

CS224n HW1

Xinyu Tan

July 19, 2017

1 Softmax

(a)

Take a look at any element i ,

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

Therefore, we have $\text{softmax}(\mathbf{x} + c) = \text{softmax}(\mathbf{x})$

(b)

Note the first case illustrate the broadcasting principle (dimension match) in numpy:

```
1 if len(x.shape) > 1:
2     #matrix
3     x = x - np.max(x, axis=1, keepdims=True)
4     x = np.exp(x) / np.sum(np.exp(x), axis=1, keepdims=True)
5 else:
6     # vector
7     x = x - x.max() # normalize
8     x = np.exp(x) / np.sum(np.exp(x))
```

2 Neural Network Basics

(a)

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

(b)

First, we have

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

For one-hot encoding, only k -th element in \mathbf{y} is *one*, so we have

$$\begin{aligned} CE(\mathbf{y}, \hat{\mathbf{y}}) &= -\log \hat{y}_k = -\log \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} \\ &= -\theta_k + \log \sum_j e^{\theta_j} \end{aligned}$$

Therefore,

$$\begin{aligned} \frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial \theta_k} &= -1 + \frac{e^{\theta_k}}{\sum_j e^{\theta_j}} \\ \frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial \theta_i} &= \frac{e^{\theta_i}}{\sum_j e^{\theta_j}}, \forall i \neq k \end{aligned}$$

Put them altogether,

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

(c)

We have

$$\begin{aligned} \hat{\mathbf{y}} &= \text{softmax}(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2) \\ &= \text{softmax}(\text{sigmoid}(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1) + \mathbf{b}_2) \end{aligned}$$

Denote $\mathbf{a}_2 = \mathbf{h}\mathbf{W}_2 + \mathbf{b}_2$ and $\mathbf{a}_1 = \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1$, then

$$\begin{aligned} \frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{x}} &= \frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{a}_2} \frac{\partial \mathbf{a}_2}{\partial \mathbf{x}} \\ &= (\hat{\mathbf{y}} - \mathbf{y}) \frac{\partial \mathbf{a}_2}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \\ &= (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^T \sigma'(\mathbf{a}_1) \mathbf{W}_1^T \end{aligned}$$

(d)

There are in total $D_x H + H + H D_y + D_y$ parameters.

3 word2vec

(a)

For a word \mathbf{o} , the loss function is

$$\begin{aligned} L &= -y_o \log \hat{y}_o = -\log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w=1}^W \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{u}_o^T \mathbf{v}_c + \log \sum_{w=1}^W \exp(\mathbf{u}_w^T \mathbf{v}_c) \end{aligned}$$

Therefore, derivative with respect to \mathbf{v}_c is:

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{v}_c} &= -\mathbf{u}_o + \frac{1}{\sum_{l=1}^W \exp(\mathbf{u}_l^T \mathbf{v}_c)} \sum_{w=1}^W \exp(\mathbf{u}_w^T \mathbf{v}_c) \mathbf{u}_w \\ &= -\mathbf{u}_o + \sum_{w=1}^W p(w|c) \mathbf{u}_w\end{aligned}$$

where

$$p(w|c) = \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{l=1}^W \exp(\mathbf{u}_l^T \mathbf{v}_c)}$$

(b)

Similarly,

$$\frac{\partial L}{\partial \mathbf{u}_w} = \begin{cases} -(1 - \hat{y}_o) \mathbf{v}_c & w = o \\ \hat{y}_w \mathbf{v}_c & w \neq o \end{cases}$$

(c)

Given

$$J_{\text{neg-sample}}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log \sigma(\mathbf{u}_o^T \mathbf{v}_c) - \sum_{k=1}^K \log \sigma(-\mathbf{u}_o^T \mathbf{v}_c)$$

we have

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

and

$$\frac{\partial J}{\partial \mathbf{u}_w} = \begin{cases} -\sigma(-\mathbf{u}_w^T \mathbf{v}_c) \mathbf{v}_c & w = o \\ \sigma(\mathbf{u}_w^T \mathbf{v}_c) \mathbf{v}_c & w \neq o \end{cases}$$

(d)

Let's denote \mathbf{U} the matrix that aggregates all the output word vectors.

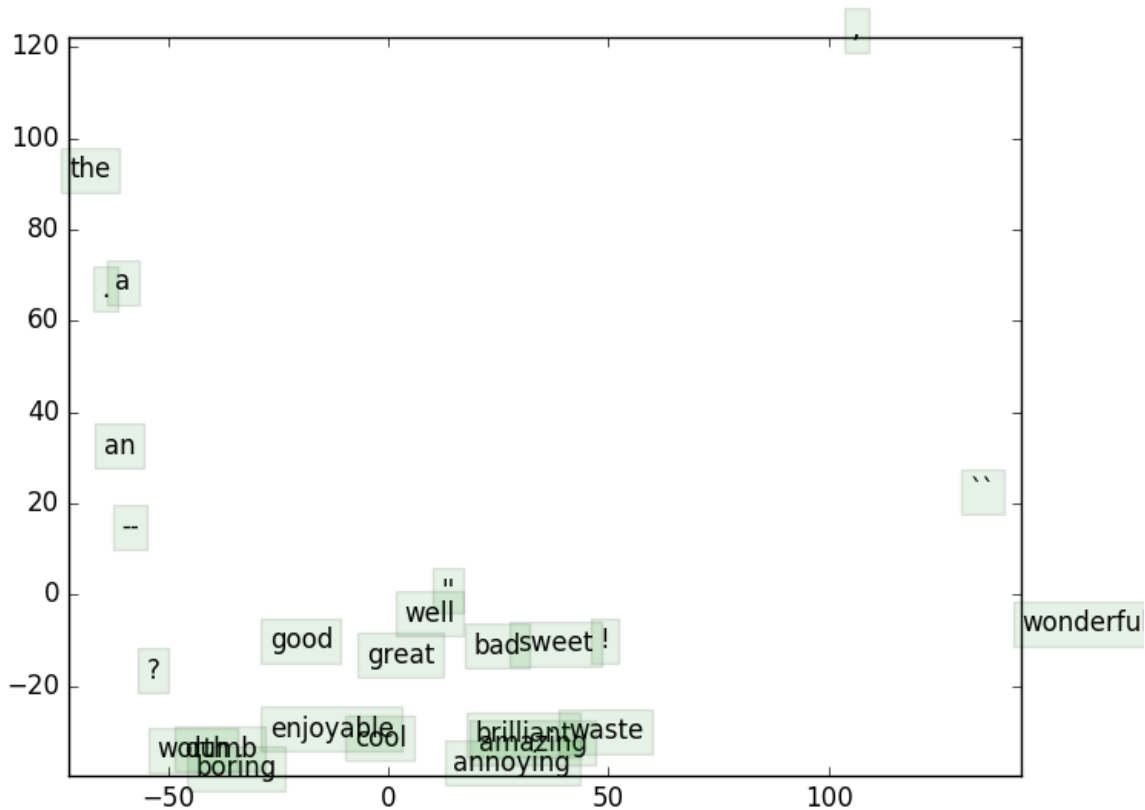
1. Skip-gram:

$$\begin{aligned}\frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{U}} &= \sum_{j \neq 0} \frac{\partial F(\mathbf{u}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{U}} \\ \frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{v}_c} &= \sum_{j \neq 0} \frac{\partial F(\mathbf{u}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c}\end{aligned}$$

2. CBOW:

$$\begin{aligned}\frac{\partial J_{\text{CBOW}}}{\partial \mathbf{U}} &= \frac{\partial F(\mathbf{u}_c, \hat{\mathbf{v}})}{\partial \mathbf{U}} \\ \frac{\partial J_{\text{CBOW}}}{\partial \mathbf{v}_j} &= \frac{\partial F(\mathbf{u}_c, \hat{\mathbf{v}})}{\partial \hat{\mathbf{v}}}\end{aligned}$$

(g)



Summary: “the”, “a”, and “an” are close to each other; most adjectives are clustered at the bottom of the graph; “wonderful” somehow is far from its similar words.

4 Sentiment Analysis

(b)

The main reason to introduce regularization is to avoid overfitting.

(c)

Code for choosing the best model:

```
1 def chooseBestModel(results):
2     """Choose the best model based on parameter tuning on the dev set
3
4     Arguments:
5     results — A list of python dictionaries of the following format:
6         {
7             "reg": regularization,
8             "clf": classifier,
9             "train": trainAccuracy,
10            "dev": devAccuracy,
11            "test": testAccuracy
```

```
12     }
13
14     Returns:
15     Your chosen result dictionary.
16     """
17     bestResult = None
18
19     ### YOUR CODE HERE
20     test_acc = []
21     for res in results:
22         test_acc.append(res["dev"])
23     ### END YOUR CODE
24     best_index = test_acc.index(max(test_acc))
25     bestResult = results[best_index]
26     return bestResult
```