

Assignment #1 Solutions

(1) SOFTMAX

(a)

For each dimension i

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

Q.E.D

(b) Please check code.

(2) NEURAL NETWORK BASICS

(c)

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

(b)

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

$$\hat{y} = \text{softmax}(o)$$

$$\Rightarrow \hat{y}_i = \frac{e^{\theta_i}}{\sum_j e^{\theta_j}}$$

$$\Rightarrow \frac{\partial \hat{y}_i}{\partial \theta_i} = \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \left(1 - \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \right) = \hat{y}_i(1 - \hat{y}_i) \quad - \textcircled{1}$$

Also

$$\frac{\partial \hat{y}_i}{\partial \theta_j} = -\frac{e^{\theta_i} e^{\theta_j}}{(\sum_j e^{\theta_j})^2} = -\hat{y}_i \hat{y}_j \quad - \textcircled{2}$$

If k is the true class,

$$\begin{aligned} CE(y, \hat{y}) &= -\sum_i y_i \log(\hat{y}_i) \\ &= -\log(\hat{y}_k) \end{aligned}$$

$$\therefore \frac{\partial CE}{\partial \theta_i} = -\frac{1}{\hat{y}_k} \frac{\partial \hat{y}_k}{\partial \theta_i} = \begin{cases} \hat{y}_i - 1 & \text{if } i = k \\ \hat{y}_i & \text{otherwise} \end{cases} \quad [\text{Using } \textcircled{1} \text{ & } \textcircled{2}]$$

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

$$(c) \quad h = \text{sigmoid}(xW_1 + b_1) \quad \hat{y} = \text{softmax}(hW_2 + b_2)$$

$$\text{Let say } z_1 = xW_1 + b_1 \quad \& \quad z_2 = hW_2 + b_2$$

$$\therefore h = \text{sigmoid}(z_1)$$

$$\hat{y} = \text{softmax}(z_2)$$

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

$$\frac{\partial CE}{\partial z_2} = \hat{y} - y \quad (\text{proved in part (b)})$$

Using backpropagation / chain rule.

$$\begin{aligned} \frac{\partial CE}{\partial h} &= \frac{\partial CE}{\partial z_2} \cdot \cancel{\frac{\partial z_2}{\partial h}} \\ &= (\hat{y} - y) \cdot W_2^T \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial CE}{\partial z_1} &= \frac{\partial CE}{\partial h} \odot \sigma'(z_1) \\ &= [(\hat{y} - y) \cdot W_2^T] \odot \sigma'(z_1) \end{aligned}$$

Following the same approach

$$\frac{\partial CE}{\partial x} = \frac{\partial CE}{\partial z_1} \cdot \frac{\partial z_1}{\partial x}$$

$$\boxed{\frac{\partial CE}{\partial x} = [((\hat{y} - y) \cdot W_2^T) \odot \sigma'(z_1)] \cdot W_1^T}$$

$$(d) \quad \text{No. of Weights in layer 1} = D_x \cdot H$$

$$\text{biases in layer 1} = H$$

$$\text{No. of weights in layer 2} = D_y \cdot H$$

$$\text{biases in layer 2} = D_y$$

$$\boxed{\text{Total number of parameters} = (D_x + 1)H + D_y(H + 1)}$$

Check code for parts (e), (f) & (g)

(3) Word 2 vec

$$(a) \hat{y}_o = p(o|c) = \frac{\exp(u_o^\top v_c)}{\sum_{w=1}^W \exp(u_w^\top v_c)}$$

$$J_{\text{softmax-CE}}(o, v_c, U) = CE(y, \hat{y})$$

Assuming $\theta = [u_1^\top v_c \ u_2^\top v_c \ \dots \ u_w^\top v_c]$ - row vector of shape $1 \times W$

$$\frac{\partial J}{\partial \theta} = \hat{y} - y \quad (\text{from part(b) Q2}) \quad \text{assume } \hat{y} - y \text{ is also a row vector of shape } 1 \times W$$

Also $\frac{\partial \theta}{\partial v_c}$ is a Jacobian and it can be written in matrix form as follows

$$\frac{\partial \theta}{\partial v_c} = [u_1, u_2, u_3, \dots, u_w] = U \quad (\text{dimension } W \times d)$$

$$\text{And } \frac{\partial J}{\partial \theta} = \left[\frac{\partial J}{\partial \theta_1}, \dots, \frac{\partial J}{\partial \theta_W} \right]_{1 \times W} \quad \text{and} \quad \frac{\partial \theta}{\partial v_c} = \begin{bmatrix} \frac{\partial \theta_1}{\partial v_{c1}}, \frac{\partial \theta_1}{\partial v_{c2}}, \dots, \frac{\partial \theta_1}{\partial v_{cd}} \\ \vdots \\ \frac{\partial \theta_W}{\partial v_{c1}}, \frac{\partial \theta_W}{\partial v_{c2}}, \dots, \frac{\partial \theta_W}{\partial v_{cd}} \end{bmatrix}_{W \times d}$$

$$\text{Since } \frac{\partial J}{\partial v_{ci}} = \sum_i \frac{\partial J}{\partial \theta_i} \frac{\partial \theta_i}{\partial v_{ci}}$$

$\therefore \frac{\partial J}{\partial v_c}$ can be written as matrix multiplication of U & $\hat{y} - y$ after arranging dimensions.

$$\boxed{\frac{\partial J}{\partial v_c} = U \underbrace{(\hat{y} - y)}_{1 \times d \text{ row vector}} U^\top}$$

Please note that ~~this~~ multiplication order could be changed depending upon how matrix dimensions are defined

(b) Assuming $\theta = [u_1^T v_c \ u_2^T v_c \ \dots \ u_w^T v_c]$ - row vector of shape $1 \times w$

$$\frac{\partial J}{\partial \theta} = (\hat{y} - y) \quad (\text{row vector of shape } 1 \times w)$$

$$\frac{\partial \theta}{\partial u_i} = \begin{bmatrix} & \stackrel{\text{row vector}}{0} \\ & \vdots \\ & v_c^T \\ & \vdots \\ & 0 \end{bmatrix} = v_i \quad \Rightarrow \text{is a Jacobian matrix of shape } w \times d$$

From part (a)

$$\begin{aligned} \frac{\partial J}{\partial u_i} &= (\hat{y} - y) v_i \\ &= [v_{c1}(\hat{y}_1 - y_1) \ v_{c2}(\hat{y}_1 - y_1) \ \dots \ v_{cd}(\hat{y}_1 - y_1)] \\ &\qquad\qquad\qquad \text{row vector } d \times 1 \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial U} &= \begin{bmatrix} v_{c1}(\hat{y}_1 - y_1) & v_{c2}(\hat{y}_1 - y_1) & \dots & v_{cd}(\hat{y}_1 - y_1) \\ \vdots & & & \\ v_{c1}(\hat{y}_w - y_w) & v_{c2}(\hat{y}_w - y_w) & \dots & v_{cd}(\hat{y}_w - y_w) \end{bmatrix} - \text{Shape } w \times d \\ &= (\hat{y} - y)^T v_c \quad (\text{assuming } v_c \text{ is a row vector } 1 \times d) \end{aligned}$$

$$\boxed{\frac{\partial J}{\partial U} = (\hat{y} - y)^T v_c} \quad \text{where } (\hat{y} - y) \text{ & } v_c \text{ are row vectors}$$

The order can be different depending upon how the vector dimensions are defined.

(c)

$$J_{\text{neg-sample}}(o, v_c, v) = -\log(\sigma(u_o^T v_c)) - \sum_{k=1}^K \log(\sigma(-u_k^T v_c))$$

We know that $\sigma' = \sigma(1 - \sigma)$

Also $\frac{\partial \bar{a}^T x}{\partial x} = a$ & ~~$\frac{\partial a^T x}{\partial x}$~~ $= a$

$$\begin{aligned} \therefore \frac{\partial J}{\partial v_c} &= -\frac{1}{\sigma(u_o^T v_c)} \times \sigma(u_o^T v_c) (1 - \sigma(u_o^T v_c)) \cdot u_o \\ &\quad - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \times \sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c)) (-u_k) \end{aligned}$$

$$\frac{\partial J}{\partial v_c} = (\sigma(u_o^T v_c) - 1) u_o - \sum_{k=1}^K (\sigma(-u_k^T v_c) - 1) u_k$$

Similarly, following the same calculations:

$$\frac{\partial J}{\partial u_o} = (\sigma(u_o^T v_c) - 1) v_c$$

$$\frac{\partial J}{\partial u_k} = -(\sigma(-u_k^T v_c) - 1) v_c \quad \text{for } k = 1, \dots, K$$

For the original softmax-CE loss, we have to calculate W sums of $\exp(u_w^T v_c)$ so it can be very expensive operation depending on the size of vocabulary whereas for negative sampling loss only K sums need to be calculated [Each summation ~~of~~ item contains a dot product]. Therefore speed up ratio is $k : W$.

(d)

$$J_{\text{skip-gram}}(\text{word } c-m, \dots, c+m) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} F(w_{t+j}, v_c)$$

We already know how to calculate $\frac{\partial F}{\partial v_c}$, $\frac{\partial F}{\partial U}$ ($\frac{\partial F}{\partial u_w} v_c$)

Therefore, for skip-gram model, gradients for the context window will be given as follows

$$\frac{\partial J_{\text{skip-gram}}(\text{word } c-m, \dots, c+m)}{\partial v_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial F(w_{t+j}, v_c)}{\partial v_c}$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word } c-m, \dots, c+m)}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial F(w_{t+j}, v_c)}{\partial U}$$

for $v_j \neq v_t$ $\frac{\partial J}{\partial v_j}$ will be zero since v_j is not part of J function

and for CBOW,

$$\frac{\partial J_{\text{CBOW}}(\text{word } c-m, \dots, c+m)}{\partial U} = \frac{\partial F(w_t, \hat{v})}{\partial U} \quad \text{where } \hat{v} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} v_{t+j}$$

$$\frac{\partial J_{\text{CBOW}}(\text{word } c-m, \dots, c+m)}{\partial v_j} = \frac{\partial F(w_t, \hat{v})}{\partial v_j}$$

for $j \in \{c-m, \dots, c-1, c+1, \dots, c+m\}$

& for $j \notin \{c-m, \dots, c-1, c+1, \dots, c+m\}$

$$\frac{\partial J}{\partial v_j} = 0$$

~~$\frac{\partial \hat{v}}{\partial v_j}$~~ $\frac{\partial F(w_t, \hat{v})}{\partial \hat{v}}$

~~$d \times d$ matrix with all columns zero~~

~~diagonal matrix with diagonal entries 1~~

$\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$

↓ Proof.

$$\begin{aligned} \frac{\partial J_{\text{CBOW}}}{\partial v_j} &= \frac{\partial F(w_t, \hat{v})}{\partial \hat{v}} \frac{\partial \hat{v}}{\partial v_j} \\ &\underbrace{\frac{\partial F(w_t, \hat{v})}{\partial \hat{v}}}_{d \times d} \underbrace{\frac{\partial \hat{v}}{\partial v_j}}_{d \times d} \\ &= \frac{\partial F(w_t, \hat{v})}{\partial U} \end{aligned}$$

diagonal matrix with diagonal entries as 1

$$\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

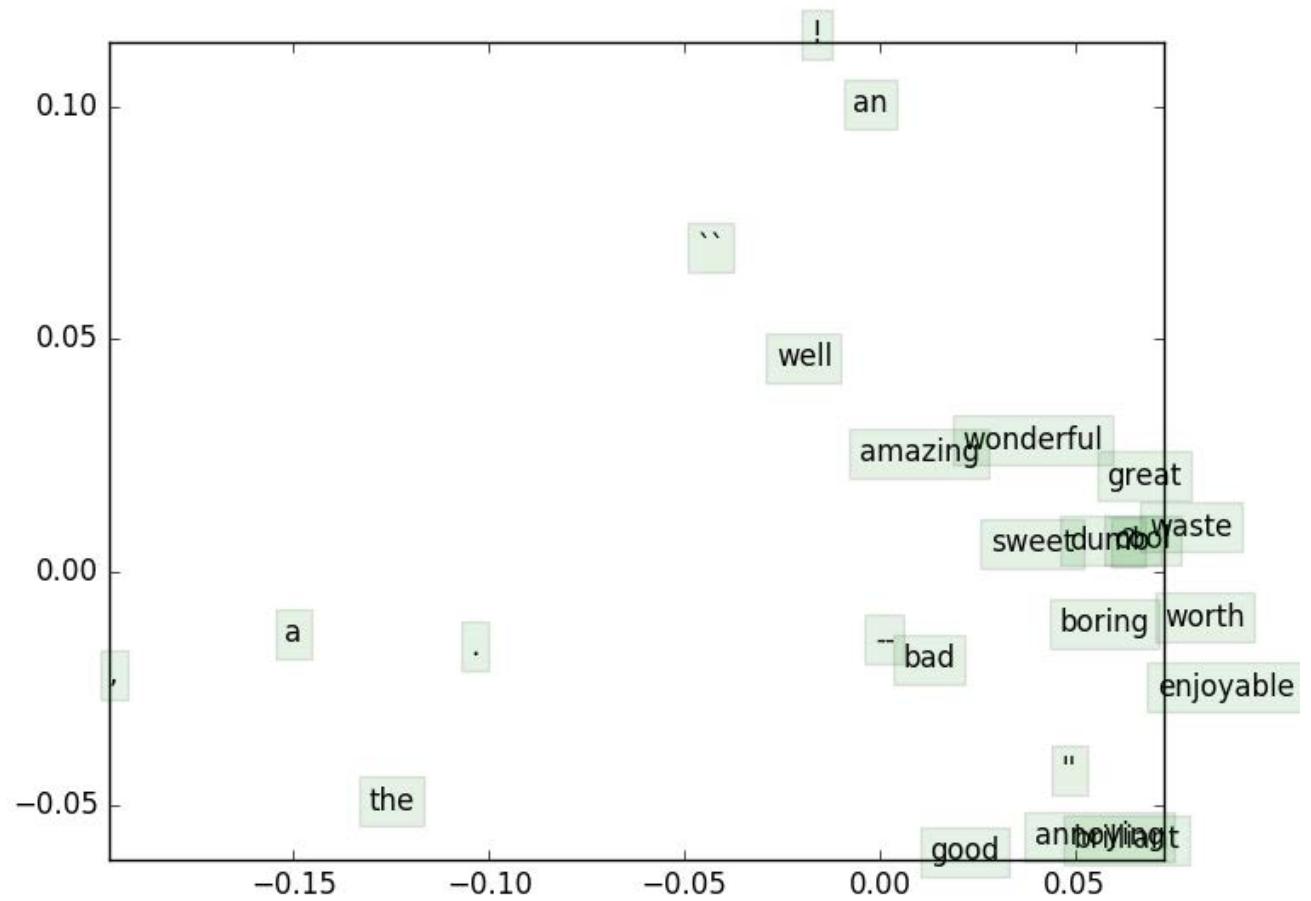
(e) Please check code

(f) Please check code

(g) Please check figure

It can be inferred from the plot that the words that occur in similar contextual meanings are located in the vicinity of each other (e.g. "amazing" & "wonderful" & "well" also "boring" & "enjoyable")

(h) Please check code.



(4) Sentiment Analysis

(a) Please check code.

(b) regularization is used to solve overfitting problem in order for the model to better generalize over unseen Examples.

(c) Please check code.

(Code attached for choose best Model)

```
def chooseBestModel(results):
    """Choose the best model based on dev set performance.

    Arguments:
    results -- A list of python dictionaries of the following format:
    {
        "reg": regularization,
        "clf": classifier,
        "train": trainAccuracy,
        "dev": devAccuracy,
        "test": testAccuracy
    }

    Each dictionary represents the performance of one model.

    Returns:
    Your chosen result dictionary.
    """
    bestResult = None

    ### YOUR CODE HERE
    best_dev_acc = 0.0
    for model in results:
        if model['dev'] > best_dev_acc:
            bestResult = model
            best_dev_acc = model['dev']
    ### END YOUR CODE
    return bestResult
```

(d)

Best train accuracy = ~~38.692~~ 39.934 %.

Best dev accuracy = 36.785 %.

Best test accuracy = 37.1945 %.

Best train accuracy = 31.168 %.

Best dev accuracy = 32.698 %.

Best test accuracy = 36.317 %.

} pretrained

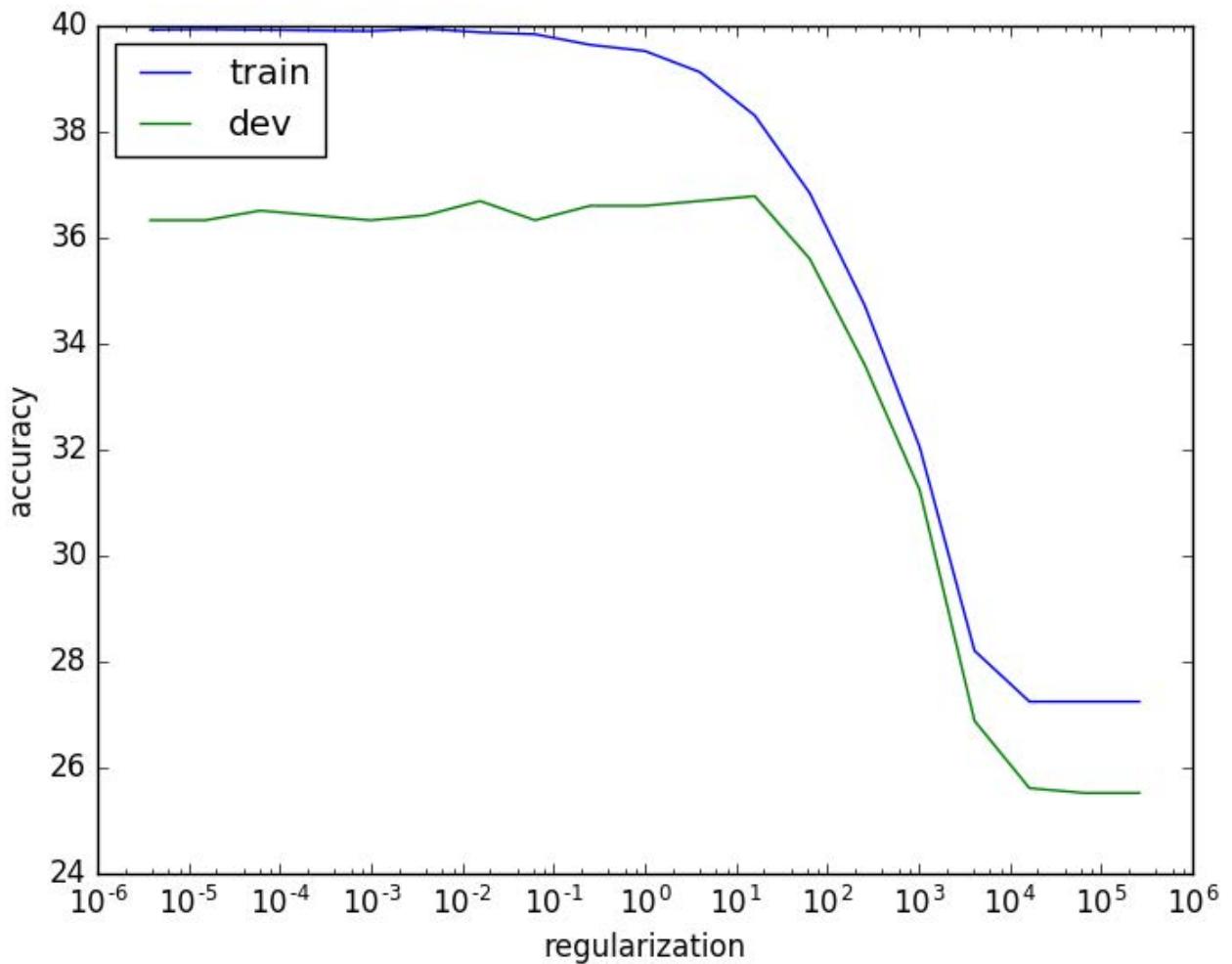
} "glove vectors"

Possible reasons why pretrained gives better performance:

- 1) pretrained glove vectors were probably trained on larger corpus
- 2) Glove vectors probably has more optimized & higher dimensions compared to word2vec.
- 3) Glove vectors have a better objective function (loss) compared to word2vec as demonstrated by authors.

(e) Please check the plot.

It can be seen from the plot that regularization parameter is increased training error gets closer and closer to dev set accuracy so in a way it better generalizes on unseen examples. Also if the regularization is increased too much, then ~~the~~ both training and dev accuracy falls.

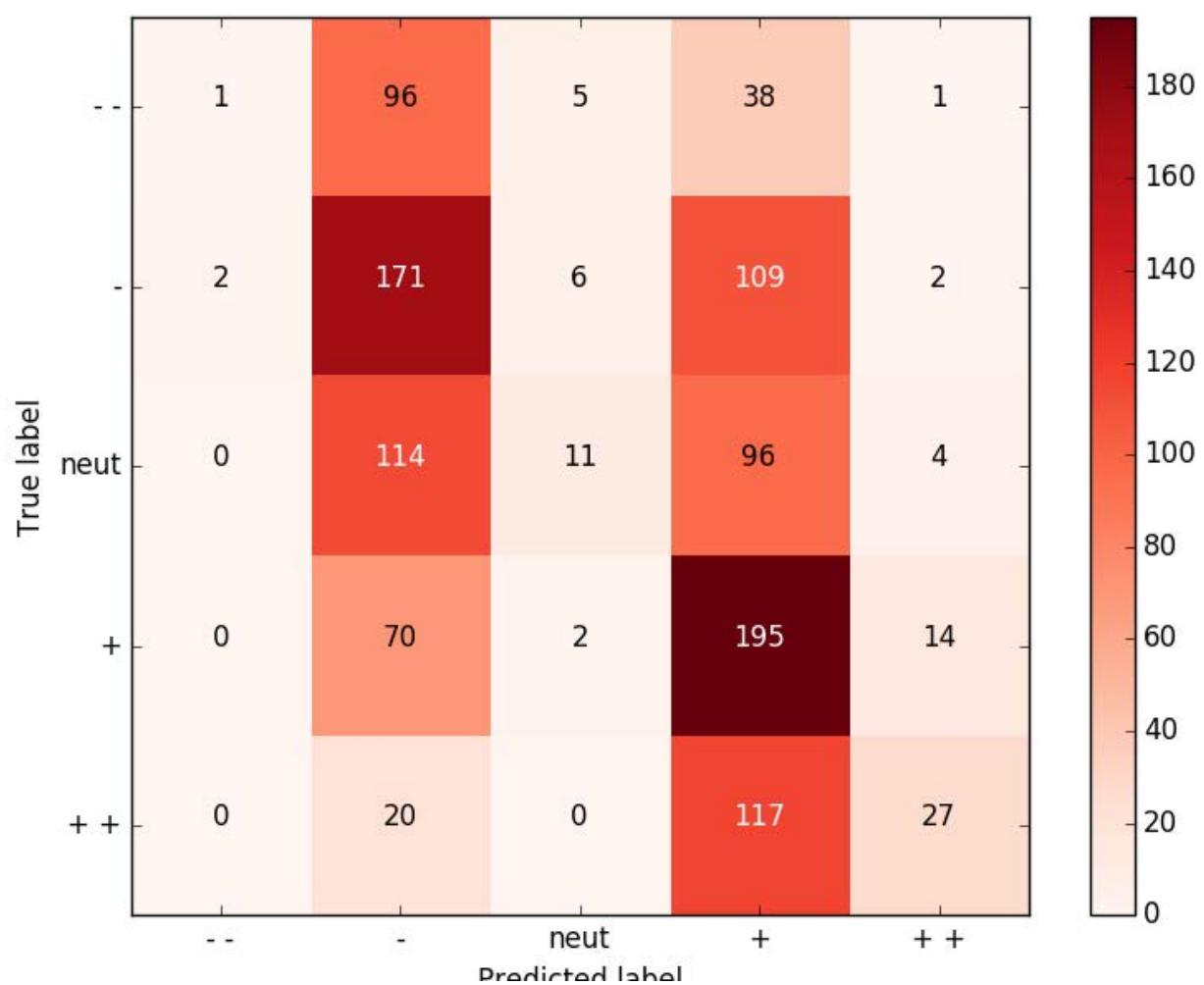


(f) Please check the plot for confusion matrix

From the confusion matrix, it can be interpreted that while in general accuracy is quite low, majority of '-' & '--' and majority of '+' & '++' are classified within their group (i.e. '-' & '--' are predicted to be either '-' & '--' and vice versa for '+' & '++')

Also neutral ~~set~~ sentences are mostly predicted to be '+' or '-' possibly due to mixture of words and no clear pattern.

One clear pattern is that very negative or very positive examples are rarely mixed with each other.



(g)

Example 1:

- ultimately feels empty and unsatisfying, like swallowing a communion wafer without the wine

True: 0

Predicted: 3

In this sentence there are positive words such as "ultimately" & "communion" etc & also negative words such as "unsatisfying" and slightly negative words such as "empty" so just averaging word vectors doesn't really clarify it as a "negative" sentence.

Example 2:

falls neatly into the category of good, stupid fun -

True: 3

Predicted: 1

Probably the word "stupid" averages out the positiveness of the ~~stuck~~ sentence.

Example 3:

the humor isn't as sharp, the effects not as innovative, nor the story as imaginative as in the original

True: 0

Predicted: 3

There are lots of positive words such as "sharp", "innovative" & "imaginative" which probably increases the positiveness of the sentence.

Averaging doesn't seem to handle negative well & also word order is lost. Perhaps instead of averaging over words, we can average over multiple context windows.