c) - 1) \mathbf{v}_c \\
&= (\sigma(\mathbf{u}_k^T \mathbf{v}_c) - 1) \mathbf{v}_c, \text{ for all } k \neq o
\end{aligned}$$

Negative sampling is faster than softmax-CE loss at speed up ratio  $\frac{V}{K}$ , where V is all words in dictionary count, and K is the sampling size.



### 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.*