

Assignment 1

frollo

CS224d: Deep Learning for Natural Language Processing

August 16, 2017

1 Softmax

(a)

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

2 Neural Network Basics

(a)

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

(b)

Assume the k-th dimension of \mathbf{y} is one.

$$\begin{aligned}
\frac{\partial}{\partial \theta_j} CE(\mathbf{y}, \hat{\mathbf{y}}) &= \frac{\partial}{\partial \theta_j} -\log \hat{\mathbf{y}}_k \\
&= -\frac{1}{\hat{\mathbf{y}}_k} \cdot \frac{e^{\theta_k} (\sum_i e^{\theta_i}) \cdot \mathbb{I}[j = k] - e^{\theta_k} \cdot e^{\theta_j}}{(\sum_i e^{\theta_i})^2} \\
&= -\frac{(\sum_i e^{\theta_i})}{e^{\theta_k}} \cdot \frac{e^{\theta_k} (\sum_i e^{\theta_i}) \cdot \mathbb{I}[j = k] - e^{\theta_k} \cdot e^{\theta_j}}{(\sum_i e^{\theta_i})^2} \\
&= \frac{e^{\theta_j}}{(\sum_i e^{\theta_i})} - \mathbb{I}[j = k] \\
&= (\hat{\mathbf{y}} - \mathbf{y})_j
\end{aligned}$$

After vectorizing the result, we have

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

(c)

Let

$$\begin{aligned}
\mathbf{z}^{(3)} &= \mathbf{h} \mathbf{W}^{(1)} + \mathbf{b}^{(1)} \\
\mathbf{z}^{(2)} &= \mathbf{x} \mathbf{W}^{(2)} + \mathbf{b}^{(2)}
\end{aligned}$$

We have

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{z}^{(3)}} J_{CE} &= \hat{\mathbf{y}} - \mathbf{y} \\
\frac{\partial}{\partial \mathbf{b}^{(2)}} J_{CE} &= \sum_i \frac{\partial J_{CE}}{\partial z_i^{(3)}} \frac{\partial z_i^{(3)}}{\partial \mathbf{b}^{(2)}} \\
&= \frac{\partial J_{CE}}{\partial \mathbf{z}^{(3)}} \\
&= \hat{\mathbf{y}} - \mathbf{y} \\
\frac{\partial}{\partial W_{xy}^{(2)}} J_{CE} &= \sum_i \frac{\partial J_{CE}}{\partial z_i^{(3)}} \frac{\partial z_i^{(3)}}{\partial W_{xy}^{(2)}} \\
&= \frac{\partial J_{CE}}{\partial z_y^{(3)}} \cdot h_x \\
\frac{\partial}{\partial \mathbf{W}^{(2)}} J_{CE} &= \mathbf{h}^T (\hat{\mathbf{y}} - \mathbf{y}) \\
\frac{\partial}{\partial a_x^{(2)}} J_{CE} &= \sum_i \frac{\partial J_{CE}}{\partial z_i^{(3)}} \frac{\partial z_i^{(3)}}{\partial a_i^{(2)}} \\
&= (\hat{\mathbf{y}} - \mathbf{y}) (\mathbf{W}_x^{(2)})^T \\
\frac{\partial}{\partial \mathbf{a}^{(2)}} J_{CE} &= (\hat{\mathbf{y}} - \mathbf{y}) (\mathbf{W}^{(2)})^T
\end{aligned}$$

$$\begin{aligned}
\frac{\partial}{\partial z_x^{(2)}} J_{CE} &= \sum_i \frac{\partial J_{CE}}{\partial z_i^{(3)}} \frac{\partial z_i^{(3)}}{\partial z_x^{(2)}} \\
&= \frac{\partial J_{CE}}{\partial a_x^{(2)}} \cdot \sigma'(z_x^{(2)}) \\
\frac{\partial}{\partial \mathbf{z}^{(2)}} J_{CE} &= \frac{J_{CE}}{\partial \mathbf{a}^{(2)}} \circ \sigma'(\mathbf{z}^{(2)}) \\
\frac{\partial}{\partial \mathbf{b}^{(1)}} J_{CE} &= \frac{\partial}{\partial \mathbf{z}^{(2)}} J_{CE} \\
\frac{\partial}{\partial \mathbf{W}^{(1)}} J_{CE} &= \mathbf{x}^T \frac{\partial J_{CE}}{\partial \mathbf{z}^{(2)}} \\
\frac{\partial}{\partial \mathbf{x}} J_{CE} &= \frac{\partial}{\partial \mathbf{a}^{(1)}} J_{CE} \\
&= \frac{\partial J_{CE}}{\partial \mathbf{z}^{(2)}} (\mathbf{W}^{(1)})^T
\end{aligned}$$

(d)

$$D_x \cdot H + H \cdot D_y + H + D_y$$

3 word2vec

(a)

Let $\mathbf{y}, \hat{\mathbf{y}} \in \mathbb{R}^{W \times 1}$, and $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{W \times d}$.

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

(b)

$$\begin{aligned}
\frac{\partial J_{CE}}{\partial (\mathbf{u}_w)_x} &= \sum_i \frac{\partial J_{CE}}{\partial \theta_i} \frac{\partial \theta_i}{\partial (\mathbf{u}_w)_k} \\
&= (\hat{\mathbf{y}} - \mathbf{y})_w \cdot (\mathbf{v}_c)_x \\
\frac{\partial J_{CE}}{\partial \mathbf{u}_w} &= (\hat{\mathbf{y}} - \mathbf{y})_w \cdot \mathbf{v}_c \\
\frac{\partial J_{CE}}{\partial \mathbf{U}} &= (\hat{\mathbf{y}} - \mathbf{y})^T (\mathbf{v}_c)
\end{aligned}$$

(c)

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

Let \mathbf{u}_w be the word vector of word w .

$$\frac{\partial J_{NS}}{\partial \mathbf{u}_w} = \begin{cases} -[1 - \sigma(\mathbf{u}_w^T \mathbf{v}_c)] \mathbf{v}_c & \text{if } w = o \\ [1 - \sigma(-\mathbf{u}_w^T \mathbf{v}_c)] \mathbf{v}_c & \text{if } w \neq o, w \in \{1, 2, \dots, K\} \\ 0 & \text{otherwise} \end{cases}$$

In J_{CE} we need to calculate gradient of all W word vectors, but the number decreases to $K + 1$ when we apply negative sampling. (speed-up ration $\approx \frac{K}{W}$)
(d)

$$\begin{aligned} \frac{\partial J_{SG}}{\partial \mathbf{v}_c} &= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c} \\ \frac{\partial J_{SG}}{\partial \mathbf{U}} &= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{U}} \end{aligned}$$

$$\begin{aligned} \frac{\partial J_{CBOW}}{\partial \mathbf{v}_{c+j}} &= \frac{\partial J_{CBOW}}{\partial \hat{\mathbf{v}}} \frac{\partial \hat{\mathbf{v}}}{\partial \mathbf{v}_{c+j}} \\ &= \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \hat{\mathbf{v}}} \end{aligned}$$

$$\frac{\partial J_{CBOW}}{\partial \mathbf{u}_w} = \begin{cases} \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{u}_w} & \text{if } w = w_c \\ 0 & \text{otherwise} \end{cases}$$

(g)

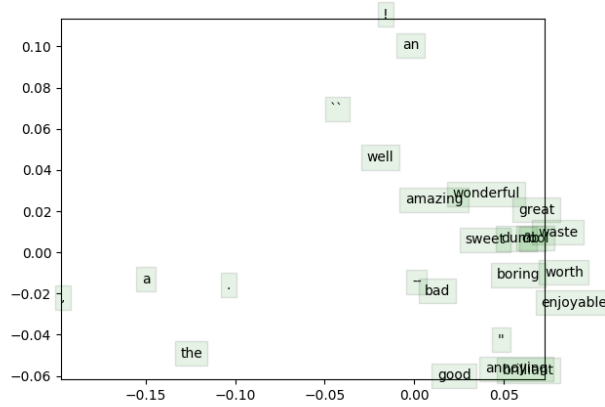


Figure 1: q3_word_vectors.png

4 Sentiment Analysis

(b)

Imposing such a constraint can be interpreted as the prior Bayesian belief that the optimal weights are close to zero.

(c)

```
bestResult = None
for res in results:
    if bestResult == None or \
        bestResult["dev"] < res["dev"]:
        bestResult = res
return res
```

(d)

option	train	dev	test
yourvectors	31.133	32.607	30.362
pretrained	38.642	36.876	37.692

Table 1: Comparison of two different word vector source

The training data used in pretrained wordvectors is larger than ours. It has larger vocabulary size. Also the pretrained model is GloVe which is more powerful than our simple skipgram model.

(e)

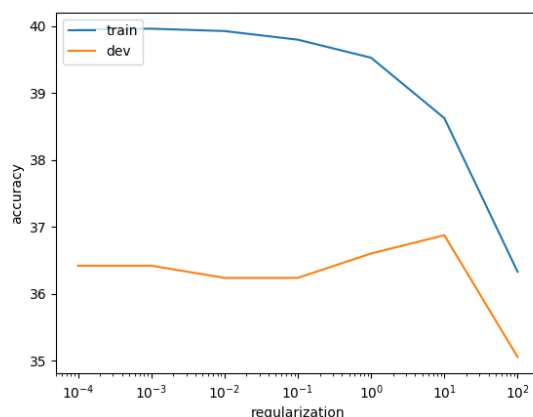


Figure 2: q4_reg_v_acc.png

When λ is small, we have the best training accuracy. That means we let the model overfit our training data. By increasing λ we may find better developing accuracy, and we believe that it is more generalized.

(f)

It seems like our model will predict + and - the most. Not the neutral one.

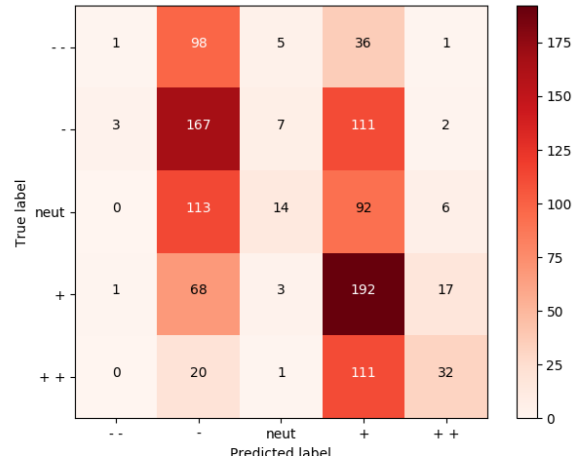


Figure 3: q4_dev_conf.png

(g)

(True: 3, Predicted: 1) and if you 're not nearly moved to tears by a couple of senses, you 've got ice water in your veins .

I guess the model doesn't know the if statement and it just gives a low score caused by some negative word such as tear.

(True: 0, Predicted: 2) i had to look away - this was god awful .

The usage of word "god" is difficult to our model.

(True: 1, Predicted: 3) the best that can be said about the work here of scottish director ritchie ... is that he obviously does n't have his heart in it .

Um... maybe our model doesn't know the meaning of sentence "the best that...".