

# Assignment 1

## 1 Softmax

(a) Proof:

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

(b) See file q1\_softmax.py

## 2 Neural Network Basic

(a) We have:

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

(b) Because  $\mathbf{y}$  is a one-hot label, we can assume that all elements in it is zero, except that the  $k$ th element is one. Therefore, by using the chain rule, we have:

$$\begin{aligned}\frac{\partial CE}{\partial \theta_i} &= \frac{\partial CE}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \theta_i} \\ &= -\frac{1}{\hat{\mathbf{y}}_k} \frac{\partial \hat{\mathbf{y}}}{\partial \theta_i}\end{aligned}$$

For  $i = k$ , we have:

$$\frac{\partial \hat{\mathbf{y}}}{\partial \theta_k} = -\frac{e^{x_k} \sum_{j \neq k} e^{x_j}}{(\sum_j e^{x_j})^2}$$

For  $i \neq k$ , we have:

$$\frac{\partial \hat{\mathbf{y}}}{\partial \theta_i} = \frac{e^{x_k} e^{x_i}}{(\sum_j e^{x_j})^2}$$

Therefore, we have:

$$\frac{\partial CE}{\partial \theta} = \hat{\mathbf{y}} - \mathbf{y}$$

(c) By using the result from the above question, we know that:

$$\frac{\partial CE}{\partial \mathbf{h}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^T$$

Let  $\mathbf{O}_1 = \mathbf{x} \mathbf{W}_1 + \mathbf{b}_1$ , and by using the result of question (a), we have:

$$\frac{\partial CE}{\partial \mathbf{O}_1} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^T * \mathbf{h}(1 - \mathbf{h})$$

Therefore, we have:

$$\frac{\partial CE}{\partial \mathbf{x}} = \left( (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^T * \mathbf{h}(1 - \mathbf{h}) \right) \mathbf{W}_1^T$$

(d) The dimension of  $\mathbf{W}_1$  is  $D_x \times H$  and the dimension of  $\mathbf{W}_2$  is  $H \times D_y$ . Therefore, the total number of parameters are  $D_x \times H + H \times D_y + H + D_y$

(e) See file q2\_sigmoid.py

(f) See file q2\_gradcheck.py

(g) See file q2\_neural.py

### 3 word2vec

(a) Let  $\mathbf{U} = [\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_v]$  where each column represents an "output" word vector for some word in the vocabulary. The gradient of loss with respect to  $\mathbf{v}_c$  is given by:

$$\begin{aligned} \frac{\partial CE}{\partial \mathbf{v}_c} &= -\frac{1}{\hat{y}_o} \frac{\partial \hat{y}_o}{\partial \mathbf{v}_c} \\ &= \frac{1}{\hat{y}_o} \frac{\left( \sum_{w=1}^V ((\boldsymbol{\mu}_w - \boldsymbol{\mu}_o) \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c)) \right) \exp(\boldsymbol{\mu}_o^T \mathbf{v}_c)}{\left( \sum_{w=1}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c) \right)^2} \\ &= \frac{\sum_{w=1}^V ((\boldsymbol{\mu}_w - \boldsymbol{\mu}_o) \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c))}{\sum_{w=1}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c)} \end{aligned}$$

(b) For  $\boldsymbol{\mu}_k$  where  $k \neq o$ , we have:

$$\begin{aligned} \frac{\partial CE}{\partial \boldsymbol{\mu}_k} &= \frac{1}{\hat{y}_o} \frac{\exp(\boldsymbol{\mu}_o^T \mathbf{v}_c) \exp(\boldsymbol{\mu}_k^T \mathbf{v}_c) \mathbf{v}_c}{\left( \sum_{w=1}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c) \right)^2} \\ &= \frac{\exp(\boldsymbol{\mu}_k^T \mathbf{v}_c) \mathbf{v}_c}{\sum_{w=1}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c)} \end{aligned}$$

For  $\boldsymbol{\mu}_o$ , we have:

$$\begin{aligned} \frac{\partial CE}{\partial \boldsymbol{\mu}_o} &= -\frac{1}{\hat{y}_o} \frac{\mathbf{v}_c \exp(\boldsymbol{\mu}_o^T \mathbf{v}_c) \sum_{w=1, w \neq o}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c)}{\left( \sum_{w=1}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c) \right)^2} \\ &= -\frac{\mathbf{v}_c \sum_{w=1, w \neq o}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c)}{\sum_{w=1}^V \exp(\boldsymbol{\mu}_w^T \mathbf{v}_c)} \end{aligned}$$

(c) The gradient of negative sampling loss with respect to  $\mathbf{v}_c$  is given by:

$$\frac{\partial J_{NS}}{\partial \mathbf{v}_c} = (\sigma(\boldsymbol{\mu}_o^T \mathbf{v}_c) - 1) \boldsymbol{\mu}_o - \sum_{k=1}^K \left( (\sigma(-\boldsymbol{\mu}_k^T \mathbf{v}_c) - 1) \boldsymbol{\mu}_k \right)$$



## 4 Sentiment Analysis

(a) See file q4\_sentiment.py

(b) The reason why we want to introduce regularization when doing classification is to prevent overfitting.

(c) See file q4\_sentiment.py. My code for *chooseBestModel* is:

```

1 current_best_dev_acc = 0
2 for result in results:
3     if result['dev'] > current_best_dev_acc:
4         bestResult = result
5         current_best_dev_acc = result['dev']

```

(d) The best results are:

	Train	Dev	Test
Our vectors	31.004%	32.698%	30.362%
Pretrained	39.934%	36.512%	37.014%

The pretrained vectors are better because:

1. These vectors are trained on a much larger data set.
2. These vectors are trained for a much longer time.
3. GloVe tends to work better than Skip-Gram and CBOW.

(e) The plot is:

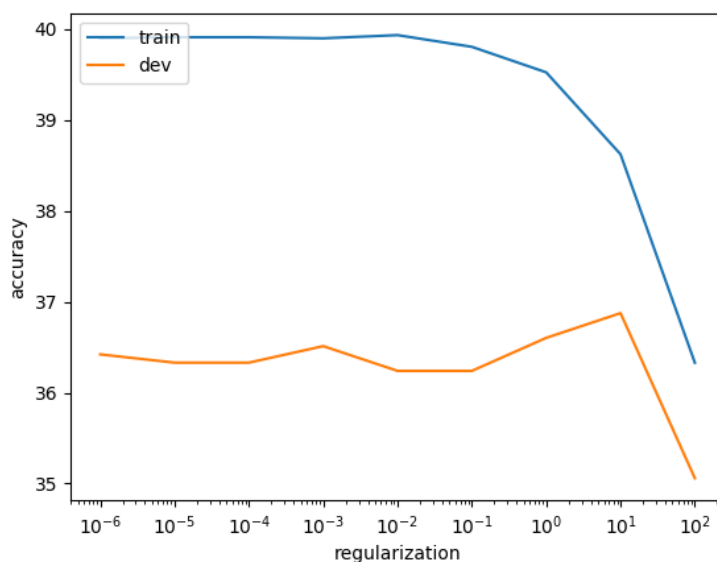


图 2: Reg Value with respect to Dev and Train Acc

We can see from above that when regularization value is small, there is a huge gap between training accuracy and dev accuracy, meaning that we have encountered overfitting. As the regularization value

increase, the gap becomes smaller, meaning that we have greatly reduced the overfitting. However, as the regularization value continues to increase, both the training accuracy and dev accuracy begin to fall, meaning that we may encounter underfitting.

(f) The plot is:

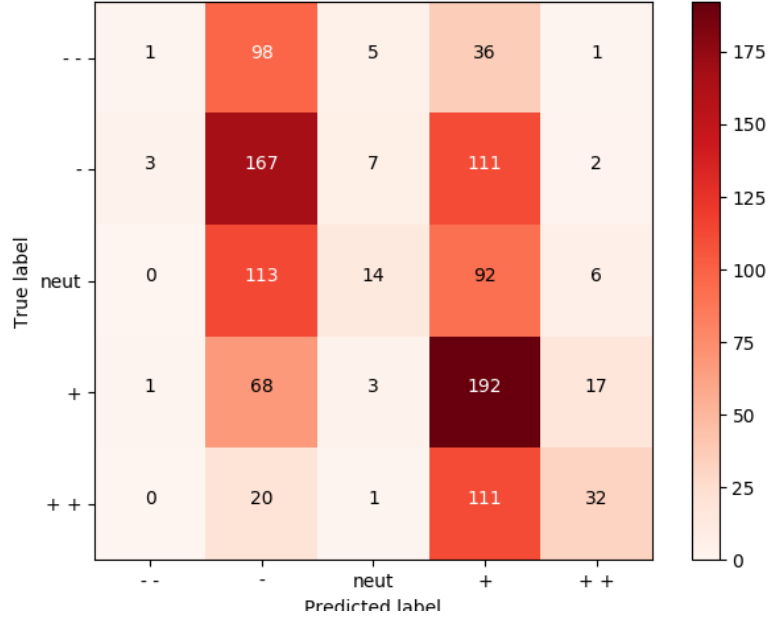


图 3: Confusion Matrix on the Dev Set (with Pretrained Vectors)

It is easy to see from above that our model tends to predict a sentence more positive than it should be. Also, our model using the pretrained vectors works quite well that it usually only makes small mistakes (like predicting a positive sentence to be very positive) rather than big mistakes (like predicting a very negative sentence to be very positive).

(g) For the following examples:

1. Comment: like mike is a winner for kids, and no doubt a winner for lil bow wow, who can now add movies to the list of things he does well. Its ground truth is very positive (4), however, we predict it as negative (1). We will need feature "list of things he does well".
2. Comment: this riveting world war ii moral suspense story deals with the shadow side of american culture: racial prejudice in its ugly and diverse forms. Its ground truth is negative (1), however we predict it as positive (3).
3. Comment: a painfully funny ode to bad behavior. Its ground truth is positive (3), while we predict it as vert negative (0). This maybe is because our model sees words like "painfully" and "bad". The feature we need to classify it right is "ode".