

CS224N: Assignment #1

Kangwei Ling

February 14, 2017

1 Softmax

(a) eliminate the term c .

$$\text{softmax}(\mathbf{x} + \mathbf{c})_i = \frac{e^{\mathbf{x}_i + c}}{\sum_j e^{\mathbf{x}_j + c}} = \frac{e^c \cdot e^{\mathbf{x}_i}}{e^c \cdot \sum_j e^{\mathbf{x}_j}} = \text{softmax}(\mathbf{x})_i$$

(b) in `q1_softmax.py`.

2 Neural Network Basics

(a)

$$\frac{d}{dx}\sigma(x) = \frac{-1}{(1 + e^{-x})^2} \cdot (e^{-x})' = \frac{-1}{(1 + e^{-x})^2} \cdot (-e^{-x}) = \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} = \sigma(x)(1 - \sigma(x))$$

(b) As \mathbf{y} is the one-hot label vector, $\sum_j y_j = 1$, and

$$\frac{\partial \hat{y}_j}{\partial \theta_i} = \begin{cases} -\hat{y}_j \hat{y}_i & \text{if } i \neq j \\ \hat{y}_j(1 - \hat{y}_j) & \text{otherwise} \end{cases}$$

thus

$$\begin{aligned} \frac{\partial}{\partial \theta_i} CE(\mathbf{y}, \hat{\mathbf{y}}) &= -y_i(1 - \hat{y}_i) - \sum_{j \neq i} y_j(-\hat{y}_i) \\ &= -y_i + y_i \hat{y}_i + \hat{y}_i \sum_{j \neq i} y_j \\ &= \hat{y}_i \sum_j y_j - y_i \\ &= \hat{y}_i - y_i \end{aligned}$$

Therefore, $\frac{\partial}{\partial \theta} CE(\mathbf{y}, \hat{\mathbf{y}}) = \hat{\mathbf{y}} - \mathbf{y}$.

- (c) First find the gradients with respect to the hidden layer \mathbf{h} . Let $\mathbf{z}_2 = \mathbf{h}\mathbf{W}_2 + \mathbf{b}_2$, $\mathbf{z}_1 = \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1$, $\delta_2 = \frac{\partial(\cdot)}{\partial \mathbf{z}_2}$, $\delta_1 = \frac{\partial(\cdot)}{\partial \mathbf{z}_1}$. ($(\cdot) = CE(\mathbf{y}, \hat{\mathbf{y}})$, \circ means piecewise multiplication). thus

$$\begin{aligned}\delta_2 &= \hat{\mathbf{y}} - \mathbf{y} \\ \frac{\partial(\cdot)}{\partial \mathbf{h}} &= \delta_2 \mathbf{W}_2^T \\ \delta_1 &= \delta_2 \mathbf{W}_2^T \circ \sigma'(\mathbf{z}_1) \\ \delta_3 &= \frac{\partial(\cdot)}{\partial \mathbf{x}} = \delta_1 \mathbf{W}_1^T\end{aligned}$$

(d) $D_x H + H D_y + H + D_y$

(e) `q2_sigmoid.py`

(f) `q2_gradcheck.py`

(g) `q2_neural.py`

3 word2vec

- (a) Let $\mathbf{y}, \hat{\mathbf{y}}$ be column vectors. All vectors are column vectors.

$$\frac{\partial J}{\partial \mathbf{v}_c} = \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y})$$

- (b)

$$\frac{\partial J}{\partial \mathbf{u}_i} = (\hat{y}_i - y_i) \mathbf{v}_c$$

$$\frac{\partial J}{\partial \mathbf{U}} = \mathbf{v}_c (\hat{\mathbf{y}} - \mathbf{y})^T$$

- (c) gradients with respect to \mathbf{v}_c :

$$\frac{\partial J}{\partial \mathbf{v}_c} = (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{u}_o + \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{u}_k$$

gradients with respect to \mathbf{U} :

$$\frac{\partial J}{\partial \mathbf{u}_i} = \begin{cases} (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{v}_c & \text{if } i = o \\ (1 - \sigma(-\mathbf{u}_i^T \mathbf{v}_c)) \mathbf{v}_c & \text{if } i \in \{1, \dots, K\} \\ 0 & \text{otherwise} \end{cases}$$

(d) For skip-gram:

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

$$\frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{v}_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c}$$

$$\frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{v}_i} = \mathbf{0} \quad \text{for } i \neq c$$

For CBOW:

$$\frac{\partial J_{\text{CBOW}}}{\partial \mathbf{U}} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{U}}$$

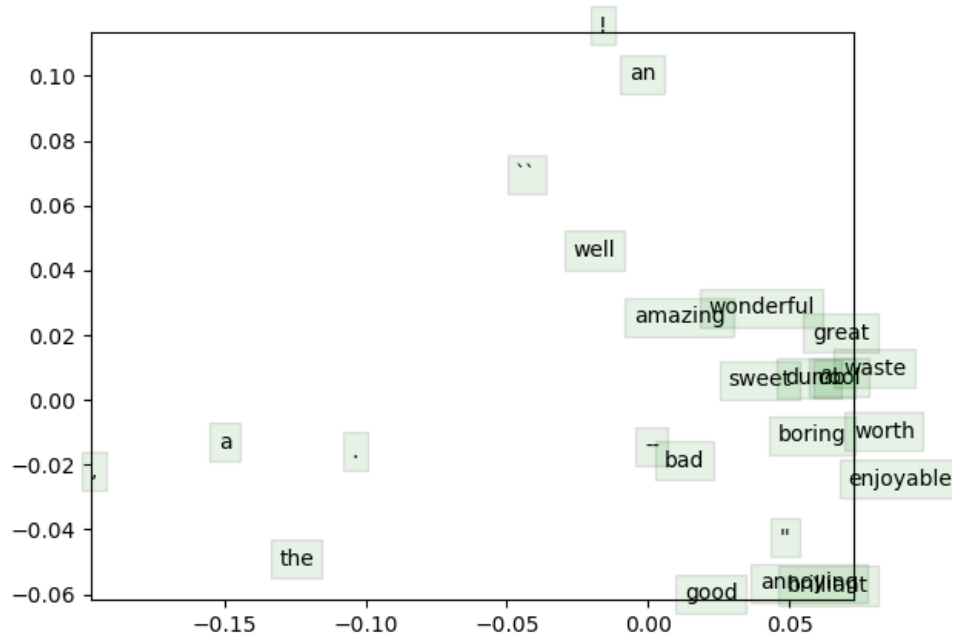
$$\frac{\partial J_{\text{CBOW}}}{\partial \mathbf{v}_j} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \hat{\mathbf{v}}} \quad \text{for } j \in \{c-m, \dots, c-1, c+1, \dots, c+m\}$$

$$\frac{\partial J_{\text{CBOW}}}{\partial \mathbf{v}_j} = \mathbf{0} \quad \text{for } j \notin \{c-m, \dots, c-1, c+1, \dots, c+m\}$$

(e) q3_word2vec.py

(f) q3_sgd.py

(g)



(h) q3_word2vec.py

4 Sentiment Analysis

- (a) `q4_sentiment.py`
- (b) To prevent out model from overfitting.
- (c) `q4_sentiment.py`
- (d)

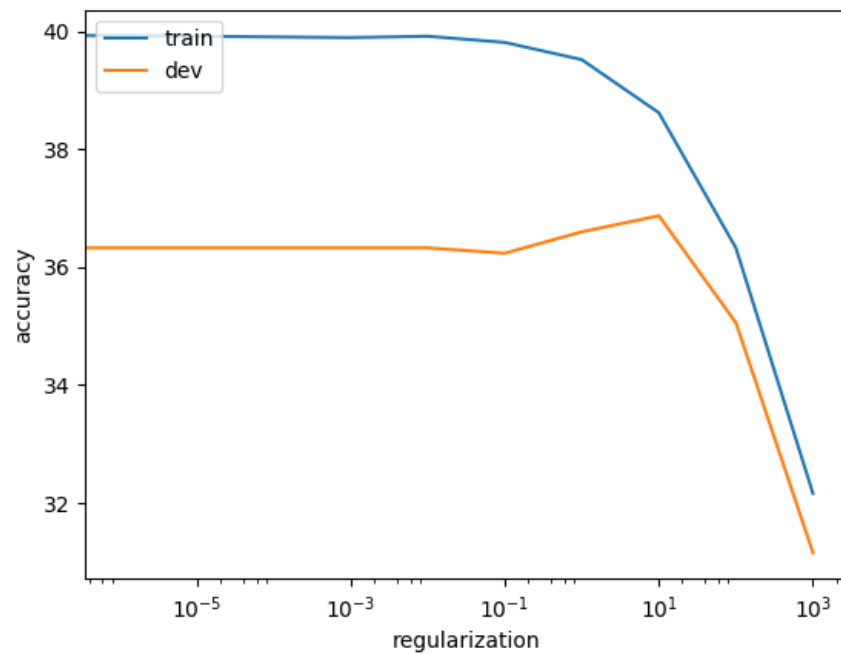
Best Accuracies	word vectors from problem 3	pretrained
train	31.121	39.934
dev	32.698	36.876
test	30.271	37.69

Table 1: best accuracies

The model using pretrained GloVe vectors outperforms for the following reasons:

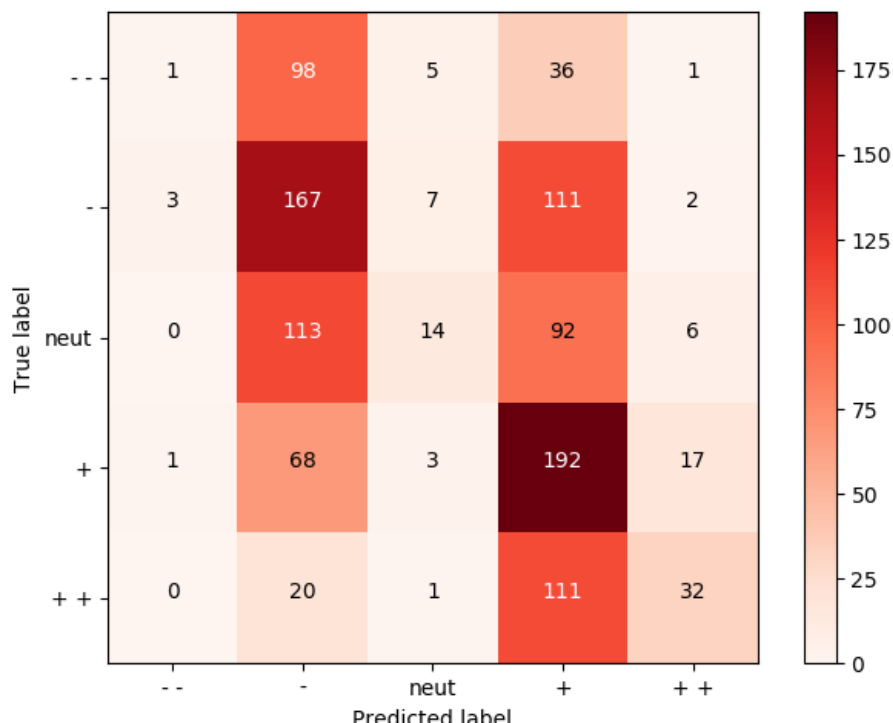
- The word vectors are trained on massive data (wikipedia), thus hold more information
- GloVe's vectors use different dimension from word2vec.
- GloVe represents words better than word2vec.

(e)



Regularization helped tackle overfitting, but regularization with too big weight harm the performance.

(f)



The mispredictions of the model mostly happened near neutral zone. For strong sentiment sentence, the prediction is quite accurate.

- (g)
- Some cold/negative words in the sentence cause the whole sentence to be predicted as negative, while the sentence has a positive context.(Or negative meaning get overrun with positive words)

nothing 's at stake , just a twisty double-cross you can smell a mile away
 – still , the derivative nine queens is lots of fun .

This sentence is with positive label(3), yet was predicted as negative(1).

- Some negation words are not exhibiting their negating effect.

scores no points for originality , wit , or intelligence .

True label(0), predicted label(3).

- Names hurt performance.

it takes a strange kind of laziness to waste the talents of robert forster ,
 anne meara , eugene levy , and reginald veljohnson all in the same movie

True label(0), predicted label(3).