

Answer to Question 1: Softmax

(1.a)

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

(1.b) q1_softmax.py is submitted.

Answer to Question 2: Neural Network Basics

(2.a)

$$\begin{aligned}
 \sigma'(\mathbf{x}) &= -(1 + e^{-\mathbf{x}})^{-2}(-e^{-\mathbf{x}}) \\
 &= \frac{e^{-\mathbf{x}}}{(1 + e^{-\mathbf{x}})^2} \\
 &= -\frac{1}{(1 + e^{-\mathbf{x}})^2} + \frac{1}{1 + e^{-\mathbf{x}}} \\
 &= \sigma(\mathbf{x}) - \sigma(\mathbf{x})^2
 \end{aligned}$$

(2.b)

$$\begin{aligned}
 \frac{d(CE(\mathbf{y}, \hat{\mathbf{y}}))}{d\theta} &= \frac{d(-\sum_i y_i \log(\hat{y}_i))}{d\theta} \\
 &= \frac{d(-\sum_i y_i \log(\frac{e^{\theta_i}}{\sum_j e^{\theta_j}}))}{d\theta} \\
 &= \frac{d(-\theta_k + \log \sum_i e^{\theta_i})}{d\theta} \\
 &= -\mathbf{e}_k + \frac{d(\log \sum_i e^{\theta_i})}{d\theta} \\
 &= -\mathbf{y} + \frac{1}{\sum_i e^{\theta_i}} \cdot \frac{d \sum_j e^{\theta_j}}{d\theta} \\
 &= -\mathbf{y} + \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \\
 &= -\mathbf{y} + \hat{\mathbf{y}}
 \end{aligned}$$

(2.c)

$$\begin{aligned}
\frac{\partial \mathbf{h}}{\partial \mathbf{x}} &= \frac{\partial(\text{sigmoid}(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1))}{\partial \mathbf{x}}, & \text{Let } \theta_1 &= \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1 \\
&= \frac{\partial(\text{sigmoid}(\theta_1))}{\partial \theta_1} \cdot \frac{\partial \theta_1}{\partial \mathbf{x}} \\
&= \sigma'(\theta_1) \cdot \frac{\partial(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)}{\partial \mathbf{x}} \\
&= \sigma'(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1) \cdot \mathbf{W}_1^T \\
\frac{d(CE(\mathbf{y}, \hat{\mathbf{y}}))}{d\mathbf{x}} &= \frac{d(CE(\mathbf{y}, \hat{\mathbf{y}}))}{d\theta_2} \cdot \frac{d\theta_2}{d\mathbf{x}}, & \text{Let } \theta_2 &= \mathbf{h}\mathbf{W}_2 + \mathbf{b}_2 \\
&= (\hat{\mathbf{y}} - \mathbf{y}) \cdot \frac{d\theta_2}{d\mathbf{x}} \\
&= (\hat{\mathbf{y}} - \mathbf{y}) \cdot \frac{d(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2)}{d\mathbf{x}} \\
&= (\hat{\mathbf{y}} - \mathbf{y}) \cdot (\mathbf{W}_2^T \cdot \frac{\partial \mathbf{h}}{\partial \mathbf{x}} + \mathbf{h} \frac{\partial \mathbf{W}_2}{\partial \mathbf{x}} + \frac{\partial \mathbf{b}_2}{\partial \mathbf{x}}) \\
&= (\hat{\mathbf{y}} - \mathbf{y}) \cdot (\mathbf{W}_2^T \cdot \frac{\partial \mathbf{h}}{\partial \mathbf{x}}) \\
&= (\hat{\mathbf{y}} - \mathbf{y}) \cdot \mathbf{W}_2^T \cdot \sigma'(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1) \cdot \mathbf{W}_1^T
\end{aligned}$$

(2.d)

from the inputs \mathbf{x} to the hidden layer: $(D_x + 1)H$
from the hidden layer to the outputs \mathbf{y} : $(H + 1)D_y$
 \therefore there are $((D_x + D_y + 1)H + D_y)$ parameters.

(2.e) q2_sidmoid.py is submitted.

(2.f) q2_gradcheck.py is submitted.

(2.g) q2_neural.py is submitted.

Answer to Question 3: word2vec

(3.a) $U = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_w]$, \mathbf{y} and $\hat{\mathbf{y}}$ are both column vectors

$$\begin{aligned}
 \text{Let : } f_o &= \mathbf{u}_o^T \mathbf{v}_c, f_w = \mathbf{u}_w^T \mathbf{v}_c \\
 p(\mathbf{o}|\mathbf{c}) &= \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w=1}^W \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\
 &= \frac{\exp(f_o)}{\sum_{w=1}^W \exp(f_w)} \\
 &= \text{softmax}(f_o) \\
 \frac{\partial J}{\partial \mathbf{v}_c} &= \sum_{w=1}^W \frac{\partial(-\log(\text{softmax}(f_o)))}{\partial \mathbf{v}_c} \\
 &= \sum_{w=1}^W \frac{\partial(-\log(\text{softmax}(f_o)))}{\partial f_w} \cdot \frac{\partial f_w}{\partial \mathbf{v}_c} \\
 &= \sum_{w=1}^W \delta_w \cdot \frac{\partial \mathbf{u}_w^T \mathbf{v}_c}{\partial \mathbf{v}_c} \\
 &= \sum_{w=1}^W \delta_w \mathbf{u}_w \\
 &= U^T \delta \\
 &= U^T (\hat{\mathbf{y}} - \mathbf{y})
 \end{aligned}$$

(3.b) Similar to (3.a), we can get:

$$\begin{aligned}
 \frac{\partial J}{\partial \mathbf{u}_w} &= \frac{\partial(-\log(\text{softmax}(f_o)))}{\partial \mathbf{u}_w} \\
 &= \frac{\partial(-\log(\text{softmax}(f_o)))}{\partial f_w} \cdot \frac{\partial f_w}{\partial \mathbf{u}_w} \\
 &= \delta_w \cdot \frac{\partial \mathbf{u}_w^T \mathbf{v}_c}{\partial \mathbf{u}_w} \\
 &= \delta_w \mathbf{v}_c \\
 \frac{\partial J}{\partial U} &= \mathbf{v}_c \delta^T \\
 &= \mathbf{v}_c (\hat{\mathbf{y}} - \mathbf{y})^T
 \end{aligned}$$

(3.c)

$$\begin{aligned}
\frac{\partial J_{neg-sample}(\mathbf{o}, \mathbf{v}_c, U)}{\partial \mathbf{v}_c} &= \frac{\partial(-\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)))}{\partial \mathbf{v}_c} \\
&= -\frac{\sigma'(\mathbf{u}_o^T \mathbf{v}_c)}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \mathbf{u}_o - \sum_{k=1}^K \frac{\sigma'(-\mathbf{u}_k^T \mathbf{v}_c)}{\sigma(-\mathbf{u}_k^T \mathbf{v}_c)} (-\mathbf{u}_k) \\
&= (-1 + \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o + \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{u}_k \\
\frac{\partial J_{neg-sample}(\mathbf{o}, \mathbf{v}_c, U)}{\partial \mathbf{u}_o} &= \frac{\partial(-\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)))}{\partial \mathbf{u}_o} \\
&= (-1 + \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{v}_c \\
\frac{\partial J_{neg-sample}(\mathbf{o}, \mathbf{v}_c, U)}{\partial \mathbf{u}_k} &= \frac{\partial(-\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)))}{\partial \mathbf{u}_k} \\
&= (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{v}_c
\end{aligned}$$

This cost function sums over K elements, but the softmax-CE loss function sums over W elements.

(3.d)

For skip-gram:

$$\begin{aligned}
\frac{\partial J_{skip-gram}(word_{c-m, \dots, c+m})}{\partial \mathbf{v}_k} &= \frac{\sum_{-m \leq j \leq m, m \neq 0} \partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_k} \\
&= \begin{cases} \frac{\sum_{-m \leq j \leq m, m \neq 0} \partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c}, & \text{for } k = c \\ \mathbf{0}, & \text{for } k \neq c \end{cases} \\
\frac{\partial J_{skip-gram}(word_{c-m, \dots, c+m})}{\partial \mathbf{u}_k} &= \frac{\sum_{-m \leq j \leq m, m \neq 0} \partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{u}_k}
\end{aligned}$$

For the CBOW:

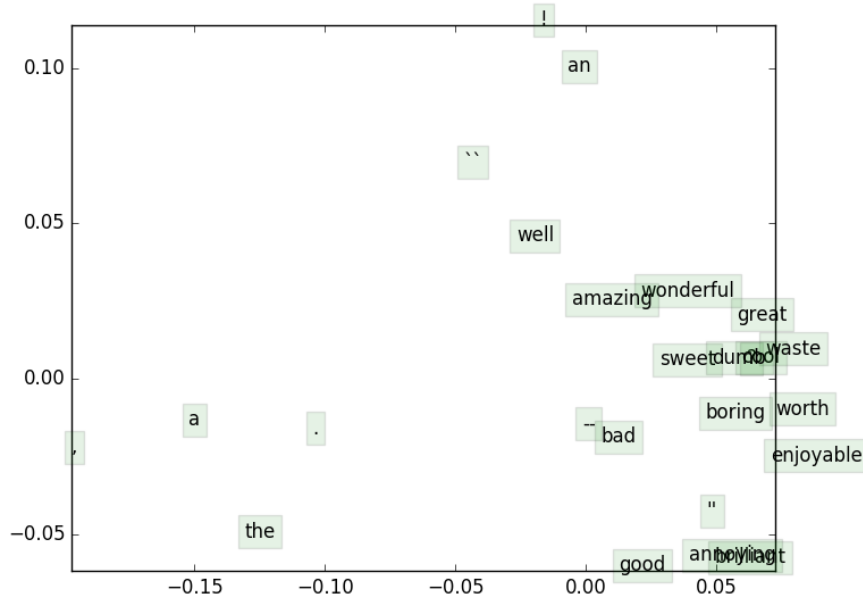
$$\begin{aligned}
\frac{\partial J_{CBOW}(word_{c-m, \dots, c+m})}{\partial \mathbf{v}_k} &= \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{v}_k} \\
&= \frac{\partial F(\mathbf{w}_c, \sum_{-m \leq j \leq m, j \neq 0} \mathbf{v}_{c+j})}{\partial \mathbf{v}_k} \\
&= \begin{cases} \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{v}_k}, & \text{for } k = c - m, \dots, c + m \\ \mathbf{0}, & \text{otherwise} \end{cases}
\end{aligned}$$

$$\frac{\partial J_{CBOW}(word_{c-m, \dots, c+m})}{\partial \mathbf{u}_k} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{u}_k}$$

(3.e) q3_word2vec.py is submitted.

(3.f) q3_sgd.py is submitted.

(3.g) q3_run.py is submitted.



Those words with similar meanings (like "wonderful" and "amazing" which are both adjectives) are located in a region. Punctuation and words like "the", "a", "an" which don't have any specific meanings are located separately, scattered all around the plane.

Answer to Question 4: Sentiment Analysis

(4.a) q4_sentiment.py is submitted.

(4.b) In order to avoid overfitting. Otherwise, the parameters can be extremely large to get high accuracy, which results in overfitting.

(4.c) q4_sentiment.py is submitted.

Best regularization value: 5.00E-04

Test accuracy (%): 37.013575

Code for chooseBestModel:

```
def chooseBestModel(results):
    bestResult = None
    ### YOUR CODE HERE
    #print results
    maxDev = 0
    for r in results:
        if maxDev <= r['dev']:
            maxDev = r['dev']
            bestResult = r
    ### END YOUR CODE
    return bestResult
```

(4.d)

Results:

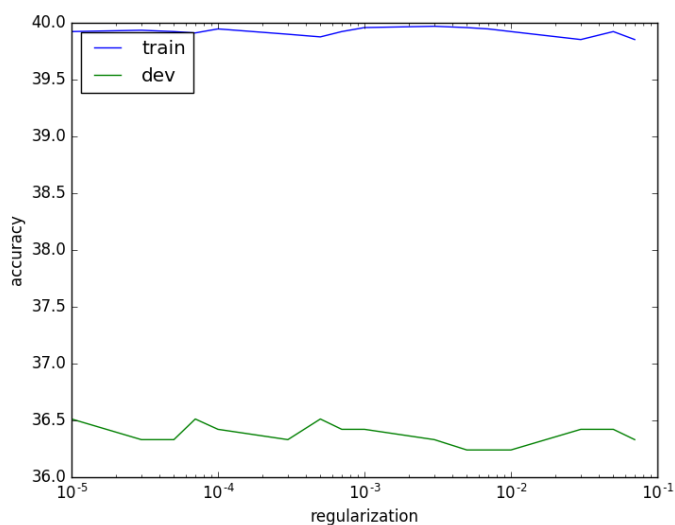
=== Recap ===				=== Recap ===			
Reg	Train	Dev	Test	Reg	Train	Dev	Test
1.00E-05	39.923	36.512	37.014	1.00E-05	31.051	32.607	30.452
3.00E-05	39.934	36.331	36.968	3.00E-05	30.993	32.698	30.271
5.00E-05	39.923	36.331	36.923	5.00E-05	31.016	32.516	30.498
7.00E-05	39.911	36.512	37.014	7.00E-05	31.086	32.516	30.362
1.00E-04	39.946	36.421	36.968	1.00E-04	31.039	32.516	30.452
3.00E-04	39.899	36.331	36.968	3.00E-04	31.133	32.698	30.452
5.00E-04	39.876	36.512	37.014	5.00E-04	31.145	32.607	30.362
7.00E-04	39.923	36.421	37.014	7.00E-04	31.098	32.516	30.407
1.00E-03	39.958	36.421	36.968	1.00E-03	31.121	32.607	30.362
3.00E-03	39.970	36.331	37.014	3.00E-03	31.098	32.243	30.362
5.00E-03	39.958	36.240	37.104	5.00E-03	30.911	32.425	30.452
7.00E-03	39.946	36.240	37.149	7.00E-03	30.887	32.516	30.317
1.00E-02	39.923	36.240	37.195	1.00E-02	30.946	32.334	29.910
3.00E-02	39.853	36.421	37.376	3.00E-02	30.536	31.608	30.045
5.00E-02	39.923	36.421	37.511	5.00E-02	30.337	31.880	30.136
7.00E-02	39.853	36.331	37.195	7.00E-02	30.372	32.243	30.045
Best regularization value: 5.00E-04				Best regularization value: 3.00E-04			
Test accuracy (%): 37.013575				Test accuracy (%): 30.452489			

(i) For the pretrained vectors, the best train accuracy is 39.978%, the best dev accuracy is 36.512%, and the best test accuracy is 37.511%.

(ii) For my word vectors, the best train accuracy is 31.145%, the best dev accuracy is 32.698%, and the best test accuracy is 30.498%.

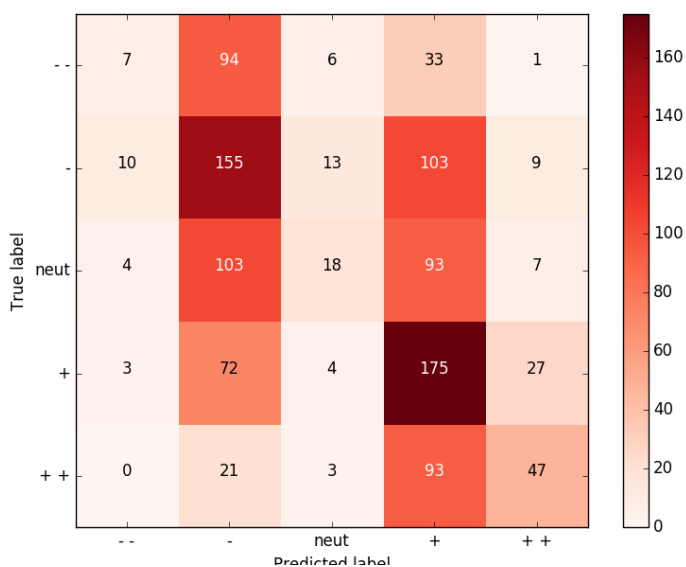
(iii) Reasons: The Wikipedia database is larger in size and it's more reliable. It's likely that the training set and the test set share similar features. GloVe also contributes to good performance.

(4.e)



Train accuracy is around 40% and development accuracy is around 36.5%. According to the plot, increasing or decreasing the regularization will not affect the accuracy much.

(4.f)



According to the diagonal elements, the prediction accuracy is relatively high for the positive(+) class and the negative(-) class. According to the off-diagonal elements, the model tends to label the most vectors as positive(+) or negative(-) instead of very negative(--), neutral or very positive(++).

(4.g)

True: 4 Predicted: 1

Text: an effectively creepy , fear-inducing -lrb- not fear-reducing -rrb- film from japanese director hideo nakata , who takes the superstitious curse on chain letters and actually applies it .

Reason: The text uses several negative adjectives to describe the director's talents, which can be confusing for the classifier. The classifier should be able to classify a polysemy according to the text.

True: 4 Predicted: 1

Text: the movie is n't just hilarious : it 's witty and inventive , too , and in hindsight , it is n't even all that dumb .

Reason: The text misspelled 'not' as 'n't'. The classifier should be able to find and correct minor misspellings.

True: 1 Predicted: 4

Text: directed in a paint-by-numbers manner .

Reason: The text uses a compound. The classifier should be able to separate and classify compounds.