

---

# CS 224N: Assignment 1

PETER888@STANFORD.EDU

SATURDAY 10<sup>TH</sup> FEBRUARY, 2018

---

## Problem 1: Softmax (10 points)

### 1.1 (a) Softmax Invariance to Constant (5 points)

Answer:

$\forall i,$

$$\begin{aligned}(\text{softmax}(x + c))_i &= \frac{e^{(x_i + c)}}{\sum_j e^{(x_j + c)}} = \frac{e^{x_i} \times e^c}{e^c \times \sum_j e^{x_j}} \\ &= \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(x)_i\end{aligned}$$

so,

$$\text{softmax}(x) = \text{softmax}(x + c)$$

### 1.2 (b) Softmax Coding (5 points)

Answer:

See code: `~/q1_softmax.py`.

## Problem 2: Neural Network Basics (30 points)

### 2.1 (a) Sigmoid Gradient (3 points)

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) Softmax Gradient w/ Cross Entropy Loss (3 points)

Answer:

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log(\hat{y}_i)$$

As  $\mathbf{y}$  is one-hot vector, assume  $y_k = 1$ , then other  $y_i = 0, \forall i \neq k$

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = -y_k \log(\hat{y}_k) = -\log(\hat{y}_k)$$

$$\text{As, } \hat{y}_k = \text{softmax}(\boldsymbol{\theta}) = \frac{e^{\theta_k}}{\sum_j e^{\theta_j}}$$

$$\begin{aligned} CE(\boldsymbol{\theta}) &= -\log\left(\frac{e^{\theta_k}}{\sum_j e^{\theta_j}}\right) \\ &= -\theta_k + \log\left(\sum_j e^{\theta_j}\right) \end{aligned}$$

So, for  $y_k = 1$

$$\begin{aligned} \frac{\partial CE(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} &= \frac{\partial}{\partial \boldsymbol{\theta}} (-\theta_k + \log(\sum_j e^{\theta_j})) \\ &= -\frac{\partial \theta_k}{\partial \boldsymbol{\theta}} + \frac{\partial}{\partial \boldsymbol{\theta}} \log(\sum_j e^{\theta_j}) \end{aligned} \quad (2.1)$$

And,

$$\frac{\partial \theta_k}{\partial \theta_i} = \begin{cases} 1 & , \text{if } i = k \\ 0 & , \text{if } i \neq k \end{cases}$$

So,

$$\frac{\partial \theta_k}{\partial \boldsymbol{\theta}} = \mathbf{y} \quad (2.2)$$

And,

$$\begin{aligned} \frac{\partial}{\partial \theta_i} \log(\sum_j e^{\theta_j}) &= \frac{1}{\sum_j e^{\theta_j}} \frac{\partial}{\partial \theta_i} (\sum_j e^{\theta_j}) \\ &= \frac{1}{\sum_j e^{\theta_j}} \frac{\partial}{\partial \theta_i} e^{\theta_i} \\ &= \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \\ &= \text{softmax}(\theta_i) \end{aligned} \quad (2.3)$$

With above equations,

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

## 2.3 (c) One Hidden Layer Gradient (6 points)

Answer:

To calculate  $\frac{\partial J}{\partial \mathbf{x}}$ , where  $J = CE(\mathbf{y}, \hat{\mathbf{y}})$

Let  $\boldsymbol{\theta} = \mathbf{h}\mathbf{W}_2 + \mathbf{b}_2$ , and  $\mathbf{Z1} = \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1$

and from 2(b), we know  $\frac{\partial}{\partial \boldsymbol{\theta}} CE(\mathbf{y}, \hat{\mathbf{y}}) = -\mathbf{y} + \hat{\mathbf{y}}$

we get,

$$\begin{aligned}\frac{\partial J}{\partial \boldsymbol{\theta}} &= \boldsymbol{\delta}_3 = \hat{\mathbf{y}} - \mathbf{y} \\ \frac{\partial J}{\partial \mathbf{h}} &= \boldsymbol{\delta}_2 = \boldsymbol{\delta}_3 \frac{\partial \boldsymbol{\theta}}{\partial \mathbf{h}} = \boldsymbol{\delta}_3 \mathbf{W}_2^T \\ \frac{\partial J}{\partial \mathbf{Z1}} &= \boldsymbol{\delta}_1 = \frac{\partial J}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{Z1}} = \boldsymbol{\delta}_2 \circ \sigma'(\mathbf{Z1}) \\ \frac{\partial J}{\partial \mathbf{x}} &= \boldsymbol{\delta}_1 \frac{\partial \mathbf{Z1}}{\partial \mathbf{x}} \\ &= \boldsymbol{\delta}_1 \mathbf{W}_1^T\end{aligned}$$

## 2.4 (d) No. Parameters (2 points)

Answer:

$\mathbf{W}_1, \mathbf{W}_2, b_1, b_2$  are parameters. And  $\mathbf{W1} \in \mathbb{R}^{D_x \times H}$ ,  $\mathbf{W2} \in \mathbb{R}^{H \times D_y}$ ,  $b_1 \in \mathbb{R}^H$ ,  $b_2 \in \mathbb{R}^{D_y}$

In total, parameter count is

$$(D_x \times H) + H + (H \times D_y) + D_y$$

## 2.5 (e) Sigmoid Activation Code (4 points)

See code: `~/q2_sigmoid.py`.

## 2.6 (f) Gradient Check Code (4 points)

See code: `~/q2_gradcheck.py`.

## 2.7 (g) Neural Net Code (8 points)

See code: `~/q2_neural.py`.

## Problem 3: Word2Vec (40 points + 2 bonus)

### 3.1 (a) Context Word Gradients (3 points)

Answer:

$$\text{set } z_i = \mathbf{u}_i^T \mathbf{v}_c, \text{ so } \mathbf{Z} = \mathbf{U}^T \mathbf{v}_c$$

From Problem 2.2  $\frac{\partial J}{\partial \theta} = (\hat{\mathbf{y}} - \mathbf{y})$ . Here, let  $\theta = \mathbf{Z}$ .

And also note that  $\mathbf{Z} \in \mathbb{R}^{V \times 1}$ ,  $\mathbf{v}_c \in \mathbb{R}^{D \times 1}$ ,  $\mathbf{U} \in \mathbb{R}^{D \times V}$ ,  $\hat{\mathbf{y}} - \mathbf{y} \in \mathbb{R}^{V \times 1}$

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{v}_c} &= \frac{\partial J}{\partial \mathbf{Z}} \frac{\partial \mathbf{Z}}{\partial \mathbf{v}_c} \\ &= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}) \end{aligned}$$

### 3.2 (b) Output Word Gradients (3 points)

Answer:

Note that  $\mathbf{Z} \in \mathbb{R}^{V \times 1}$ ,  $\mathbf{v}_c \in \mathbb{R}^{D \times 1}$ ,  $\mathbf{U} \in \mathbb{R}^{D \times V}$ ,  $\hat{\mathbf{y}} - \mathbf{y} \in \mathbb{R}^{V \times 1}$

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{U}} &= \frac{\partial J}{\partial \mathbf{Z}} \frac{\partial \mathbf{Z}}{\partial \mathbf{U}} \\ &= \mathbf{v}_c (\hat{\mathbf{y}} - \mathbf{y})^T \end{aligned}$$

### 3.3 (c) Repeat Gradients with Negative Sampling Loss (6 points)

Answer:

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

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

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.

### 3.4 (d) Cost Function's Gradients (8 points)

Answer:

Derivatives for the skip-gram model

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

Derivatives for the CBOW model

$$\begin{aligned}\frac{J_{\text{CBOW}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{v}_j} &= \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \hat{\mathbf{v}}}, \forall (j \neq c) \in \{c-m\dots c+m\} \\ \frac{J_{\text{CBOW}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{v}_j} &= 0, \forall (j \neq c) \notin \{c-m\dots c+m\} \\ \frac{J_{\text{CBOW}}(\text{word}_{c-m\dots c+m})}{\partial \mathbf{U}} &= \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{U}}\end{aligned}$$



### 3.5 (e) Word2Vec with SGD Code (12 points)

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

### 3.6 (f) SGD Code (4 points)

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

### 3.7 (g) Train Vectors on Stanford Sentiment Treebank (4 points)

**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.

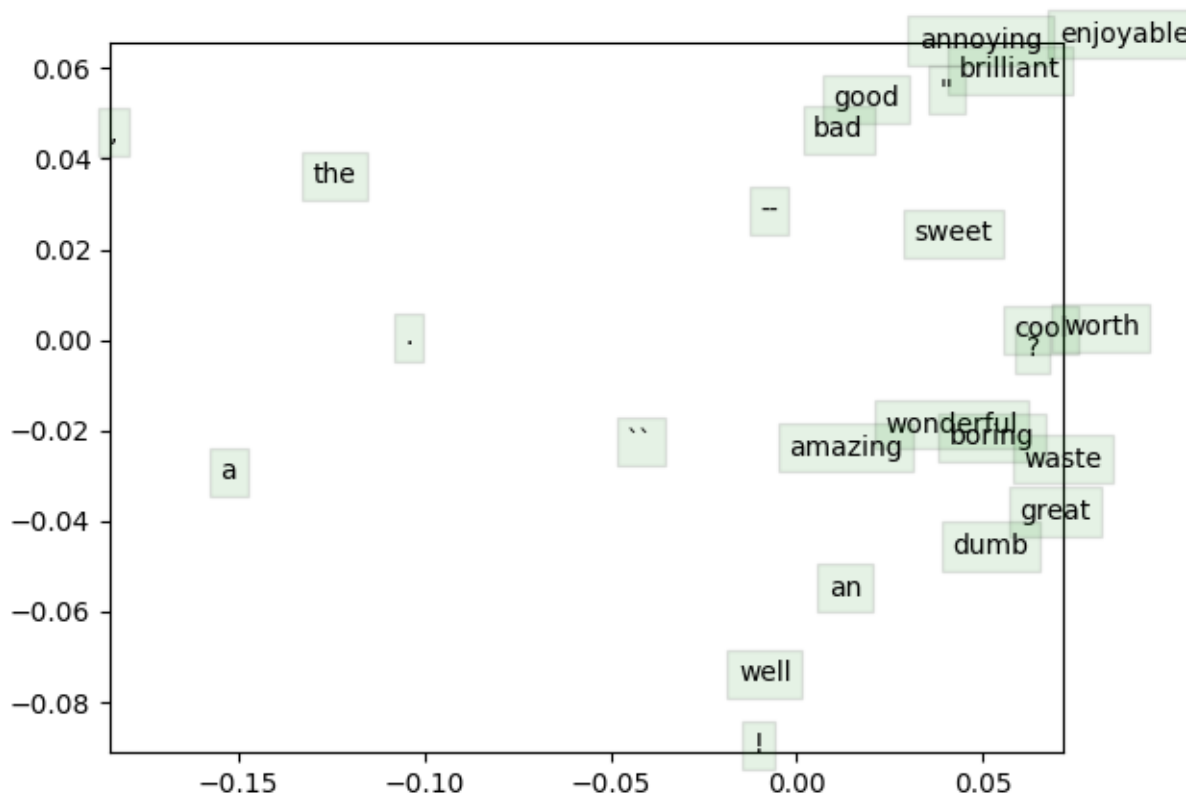


Figure 3.1: Word Vectors

### 3.8 (h) CBOW (Extra Credit: 2 points)

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

## Problem 4: Sentiment Analysis (20 points)

### 4.1 (a) Sentence Features (2 points)

**Answer:** See code: `~/q4_sentiment.py`.

### 4.2 (b) Purpose of Regularization (1 points)

**Answer:**

Purpose of regularization to constrain parameters, to avoid overfit on training dataset. And this can increase the generalizability of the model for new data.

### 4.3 (c) Hyperparameter selection (2 points)

**Answer:** See code: `~/q4_sentiment.py`.

*The chooseBestModel code*

```
### YOUR CODE HERE
best_dev = 0.0
for r in results:
    if r["dev"] > best_dev:
        bestResult = r
        best_dev = r["dev"]
### END YOUR CODE
```

## 4.4 (d) Pretrained vs Yourpoints Vectors (3 points)

Answer:

*The Best train, dev and test accuracies using "Your Vectors" are around 31.2%, 30.7%, 29.6%*

Table 4.1: Sentiment Accuracy with different Regularization rate– "Your Vectors" - on Stanford Sentiment Treebank dataset

Regularization Rate $\lambda$ Use $\frac{1}{\lambda}$	Accuracy		
	Train	Dev	Test
0.000115	<b>31.262</b>	30.699	29.593
0.007920	31.203	<b>31.153</b>	<b>29.819</b>
0.011000	31.250	30.790	29.593
0.107000	30.466	30.881	29.276
1.050000	28.956	28.247	27.059

*The Best train, dev and test accuracies using "Pretrained" are around 39.97%, 36.6%, 37.3%*

Table 4.2: Sentiment Accuracy with different Regularization rate – Pretrained GloVe

Regularization Rate $\lambda$ Use $\frac{1}{\lambda}$	Accuracy		
	Train	Dev	Test
0.000115	39.923	36.331	37.014
0.000792	<b>39.970</b>	36.240	37.104
0.011000	39.899	36.240	37.195
0.107000	39.782	<b>36.603</b>	37.195
1.050000	39.478	36.512	<b>37.330</b>

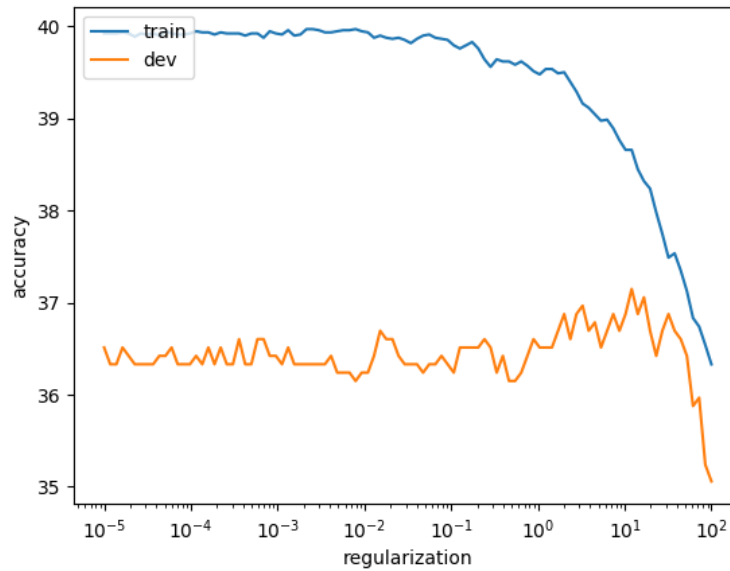
*I think there below reasons making pretrained vectors did better*

- 1. GloVe vector has higher dimensions ( $D = 50$ )*
- 2. Wikipedia dataset (billion tokens), which GloVe (pretrained) using is much larger than SST (215 thousands phrases)*
- 3. GloVe use the global static information, while Word2Vec (skip-gram) not.*

## 4.5 (e) Accuracy vs. Regularization (4 points)

Answer:

Figure 4.1: Accuracy vs. Regularization



### Explain:

In figure Accuracy vs. Regularization, the train accuracy and dev accuracy are better at  $10^{-3}$ ,  $10^{-4}$ , and  $10^{-5}$  regularization range, are around 40% and 36.3%. But the accuracy drops to 36% and 35% when regularization use  $10^1$  and  $10^2$ . And also notice that the accuracy gap are around 4% for regularization  $10^{-3}$ ,  $10^{-4}$ , and  $10^{-5}$ , which suggesting overfit on model, there is room to improve.

## 4.6 (f) Confusion Matrix (4 points)

Answer:



### Explain:

Interpret of the Confusion Matrix

The larger number in diagonal the better. We can see the model hard to tell "neutral" sentiment (value 13). And turn to put more "neutral" to "negative" (113 to 95). But model doesn't make huge mistake, like put "very negative" (-) to "very positive" (++) or vice versa.

## 4.7 (g) Error Types & Features (4 points)

Answer:

1. “and if you ’re not nearly moved to tears by a couple of scenes ,  
you ’ve got ice water in your veins .”
  - True = +; Predicted = –
  - Model not understand the idiom ”moved to tears”, so not understand whole sentence
  
2. “it ’s slow -- very , very slow . .”
  - True = --; Predicted = *netural*
  - Model need to have more context know this is a very positive sentence”.
  
3. “nothing is sacred in this gut-buster . .”
  - True = –; Predicted = +
  - Say something in converse way. That’s common error makes by word-based model.