# CS 224N: Assignment #1

writer: 纪焘

## 1、Softmax

**(a)**

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

**(b)**

```python
if len(x.shape) > 1:
    # Matrix
    c = np.max(x, axis=1).reshape(x.shape[0], 1)
    x = np.exp(x - c)
    norm = np.sum(x, axis=1).reshape(x.shape[0], 1)
    x = x / norm
else:
    # Vector
    c = np.max(x)
    x = np.exp(x - c)
    x = x / x.sum()
```

## 2、Neural Network Basics

**(a)**

$$\nabla_x \sigma(x) = \nabla_x \left( \frac{1}{1 + e^{-x}} \right)$$
$$= \nabla_x \left[ (1 + e^{-x})^{-1} \right]$$
$$= -(1 + e^{-x})^{-2} \cdot \nabla_x (1 + e^{-x})$$
$$= -(1 + e^{-x})^{-2} \cdot e^{-x} \cdot \nabla_x -x$$
$$= \frac{1 + e^{-x} - 1}{(1 + e^{-x})^2}$$
$$= \frac{1}{(1 + e^{-x})} \cdot -\frac{1}{(1 + e^{-x})^2}$$
$$= \sigma(x) \cdot (1 - \sigma(x))$$

**(b)**

*Assume that only the $k$-th dimension of $y$ is one.*

$$\nabla_\theta CE(y, \hat{y}) = \nabla_\theta \left[ -\sum_i y_i \log(\hat{y}_i) \right]$$
$$= \nabla_\theta \left[ - y_k \log(\hat{y}_k) \right]$$
$$= \nabla_\theta \left[ - \log\left( \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} \right) \right]$$
$$= \nabla_\theta \left[ - \log(e^{\theta_k}) + \log\left( \sum_j e^{\theta_j} \right) \right]$$
$$= \nabla_\theta \left[ \log\left( \sum_j e^{\theta_j} \right) - \theta_k \right]$$

$$if \quad t = k:$$
$$= \nabla_{\theta_t} \left[ \log\left( \sum_j e^{\theta_j} \right) \right] - \nabla_{\theta_t} \theta_k$$
$$= \frac{\nabla_{\theta_t} \left( \sum_j e^{\theta_j} \right)}{\sum_j e^{\theta_j}} - 1$$
$$= \frac{e^{\theta_t}}{\sum_j e^{\theta_j}} - 1$$
$$= \hat{y}_t - 1$$

$$if \quad t \neq k:$$
$$= \nabla_{\theta_t} \left[ \log\left( \sum_j e^{\theta_j} \right) \right] - \nabla_{\theta_t} \theta_k$$
$$= \frac{\nabla_{\theta_t} \left( \sum_j e^{\theta_j} \right)}{\sum_j e^{\theta_j}}$$
$$= \frac{e^{\theta_t}}{\sum_j e^{\theta_j}}$$
$$= \hat{y}_t$$
$$\therefore \quad \nabla_\theta CE(y, \hat{y}) = \hat{y} - y$$

**(c)**

*symbol definition:*

$$a = xW_1 + b_1$$

$$h = \sigma(xW_1 + b_1) = \sigma(a)$$

$$z = hW_2 + b_2$$

$$\hat{y} = softmax(hW_2 + b_2) = softmax(z)$$

$$\nabla_x CE(y, \hat{y}) = \frac{\partial J}{\partial z} CE(y, \hat{y}) \cdot \frac{\partial z}{\partial h}(hW_2 + b_2) \cdot \frac{\partial h}{\partial a}\sigma(a) \cdot \frac{\partial a}{\partial x}(xW_1 + b_1)$$
$$= (\hat{y} - y) \cdot W_2^\top \cdot \sigma(a) \cdot (1 - \sigma(a)) \cdot W_1^\top$$

## (d)

**Input $\rightarrow$ Hidden**:

$$D_x * H + H$$

**Hidden $\rightarrow$ Output**:

$$H * D_y + D_y$$

**Total**:

$$(D_x + 1) * H + (H + 1) * D_y$$

## (e)

```python
def sigmoid(x):
    s = 1. / (1. + np.exp(-x))
    return s



def sigmoid_grad(s):
    ds = s * (1. - s)
    return ds
```

## (f)

```
def gradcheck_naive(f, x):
    ### ...
    x[ix] += h
    random.setstate(rndstate)
    fx1, _ = f(x)
    x[ix] -= 2 * h
    random.setstate(rndstate)
    fx2, _ = f(x)
    numgrad = (fx1-fx2) / (2.0*h)
    x[ix] += h
    ### ...
```

**(g)**

```python
def forward_backward_prop(data, labels, params, dimensions):

    ### YOUR CODE HERE: forward propagation
    M = data.shape[0]
    # (M, H)
    a = np.dot(data, W1) + b1
    hiddens = sigmoid(a)
    # (M, Dy)
    z = np.dot(hiddens, W2) + b2
    outputs = softmax(z)

    ### END YOUR CODE

    cost = -1 * labels * np.log(outputs)
    cost = cost.sum() / M

    ### YOUR CODE HERE: backward propagation

    # (M, Dy)
    gradZs = outputs - labels
    # (M, H, Dx)
    gradW2 = np.array([np.dot(hiddens[i].reshape(1, H).T,
gradZs[i].reshape(1, Dy)) for i in xrange(M)])
    # (H, Dx)
    gradW2 = gradW2.sum(axis=0) * (1.0/M)
    # (1, Dx)
    gradb2 = (gradZs.sum(axis=0) * (1.0/M)).reshape(1, Dy)
    # (M, H)
    gradAs = np.array([np.dot(gradZs[i].reshape(1, Dy),
W2.T)*sigmoid_grad(hiddens[i]) for i in xrange(M)])
    # (M, Dx, H)
    gradW1 = np.array([np.dot(data[i].reshape(1, Dx).T, gradAs[i].reshape(1,
H)) for i in xrange(M)])
    # (Dx, H)
    gradW1 = gradW1.sum(axis=0) * (1.0/M)
    # (1, H)
    gradb1 = gradAs.sum(axis=0) * (1.0/M)

    ### END YOUR CODE
```

# 3、word2vec

**(a)**

$$\nabla_{v_c} CE(y, \hat{y}) = \nabla_{v_c}\left[-\sum_i y_i \log(\hat{y}_i)\right]$$

$$= \nabla_{v_c}\left[-y_o \log(\hat{y}_o)\right]$$

$$= \nabla_{v_c}\left[-\log\left(\frac{\exp(u_o^\top v_c)}{\sum_{w=1}^W \exp(u_w^\top v_c)}\right)\right]$$

$$= \nabla_{v_c}\left[-\log\left(\exp(u_o^\top v_c)\right) + \log\left(\sum_{w=1}^W \exp(u_w^\top v_c)\right)\right]$$

$$= \nabla_{v_c}\left(-u_o^\top v_c\right) + \nabla_{v_c}\left[\log\left(\sum_{w=1}^W \exp(u_w^\top v_c)\right)\right]$$

$$= -u_o + \frac{1}{\sum_{w=1}^W \exp(u_w^\top v_c)} \cdot \nabla_{v_c}\left[\sum_{w=1}^W \exp(u_w^\top v_c)\right]$$

$$= -u_o + \frac{1}{\sum_{w=1}^W \exp(u_w^\top v_c)} \cdot \sum_{j=1}^W \exp(u_j^\top v_c) \cdot \nabla_{v_c}(u_j^\top v_c))$$

$$= -u_o + \sum_{j=1}^W \frac{\exp(u_j^\top v_c)}{\sum_{w=1}^W \exp(u_w^\top v_c)} \cdot u_j$$

$$= -u_o + \sum_{j=1}^W \hat{y}_j \cdot u_j$$

**(b)**

$$\nabla_{u_w} CE(y, \hat{y}) = \nabla_{u_w}\left(-u_o^\top v_c\right) + \nabla_{u_w}\left[\log\left(\sum_{w=1}^W \exp(u_w^\top v_c)\right)\right]$$

$if \quad w = o:$

$$= -v_c + \frac{1}{\sum_{w=1}^W \exp(u_w^\top v_c)} \cdot \nabla_{u_w}\left[\sum_{w=1}^W \exp(u_w^\top v_c)\right]$$

$$= -v_c + \frac{1}{\sum_{w=1}^W \exp(u_w^\top v_c)} \cdot \sum_{j=1}^W \exp(u_j^\top v_c) \cdot \nabla_{u_j}(u_j^\top v_c)$$

$$= -v_c + \sum_{j=1}^W \frac{\exp(u_j^\top v_c)}{\sum_{w=1}^W \exp(u_w^\top v_c)} \cdot v_c$$

$$= -v_c + \hat{y}_j \cdot v_c$$

$$= (\hat{y}_j - 1) \cdot v_c$$

$else \quad w \neq o:$

$$= \frac{1}{\sum_{w=1}^W \exp(u_w^\top v_c)} \cdot \nabla_{u_w}\left[\sum_{w=1}^W \exp(u_w^\top v_c)\right]$$

$$= \hat{y}_j \cdot v_c$$

**(c)**

$$\nabla_{v_c} J_{neg-sample}(o, v_c, U) = \nabla_{v_c} \left[ -\log\left(\sigma(u_o^\top v_c)\right) - \sum_{k=1}^{K} \log\left(\sigma(-u_k^\top v_c)\right) \right]$$

$$= -\nabla_{v_c} \left[ \log\left(\sigma(u_o^\top v_c)\right) \right] - \nabla_{v_c} \left[ \sum_{k=1}^{K} \log\left(\sigma(-u_k^\top v_c)\right) \right]$$

$$= -\frac{1}{\sigma(u_o^\top v_c)} \cdot \nabla_{v_c} \left[\sigma(u_o^\top v_c)\right] - \sum_{k=1}^{K} \frac{1}{\sigma(-u_k^\top v_c)} \cdot \nabla_{v_c} \left[\sigma(-u_k^\top v_c)\right]$$

$$= -\frac{1}{\sigma(u_o^\top v_c)} \cdot \sigma(u_o^\top v_c) \cdot \left(1 - \sigma(u_o^\top v_c)\right) \cdot \nabla_{v_c}(u_o^\top v_c)$$

$$- \sum_{k=1}^{K} \frac{1}{\sigma(-u_k^\top v_c)} \cdot \sigma(-u_k^\top v_c) \cdot \left(1 - \sigma(-u_k^\top v_c)\right) \cdot \nabla_{v_c}(-u_k^\top v_c)$$

$$= -\left(1 - \sigma(u_o^\top v_c)\right) \cdot u_o - \sum_{k=1}^{K} \left(1 - \sigma(-u_k^\top v_c)\right) \cdot -u_k$$

$$= \left(\sigma(u_o^\top v_c) - 1\right) \cdot u_o - \sum_{k=1}^{K} \left(\sigma(-u_k^\top v_c) - 1\right) \cdot u_k$$

$$\nabla_{u_w} J_{neg-sample}(o, v_c, U) = \nabla_{u_w} \left[ -\log\left(\sigma(u_o^\top v_c)\right) - \sum_{k=1}^{K} \log\left(\sigma(-u_k^\top v_c)\right) \right]$$

$$= -\nabla_{u_w} \left[ \log\left(\sigma(u_o^\top v_c)\right) \right] - \nabla_{u_w} \left[ \sum_{k=1}^{K} \log\left(\sigma(-u_k^\top v_c)\right) \right]$$

$$if \quad w = o:$$

$$= -\frac{1}{\sigma(u_o^\top v_c)} \cdot \nabla_{u_w} \left[\sigma(u_o^\top v_c)\right]$$

$$= -\frac{1}{\sigma(u_o^\top v_c)} \cdot \sigma(u_o^\top v_c) \cdot \left(1 - \sigma(u_o^\top v_c)\right) \cdot \nabla_{u_o}(u_o^\top v_c)$$

$$= -\left(1 - \sigma(u_o^\top v_c)\right) \cdot v_c$$

$$= \left(\sigma(u_o^\top v_c) - 1\right) \cdot v_c$$

$$if \quad w \in \{1, \cdots, K\}:$$

$$= -\frac{1}{\sigma(-u_k^\top v_c)} \cdot \nabla_{u_w} \left[\sigma(-u_k^\top v_c)\right]$$

$$= -\frac{1}{\sigma(-u_k^\top v_c)} \cdot \sigma(-u_k^\top v_c) \cdot \left(1 - \sigma(-u_k^\top v_c)\right) \cdot \nabla_{u_w}(-u_k^\top v_c)$$

$$= -\left(1 - \sigma(-u_k^\top v_c)\right) \cdot -v_c$$

$$= \left(1 - \sigma(-u_k^\top v_c)\right) \cdot v_c$$

$$else:$$

$$= None$$

$$\frac{\nabla_{u_w} J_{softmax-CE}(o, v_c, U)}{\nabla_{u_w} J_{neg-sample}(o, v_c, U)} = \frac{W}{K}$$

**(d)**

$$\frac{\partial J_{skip-gram}\left(w_{c-m\cdots c+m}\right)}{\partial U} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\boldsymbol{w}_{c+j}, \boldsymbol{v}_c)}{\partial U}$$

$if \quad i = c:$

$$\frac{\partial J_{skip-gram}\left(w_{c-m\cdots c+m}\right)}{\partial \boldsymbol{v}_i} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\boldsymbol{w}_{c+j}, \boldsymbol{v}_c)}{\partial \boldsymbol{v}_i}$$

$else \quad i \neq c:$

$$\frac{\partial J_{skip-gram}\left(w_{c-m\cdots c+m}\right)}{\partial \boldsymbol{v}_i} = 0$$

$$\frac{\partial J_{CBOW}\left(w_{c-m\cdots c+m}\right)}{\partial U} = \frac{\partial F(\boldsymbol{w}_c, \hat{\boldsymbol{v}})}{\partial U}$$

$if \quad i \in \{c_{window}\}:$

$$\frac{\partial J_{CBOW}\left(w_{c-m\cdots c+m}\right)}{\partial \boldsymbol{v}_i} = \frac{\partial F(\boldsymbol{w}_c, \hat{\boldsymbol{v}})}{\partial \boldsymbol{v}_i}$$

$else \quad i \notin \{c_{window}\}:$

$$\frac{\partial J_{CBOW}\left(w_{c-m\cdots c+m}\right)}{\partial \boldsymbol{v}_i} = 0$$

## (e)

```
def normalizeRows(x):
    ### YOUR CODE HERE
    x = x / np.sqrt(np.sum(x**2, axis=1)).reshape(x.shape[0], 1)
    ### END YOUR CODE


def softmaxCostAndGradient(predicted, target, outputVectors, dataset):
    ### YOUR CODE HERE
    v_c = predicted
    o = target
    u_o = outputVectors[target]
    y_ = softmax(v_c.dot(outputVectors.T))
    cost = -np.log(y_[o])
    y_[o] -= 1
    gradPred = (y_.reshape(1, y_.shape[0]).dot(outputVectors)).flatten()
    grad = y_.reshape(y_.shape[0], 1) * v_c.reshape(1, v_c.shape[0])
    ### END YOUR CODE


def negSamplingCostAndGradient(predicted, target, outputVectors, dataset,
                               K=10):
    ### YOUR CODE HERE
    v_c = predicted
    o = target
    u_o = outputVectors[target]
    c_o = sigmoid(v_c.dot(u_o))
```

```
        cost = -np.log(c_o)
        gradPred = (c_o-1) * u_o

        grad = np.zeros(outputVectors.shape)
        grad[o] = (c_o-1) * v_c

        for k in indices[1:]:
            u_k = outputVectors[k]
            c_k = sigmoid(-v_c.dot(u_k))
            cost -= np.log(c_k)
            gradPred -= (c_k-1) * u_k
            grad[k] += (1-c_k) * v_c
        ### END YOUR CODE


def skipgram(currentWord, C, contextWords, tokens, inputVectors,
outputVectors,
             dataset, word2vecCostAndGradient=softmaxCostAndGradient):
    ### YOUR CODE HERE
    c = tokens[currentWord]
    v_c = inputVectors[c]
    for wordo in contextWords:
        o = tokens[wordo]
        costo, gradv_c, gradopv = word2vecCostAndGradient(v_c, o,
outputVectors, dataset)
        cost += costo
        gradIn[c] += gradv_c
        gradOut +=  gradopv
    ### END YOUR CODE
```
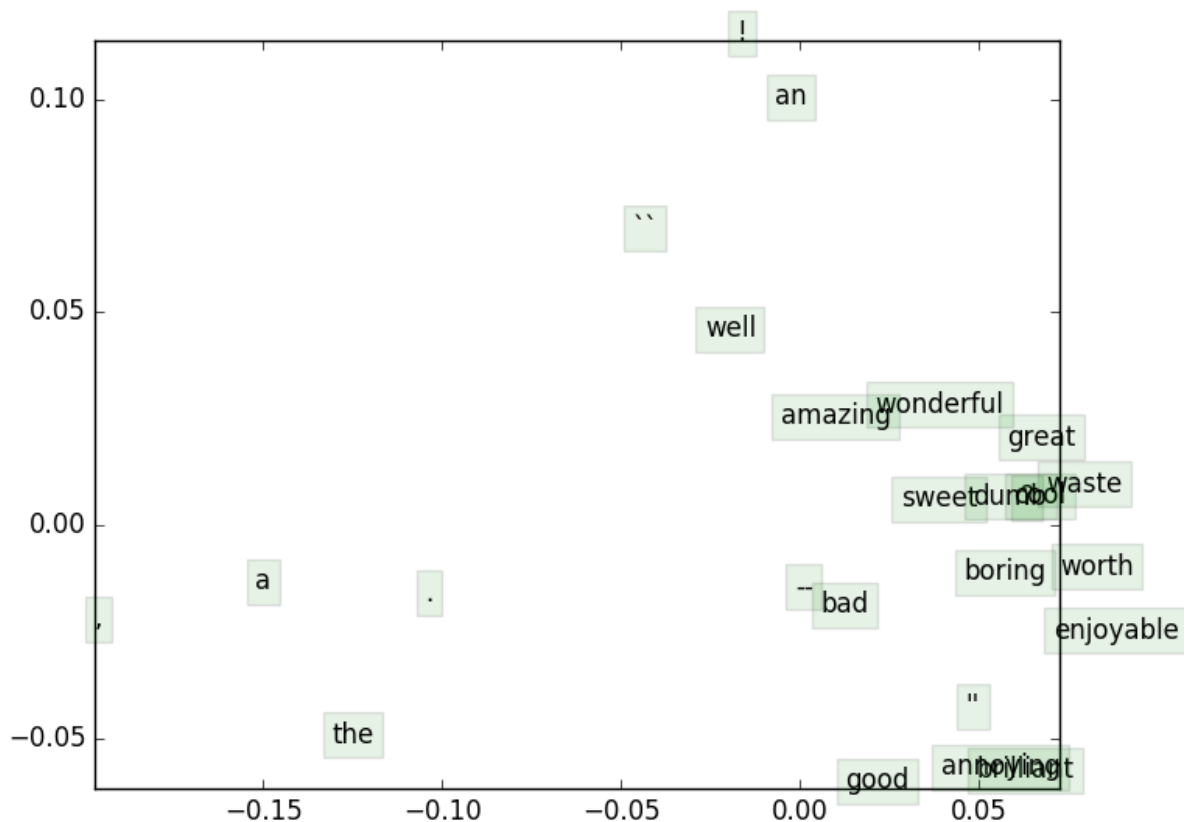
**(f)**

```
def sgd(f, x0, step, iterations, postprocessing=None, useSaved=False,
        PRINT_EVERY=10):
    ### YOUR CODE HERE
    cost, grad = f(x)
    x -= step * grad
    x = postprocessing(x)
    ### END YOUR CODE
```

**(g)**

**(h)**

```python
def cbow(currentWord, C, contextWords, tokens, inputVectors, outputVectors,
         dataset, word2vecCostAndGradient=softmaxCostAndGradient):
    ### YOUR CODE HERE
    v_c = np.zeros(inputVectors.shape[1])
    for word in contextWords:
        v_c += inputVectors[tokens[word]]
    o = tokens[currentWord]
    cost, grad, gradOut = word2vecCostAndGradient(v_c, o, outputVectors,
dataset)
    for word in contextWords:
        gradIn[tokens[word]] += grad
    ### END YOUR CODE
```

# 4、Sentiment Analysis

**(a)**

```
def getSentenceFeatures(tokens, wordVectors, sentence):
    ### YOUR CODE HERE
    for word in sentence:
        sentVector += wordVectors[tokens[word]]
    sentVector /= len(sentence)
    ### END YOUR CODE
```
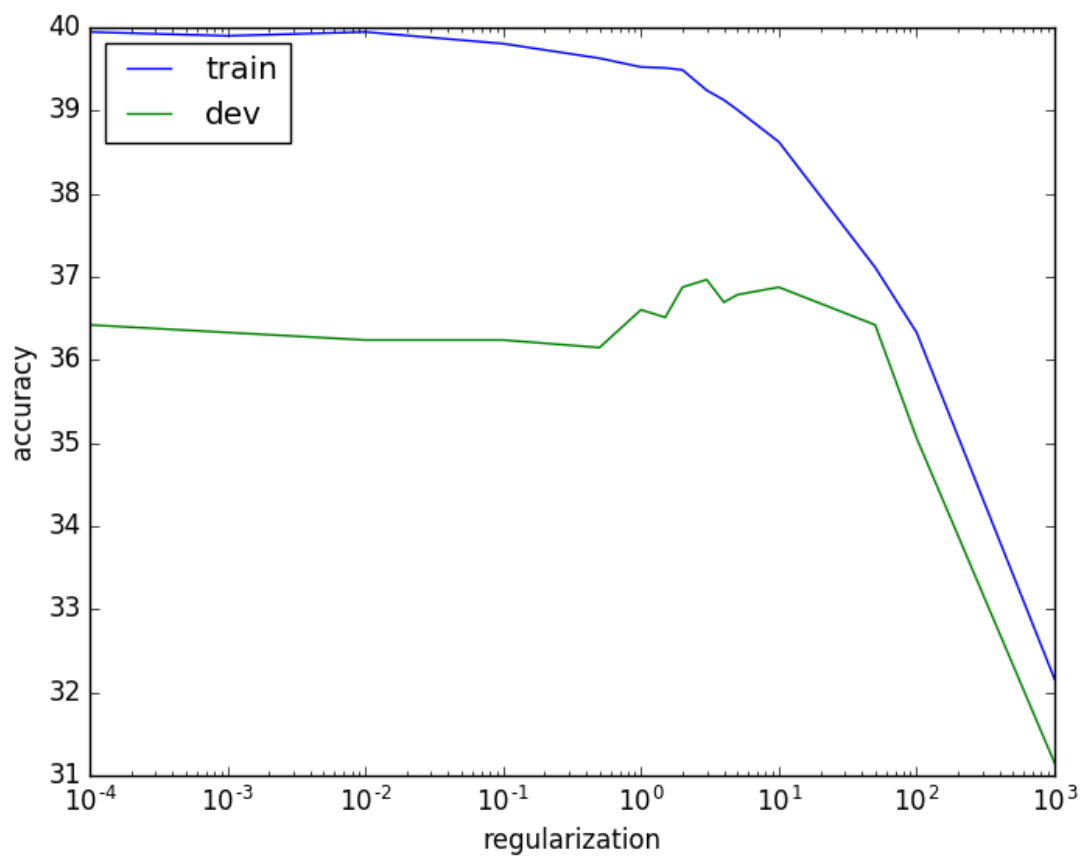
## (b)

> Reduced overfitting

## (c)

```
def getRegularizationValues():
    values = [0.0001, 0.001, 0.01, 0.1, 0.5, 1, 1.5, 2, 3, 4, 5, 10, 50,
100, 1000]
```
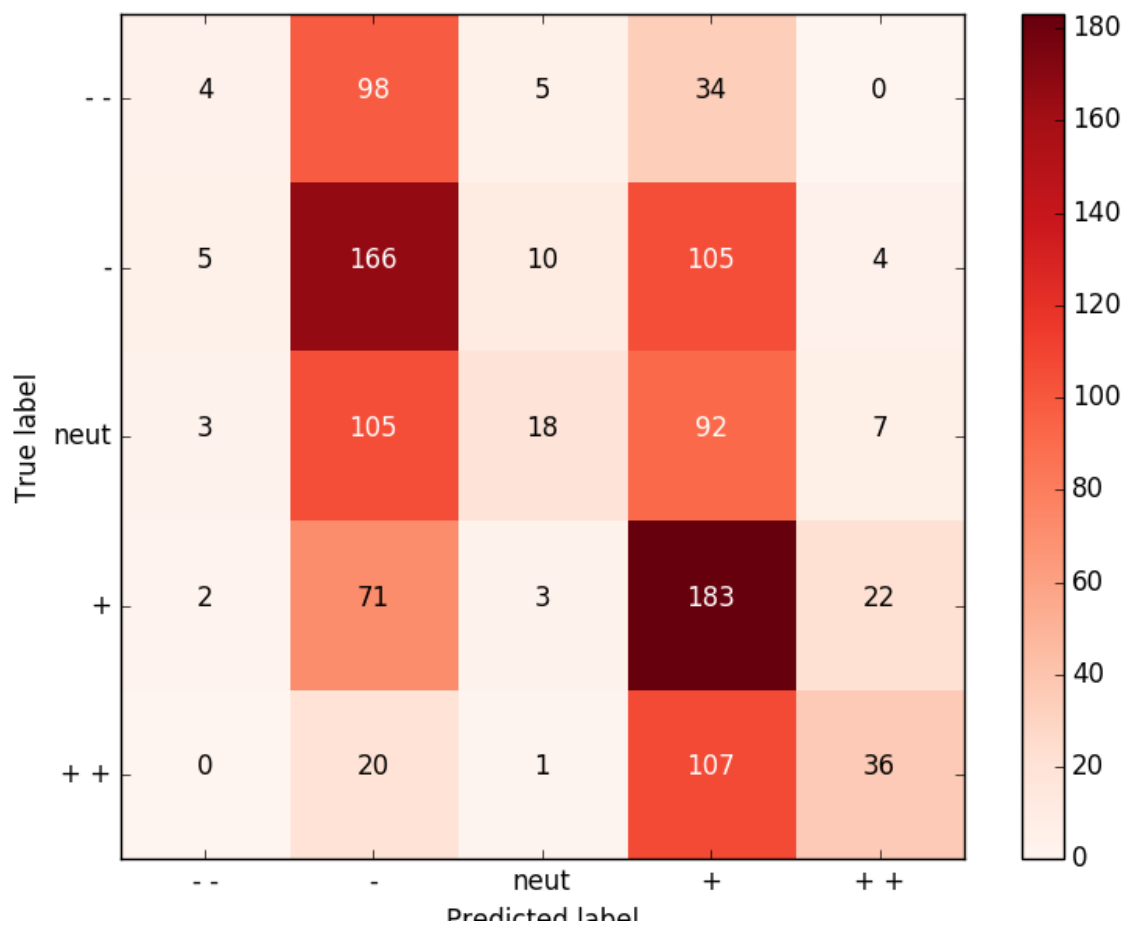
## (d)

> 1. python q4 sentiment.py --yourvectors
>    Train: 31.110
>    Dev: 32.698
>    Test: 30.407
> 2. python q4 sentiment.py --pretrained
>    Train: 39.244
>    Dev: 36.966
>    Test: 37.195

## (e)

**(f)**

**(g)**