

CS224n HW1

Xinyu Tan

July 4, 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.

(e)

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