

Assignment 2: Word Vectors

Junzhe Wang - September 19, 2019

1. Basics

a. Prove softmax is invariant to constant offset in the input.

$$\begin{aligned} ((\text{softmax}(x + c))_i &= \frac{\exp(x_i + c)}{\sum_{j=1}^{\dim(x)} \exp(x_j + c)} = \frac{\exp(x_i) * \exp(c)}{\sum_{j=1}^{\dim(x)} (\exp(x_j) * \exp(c))} \\ &= \frac{\exp(x_i) * \exp(c)}{\exp(c) * \sum_{j=1}^{\dim(x)} \exp(x_j)} = \frac{\exp(x_i)}{\sum_{j=1}^{\dim(x)} \exp(x_j)} = \text{softmax}(x)_i \end{aligned}$$

b. See softmax.py

c.

Sigmoid :	$\sigma(x) = \frac{1}{1 + e^{-x}}$
Derivative of sigmoid:	$\begin{aligned} \frac{d}{dx}(1 + e^{-x})^{-1} &= (-1) * (1 + e^{-x})^{-2} * \frac{d}{dx}(1 + e^{-x}) = \frac{-1}{(1 + e^{-x})^2} \frac{d}{dx}(e^{-x}) \\ &= \frac{-1}{(1 + e^{-x})^2} * (-1) * e^{-x} = \frac{e^{-x}}{(1 + e^{-x})^2} \end{aligned}$
$\sigma(x)(1 - \sigma(x))$	$\frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}}\right) = \frac{e^{-x}}{(1 + e^{-x})^2}$

2. Word2vec

a	$J = \sum_{i=1}^W y_i \log\left(\frac{\exp(u_i^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}\right) = \sum_{i=1}^W y_i [u_i^T v_c - \log(\sum_{w=1}^W \exp(u_w^T v_c))]$
	<p>since y is one-hot vector, there is only one non-zero item in the summation. If say oth is the current word, $y_o = 1$, and $y_{i \neq o} = 0$. The loss could be written as :</p>
	$J = (u_o^T v_c - \log(\sum_{w=1}^W \exp(u_w^T v_c)))$
	$\frac{\partial J}{\partial v_c} = -[u_o - \frac{\partial(\log(\sum_{w=1}^W \exp(u_w^T v_c)))}{\partial v_c}] = -[u_o - \frac{\sum_{w=1}^W \exp(u_w^T v_c) u_w}{\sum_{w=1}^W \exp(u_w^T v_c)}]$ $= \sum_{w=1}^W \frac{\exp(u_w^T v_c) u_w}{\sum_{w=1}^W \exp(u_w^T v_c)} - u_o = \sum_{w=1}^W \hat{y}_w u_w - u_o = U \hat{y} - U y = U(\hat{y} - y)$
b	$\frac{\partial J}{\partial u_w} = -[v_c - \frac{\partial(\log(\sum_{w=1}^W \exp(u_w^T v_c)))}{\partial u_w}] = -[v_c - \frac{\sum_{w=1}^W \exp(u_w^T v_c) v_c}{\sum_{w=1}^W \exp(u_w^T v_c)}]$ $\sum_{w=1}^W \frac{\exp(u_w^T v_c) v_c}{\sum_{w=1}^W \exp(u_w^T v_c)} - v_c = \sum_{w=1}^W \hat{y}_w v_c - v_c = (\hat{y} - 1)v_c = (\hat{y} - y)v_c$
c	<p>Since $o \in \{1, \dots, K\}$, the second part is 0</p> $\frac{\partial J}{\partial u_o} = -\frac{\sigma(u_o^T v_c)(1 - \sigma(u_o^T v_c))v_c}{\sigma(u_o^T v_c)}$ $= -(1 - \sigma(u_o^T v_c))v_c$
d	$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m, \dots, c+m})}{\partial U} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_j, v)}{\partial U}$
	$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m, \dots, c+m})}{\partial v_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_j, v)}{\partial v_c}$
	$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m, \dots, c+m})}{\partial v_j} = 0, \text{ for all } j \neq c$

- g. In the run.py file, I tested the KNN function by computing the 5 nearest neighbor of the first word('great') among the visualizeWords, which are:
 ["great", "cool", "brilliant", "wonderful", "well", "amazing", "worth", "sweet",
 "enjoyable", "boring", "bad", "dumb", "annoying", "female", "male", "queen", "king",
 "man", "woman", "rain", "snow", "hail", "coffee", "tea"], and the results are ['well', 'annoying',
 'amazing', 'bad', 'queen']

```
(py37) MacBook-Pro-3:a2 joey$ python run.py
sanity check: cost at convergence should be around or below 10
training took 0 seconds
the nearest words to great are ['well', 'annoying', 'amazing', 'bad', 'queen']
```

