

Author: Benjamin Smidt  
Created: October 3rd, 2022  
Last Updated: January 17th, 2023

## CS 224N A2: Word Vectors

*Note to the reader.* This is my work for assignment 2 of Stanford’s course CS 224N: Natural Language Processing with Deep Learning. You can find the Winter 2021 lectures on YouTube here. This document is meant to be a reference, explanation, and resource for assignment 2. If there’s a typo or a error, please email me at [benjamin.smidt@utexas.edu](mailto:benjamin.smidt@utexas.edu) so I can fix it. Finally, here is a link to my GitHub repo.

## Contents

<b>1</b>	<b>Written Understanding</b>	<b>2</b>
1.1	Problem A . . . . .	2
1.2	Problem B . . . . .	3
1.3	Problem C . . . . .	4
1.4	Problem D . . . . .	6
1.5	Problem E . . . . .	6
1.6	Problem F . . . . .	7
1.7	Problem G . . . . .	7
1.7.1	(i) . . . . .	8
1.7.2	(ii) . . . . .	10
1.7.3	(iii) . . . . .	10
1.8	Problem H . . . . .	11
1.9	Problem I . . . . .	12
<b>2</b>	<b>Programming</b>	<b>12</b>
2.1	Sigmoid . . . . .	12
2.2	Naive Softmax Loss and Gradient . . . . .	13
2.3	Negative Sampling Loss and Gradient . . . . .	13
2.4	Skipgram . . . . .	14
2.5	Stochastic Gradient Descent . . . . .	14
<b>3</b>	<b>References</b>	<b>15</b>

# 1 Written Understanding

## 1.1 Problem A

### Instructions

Prove that the naive-softmax loss is the same as the cross-entropy loss between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$ , i.e. (note that  $\mathbf{y}, \hat{\mathbf{y}}$  are vectors and  $\hat{y}_o$  is a scalar).

### Solution

Remember, we're using the skip-gram model (see lecture 1 notes, very helpful for definitions and general understanding). To start with, our naive-softmax loss is defined as

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c) \quad (1)$$

where

$$P(O = o | C = c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (2)$$

Let's define our variables. Recall that  $\mathbf{v}_c$  is the vector embedding for our given center word ( $\mathbf{v}_c = Vx$  where  $V \in \mathbb{R}^{n \times |\text{Vocab}|}$  and  $x$  is the one hot vector for the center word where  $x \in \mathbb{R}^{|\text{Vocab}|}$ ).  $\hat{\mathbf{y}}$  is our score vector and  $\hat{y}_w$  is the scalar at index  $w$  (denoting a given word) within  $\hat{\mathbf{y}}$ .

$$\hat{\mathbf{y}} = \frac{\exp(\mathbf{U}\mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad \text{and} \quad \hat{y}_w = \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (3)$$

Note that  $\hat{\mathbf{y}} \in \mathbb{R}^{|\text{Vocab}|}$  (it is a vector of length equal to the vocabulary size). Additionally, it's numerator is of the same dimension as  $\hat{\mathbf{y}}$  but its denominator is a scalar. You can interpret the expression to mean that each index of vector  $\exp(\mathbf{U}\mathbf{v}_c)$  is divided by  $\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)$  to yield  $\hat{\mathbf{y}}$  (in fact this is exactly how you might think of programming something like this). Thus, each index  $w$  in the vector  $\hat{\mathbf{y}}$  (scalar  $\hat{y}_w$ ) can be interpreted as the probability that the corresponding word (using the the index  $w$  and its one hot vector) is a context (or "outside") word given the center word (by eqn. 2)

$\mathbf{y}$  is defined as the one hot vector of the true outside word such that index  $o$  in  $\mathbf{y}$  is 1 and all other indices are 0 for this particular problem. Then

$$-\sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{y}_w) = -\mathbf{y}_1 \log(\hat{y}_1) + \dots + -\mathbf{y}_o \log(\hat{y}_o) + \dots + -\mathbf{y}_w \log(\hat{y}_w) \quad (4)$$

$$-\sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{y}_w) = -(0) \log(\hat{y}_1) + \dots + -(1) \log(\hat{y}_o) + \dots + -(0) \log(\hat{y}_w) \quad (5)$$

$$-\sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{y}_w) = -\log(\hat{y}_o) \quad (6)$$

Remember that  $\mathbf{y}_o$  represents the word at index  $o$ , who's one hot vector is a 1 at index  $o$ , indicating that it is a context word given  $\mathbf{v}_c$ . That is,

$$\hat{y}_o = P(O = o | C = c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (7)$$

Hence, by substitution

$$- \sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) = -\log\left(\frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}\right) \quad (8)$$

$$- \sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) = -\log P(O = o | C = c) \quad (9)$$

$$- \sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) = \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) \quad (10)$$

I know the instructions said one line but I was going for clarity here. Obviously, I could not define the variables (leave it to you to figure out what they mean) and just write some one liner that connects the dots. However, I wanted to make this as clear to understand as possible. Hopefully this leaves no room for ambiguity.

## 1.2 Problem B

### Instructions

Compute the partial derivative of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to  $\mathbf{v}_c$ . Please write your answer in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{U}$ . Additionally, answer the following two questions with one sentence each: (1) When is the gradient zero? (2) Why does subtracting this gradient, in the general case when it is nonzero, make  $\mathbf{v}_c$  a more desirable vector (namely, a vector closer to outside word vectors in its window)?

### Solution

We start with our definition of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ .

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (11)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = \frac{\partial}{\partial \mathbf{v}_c} -\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (12)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = \frac{\partial}{\partial \mathbf{v}_c} -\mathbf{u}_o^\top \mathbf{v}_c + \frac{\partial}{\partial \mathbf{v}_c} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \quad (13)$$

Everything up to this point should be easy to follow if you remember algebra (although see the first lecture where he goes through these exact steps in detail if you are confused). Note that in the coming step, I do a change of variables to ensure I know what I'm taking my partial derivative with respect to.

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\mathbf{u}_o^\top + \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sum_{j \in \text{Vocab}} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \quad (14)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\mathbf{u}_o^\top + \frac{\sum_{j \in \text{Vocab}} \mathbf{u}_j^\top \exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (15)$$

Recall that

$$\hat{\mathbf{y}}_w = \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (16)$$

By substitution, (I change back the  $j$  to  $w$  since there's only one sum now)

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\mathbf{u}_o^\top + \sum_{j \in \text{Vocab}} \mathbf{u}_w^\top \hat{\mathbf{y}}_w \quad (17)$$

Note that the expression to the right of the addition (above) is each index  $w$  of matrix  $\mathbf{U} \in \mathbb{R}^{|\text{Vocab}| \times n}$  (vector  $\mathbf{u}_w$ ) being weighted by each index  $w$  of vector  $\hat{\mathbf{y}}$  (scalar  $\hat{\mathbf{y}}_w$ ). Rewriting this using matrices (remember  $\mathbf{y}$  is one hot where  $\mathbf{y}_o = 1$ ), we find

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\mathbf{U}^\top \mathbf{y} + \mathbf{U}^\top \hat{\mathbf{y}}_w \quad (18)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = \mathbf{U}^\top (\hat{\mathbf{y}} - \mathbf{y}) \quad (19)$$

You can verify the shape of this vector is the same shape as  $\mathbf{v}_c \in \mathbb{R}^n$ .

1. The gradient is zero when  $\hat{\mathbf{y}} = \mathbf{y}$ . Obviously if our predicted and correct vectors are equivalent then our accuracy is perfect and there's no update that could improve the loss.
2. Because we're doing gradient descent (as opposed to ascent), the update adds some portion of  $\mathbf{U}^\top \mathbf{y}$  and subtracts some portion of  $\mathbf{U}^\top \hat{\mathbf{y}}$ . This makes intuitive sense because adding the correct vector  $\mathbf{u}_o^\top$  makes it more like that vector (which is what we want) and subtracting by  $\mathbf{U}^\top \hat{\mathbf{y}}$ , the weighted average of our incorrect vectors that are producing the loss, makes it less like those vectors (again, what we want).

### 1.3 Problem C

#### Instructions

Compute the partial derivatives of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to each of the 'outside' word vectors,  $\mathbf{u}_w$ 's. There will be two cases: when  $w = o$ , the true 'outside' word vector, and  $w \neq o$ , for all other words. Please write your answer in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{v}_c$ . In this subpart, you may use specific elements within these terms as well (such as  $\mathbf{y}_1, \mathbf{y}_2, \dots$ ). Note that  $\mathbf{u}_w$  is a vector while  $\mathbf{y}_1, \mathbf{y}_2, \dots$  are scalars.

#### Solution

**Case 1:**  $\mathbf{u}_w = \mathbf{u}_o$

We start with the case that  $\mathbf{u}_w = \mathbf{u}_o$ .

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (20)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = \frac{\partial}{\partial \mathbf{u}_o} - \log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (21)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = \frac{\partial}{\partial \mathbf{u}_o} - \mathbf{u}_o^\top \mathbf{v}_c + \frac{\partial}{\partial \mathbf{u}_o} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \quad (22)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = -\mathbf{v}_c + \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{u}_o} \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \quad (23)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = -\mathbf{v}_c + \frac{\mathbf{v}_c \exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (24)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = -\mathbf{v}_c + \mathbf{v}_c(\hat{\mathbf{y}}_o) \quad (25)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = \mathbf{v}_c(\hat{\mathbf{y}}_o - 1) \quad (26)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = \mathbf{v}_c(\hat{\mathbf{y}}^\top \mathbf{y} - 1) \quad (27)$$

Revisit eqn. 3 if going from step 24 to step 25 isn't clicking. We choose to represent  $\hat{\mathbf{y}}_o$  as  $\hat{\mathbf{y}}^\top$  from the direction's constraints of how to represent the equation. I wasn't sure if eqn. 26 was good enough or it was meant to be more generalized.

### Case 2: $\mathbf{u}_w \neq \mathbf{u}_o$

Now we move onto the case that  $\mathbf{u}_w \neq \mathbf{u}_o$ . That is, the gradient with respect to any vector  $\mathbf{u}_w$  that isn't  $\mathbf{u}_o$ . I'll use the notation  $\mathbf{u}_j$  to indicate the particular  $\mathbf{u}_w$  we want to find the gradient with respect to and (hopefully) prevent any confusion.

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (28)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_j} = \frac{\partial}{\partial \mathbf{u}_j} - \log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (29)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_j} = \frac{\partial}{\partial \mathbf{u}_j} - \mathbf{u}_o^\top \mathbf{v}_c + \frac{\partial}{\partial \mathbf{u}_j} \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \quad (30)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_j} = \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{u}_j} \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \quad (31)$$

$$= \frac{1}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{u}_j} (\exp(\mathbf{u}_1^\top \mathbf{v}_c) + \dots + \exp(\mathbf{u}_j^\top \mathbf{v}_c) + \dots + \exp(\mathbf{u}_{|\text{Vocab}|}^\top \mathbf{v}_c)) \quad (32)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_j} = \frac{\mathbf{v}_c \exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (33)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_j} = \mathbf{v}_c \hat{\mathbf{y}}_j \quad (34)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_j} = \mathbf{v}_c(\hat{\mathbf{y}}_j - \mathbf{y}_j) \quad (35)$$

Note that  $\hat{\mathbf{y}}_j$  is a scalar and  $\mathbf{y}_j$  is 0 since the word  $j$  is defined as not being an outside word.

## 1.4 Problem D

### Instructions

Write down the partial derivative of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to  $\mathbf{U}$ . Please break down your answer in terms of  $\frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_1}, \frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_2}, \dots, \frac{\partial \mathbf{J}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{|\text{Vocab}|}}$ . The solution should be one or two lines long.

### Solution

We already know

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_j} = \mathbf{v}_c(\hat{\mathbf{y}}_j - \mathbf{y}_j) \quad \text{and} \quad \frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = \mathbf{v}_c(\hat{\mathbf{y}}_o - 1) \quad (36)$$

Since  $\mathbf{y}$  is a one hot vector, then With that we can write

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = \frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_1}, \dots, \frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_V} \quad (37)$$

where  $V$  is the index of the last word vector in  $\mathbf{U}$ . It follows then that

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = \mathbf{v}_c(\hat{\mathbf{y}}_1 - \mathbf{y}_1), \mathbf{v}_c(\hat{\mathbf{y}}_2 - \mathbf{y}_2), \dots, \mathbf{v}_c(\hat{\mathbf{y}}_o - 1), \dots, \mathbf{v}_c(\hat{\mathbf{y}}_{|V|} - \mathbf{y}_{|V|}) \quad (38)$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = \mathbf{v}_c(\hat{\mathbf{y}} - \mathbf{y})^\top \quad (39)$$

Our current equation suggests  $\mathbf{U} \in \mathbb{R}^{n \times |\text{Vocab}|}$ . However,  $\mathbf{U} \in \mathbb{R}^{|\text{Vocab}| \times n}$ , the transpose of what we currently have. This comes from the fact that  $\mathbf{u}_w$  is thought of as a column vector when alone (or at least that's how I interpret it) but it is represented as a row vector in  $\mathbf{U}$ . So, we fix it with the following to achieve the correct shape.

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{v}_c^\top \quad (40)$$

## 1.5 Problem E

### Instructions

The ReLU (Rectified Linear Unit) activation function is given by the Equation:

$$f(x) = \max(0, x) \quad (41)$$

Please compute the derivative of  $f(x)$  with respect to  $x$ , where  $x$  is a scalar. You may ignore the case that the derivative is not defined at 0.

### Solution

$$f(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x > 0 \end{cases} \quad (42)$$

$$f'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases} \quad (43)$$

## 1.6 Problem F

### Instructions

The sigmoid function is given by:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (44)$$

Please compute the derivative of  $\sigma(x)$  with respect to  $x$ , where  $x$  is a scalar. Hint: you may want to write your answer in terms of  $\sigma(x)$ .

### Solution

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad \text{Multiply by } \frac{e^x}{e^x} \quad (45)$$

$$\sigma'(x) = \frac{e^x(e^x + 1) - e^x e^x}{(e^x + 1)^2} \quad \text{Product Rule} \quad (46)$$

$$\sigma'(x) = \frac{e^{2x} + e^x - e^{2x}}{(e^x + 1)^2} \quad (47)$$

$$\sigma'(x) = \frac{e^x}{(e^x + 1)^2} \quad (48)$$

$$\sigma'(x) = \frac{e^x}{e^x + 1} \frac{1}{e^x + 1} \quad (49)$$

$$\sigma'(x) = \frac{1}{1 + e^{-x}} \frac{1}{e^x + 1} \quad \text{Multiplied left by } \frac{\frac{1}{e^x}}{\frac{1}{e^x}} \quad (50)$$

$$\sigma'(x) = \sigma(x) \frac{1}{e^x + 1} \quad (51)$$

$$\sigma'(x) = \sigma(x) \frac{e^x + 1 - e^x}{e^x + 1} \quad (52)$$

$$\sigma'(x) = \sigma(x) \left( \frac{e^x + 1}{e^x + 1} - \frac{e^x}{e^x + 1} \right) \quad (53)$$

$$\sigma'(x) = \sigma(x) \left( 1 - \frac{1}{1 + e^{-x}} \right) \quad (54)$$

$$\sigma'(x) = \sigma(x) (1 - \sigma(x)) \quad (55)$$

## 1.7 Problem G

### Instructions

Now we shall consider the Negative Sampling loss, which is an alternative to the Naive Softmax loss. Assume that  $K$  negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as  $w_1, w_2, \dots, w_K$ , and their outside vectors as  $\mathbf{u}_{w_1}, \mathbf{u}_{w_2}, \dots, \mathbf{u}_{w_K}$ . For this question, assume that the  $K$  negative samples are distinct. In

other words,  $i \neq j$  implies  $w_i \neq w_j$  for  $i, j \in \{1, \dots, K\}$ . Note that  $o \notin \{w_1, \dots, w_K\}$ . For a center word  $c$  and an outside word  $o$ , the negative sampling loss function is given by:

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (56)$$

for a sample  $w_1, \dots, w_K$ , where  $\sigma(\cdot)$  is the sigmoid function.

### 1.7.1 (i)

#### Instructions

Please repeat parts (b) and (c), computing the partial derivatives of  $\mathbf{J}_{\text{neg-sample}}$  with respect to  $\mathbf{v}_c$ , with respect to  $\mathbf{u}_o$ , and with respect to the  $s^{\text{th}}$  negative sample  $\mathbf{u}_{w_s}$ . Please write your answers in terms of the vectors  $\mathbf{v}_c$ ,  $\mathbf{u}_o$ , and  $\mathbf{u}_{w_s}$ , where  $s \in [1, K]$ . **Note:** you should be able to use your solution to part (f) to help compute the necessary gradients here.

#### Part B Recap:

(b) Compute the partial derivative of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to  $\mathbf{v}_c$ . Please write your answer in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{U}$ . Additionally, answer the following two questions with one sentence each: (1) When is the gradient zero? (2) Why does subtracting this gradient, in the general case when it is nonzero, make  $\mathbf{v}_c$  a more desirable vector (namely, a vector closer to outside word vectors in its window)?

#### Solution

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (57)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = \frac{\partial}{\partial \mathbf{v}_c} -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \frac{\partial}{\partial \mathbf{v}_c} \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (58)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sigma(\mathbf{u}_o^\top \mathbf{v}_c) - \sum_{s=1}^K \frac{\partial}{\partial \mathbf{v}_c} \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (59)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\frac{1}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \frac{\partial}{\partial \mathbf{v}_c} \sigma(\mathbf{u}_o^\top \mathbf{v}_c) - \sum_{s=1}^K \frac{1}{(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))} \frac{\partial}{\partial \mathbf{v}_c} (\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (60)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} [1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{u}_o - \sum_{s=1}^K \frac{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))} [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] (-\mathbf{u}_{w_s}) \quad (61)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] \mathbf{u}_{w_s} \quad (62)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K \left[ 1 - \frac{1}{1 + \exp(-(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))} \right] \mathbf{u}_{w_s} \quad (63)$$



$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K \left[ \frac{1 + \exp(\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{1 + \exp(\mathbf{u}_{w_s}^\top \mathbf{v}_c)} - \frac{1}{1 + \exp(\mathbf{u}_{w_s}^\top \mathbf{v}_c)} \right] \mathbf{u}_{w_s} \quad (64)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K \left[ \frac{\exp(\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{1 + \exp(\mathbf{u}_{w_s}^\top \mathbf{v}_c)} \right] \mathbf{u}_{w_s} \quad (65)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K \left[ \frac{1}{1 + \exp(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)} \right] \mathbf{u}_{w_s} \quad (66)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K \sigma(\mathbf{u}_{w_s}^\top \mathbf{v}_c) \mathbf{u}_{w_s} \quad (67)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{u}_o + \sum_{s=1}^K \sigma(\mathbf{u}_{w_s}^\top \mathbf{v}_c) \mathbf{u}_{w_s} \quad (68)$$

You can see that the gradient is zero when  $[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{u}_o = \sum_{s=1}^K \sigma(\mathbf{u}_{w_s}^\top \mathbf{v}_c) \mathbf{u}_{w_s}$ . For our first term, if  $\mathbf{u}_o^\top \mathbf{v}_c$  is a large negative number (dissimilar) then  $1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)$  evaluates to be close to 1. Thus, in our gradient descent update, we'll be adding  $\approx \mathbf{u}_o$  to  $\mathbf{v}_c$  (recall we're doing gradient descent so  $-(-[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{u}_o) = [1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{u}_o \approx \mathbf{u}_o$ ). Intuitively this makes sense because adding  $\mathbf{u}_o$  to  $\mathbf{v}_c$  will make  $\mathbf{v}_c$  more like  $\mathbf{u}_o$ , which is what we want.

If  $\mathbf{u}_o^\top \mathbf{v}_c$  is a large positive number (similar) then  $1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)$  evaluates to a small number and our gradient update for the first term adds approximately 0. Again, this is intuitive because if the vectors are similar already then we don't need to change  $\mathbf{v}_c$  much.

For our second term, a similar line of reasoning can be invoked to see that we're subtracting out  $\approx \mathbf{u}_{w_s}$  from  $\mathbf{v}_c$  if  $\sigma(\mathbf{u}_{w_s}^\top \mathbf{v}_c)$  is large positive (vectors are similar) and doing nothing if  $\sigma(\mathbf{u}_{w_s}^\top \mathbf{v}_c)$  is large negative (vectors dissimilar). Again, this makes sense because we want  $\mathbf{v}_c$  to be less like the vectors that we negative sample and don't need to change  $\mathbf{v}_c$  if that's already the case.

(c) Compute the partial derivatives of  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$  with respect to each of the 'outside' word vectors,  $\mathbf{u}_w$ 's. There will be two cases: when  $w = o$ , the true 'outside' word vector, and  $w \neq o$ , for all other words. Please write your answer in terms of  $\mathbf{y}$ ,  $\hat{\mathbf{y}}$ , and  $\mathbf{v}_c$ . In this subpart, you may use specific elements within these terms as well (such as  $\mathbf{y}_1, \mathbf{y}_2, \dots$ ). Note that  $\mathbf{u}_w$  is a vector while  $\mathbf{y}_1, \mathbf{y}_2, \dots$  are scalars.

**Case 1:**  $\frac{\partial}{\partial \mathbf{u}_o}$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = \frac{\partial}{\partial \mathbf{u}_o} - \log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \frac{\partial}{\partial \mathbf{u}_o} \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (69)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} [1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{v}_c \quad (70)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = -[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)] \mathbf{v}_c \quad (71)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{v}_c \quad (72)$$

**Case 2:**  $\frac{\partial}{\partial \mathbf{u}_{w_s}}$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} = \frac{\partial}{\partial \mathbf{u}_{w_s}} - \log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (73)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} = -\frac{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)}{\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)} [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] (-\mathbf{v}_c) \quad (74)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} = [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] \mathbf{v}_c \quad (75)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} = \sigma(\mathbf{u}_{w_s}^\top \mathbf{v}_c) \mathbf{v}_c \quad (76)$$

See part (b) of this problem for that last jump, it's the exact same as before.

### 1.7.2 (ii)

#### Instructions

In lecture, we learned that an efficient implementation of backpropagation leverages the reuse of previously-computed partial derivatives. Which quantity could you reuse between the three partial derivatives to minimize duplicate computation? Write your answer in terms of  $\mathbf{U}_{o,\{w_1, \dots, w_K\}} = [\mathbf{u}_o, -\mathbf{u}_{w_1}, \dots, -\mathbf{u}_{w_K}]$ , a matrix with the outside vectors stacked as columns, and  $\mathbf{1}$ , a  $(K+1) \times 1$  vector of 1's.

#### Solution

Since we compute the sigmoid function so much we could reuse the following:

$$\sigma(\mathbf{U}_o \mathbf{v}_c)$$

### 1.7.3 (iii)

#### Instructions

Describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss.

Caveat: So far we have looked at re-using quantities and approximating softmax with sampling for faster gradient descent. Do note that some of these optimizations might not be necessary on modern GPUs and are, to some extent, artifacts of the limited compute resources available at the time when these algorithms were developed.

#### Solution

This loss function is much more efficient because, provided we have enough negative samples, the negative sampling we use will approximate the gradient update we would've gotten for

the entire vocabulary but with literally a fraction of the vocabulary used. Instead of iterating through the entire vocabulary we can just iterate through negative samples which can be orders of magnitude smaller in size.

## 1.8 Problem H

### Instructions

Now we will repeat the previous exercise, but without the assumption that the  $K$  sampled words are distinct. Assume that  $K$  negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as  $w_1, w_2, \dots, w_K$  and their outside vectors as  $\mathbf{u}_{w_1}, \dots, \mathbf{u}_{w_K}$ . In this question, you may not assume that the words are distinct. In other words,  $w_i = w_j$  may be true when  $i \neq j$  is true. Note that  $o \notin \{w_1, \dots, w_K\}$ . For a center word  $c$  and an outside word  $o$ , the negative sampling loss function is given by:

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (77)$$

for a sample  $w_1, \dots, w_K$ , where  $\sigma(\cdot)$  is the sigmoid function.

Compute the partial derivative of  $\mathbf{J}_{\text{neg-sample}}$  with respect to a negative sample  $\mathbf{u}_{w_s}$ . Please write your answers in terms of the vectors  $\mathbf{v}_c$  and  $\mathbf{u}_{w_s}$ , where  $s \in [1, K]$ . Hint: break up the sum in the loss function into two sums: a sum over all sampled words equal to  $w_s$  and a sum over all sampled words not equal to  $w_s$ . Notation-wise, you may write ‘equal’ and ‘not equal’ conditions below the summation symbols.

### Solution

I’m going to use  $j$  to iterate over the sum in this equation instead of  $s$  to make the notation more clear.

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} &= \frac{\partial}{\partial \mathbf{u}_{w_s}} -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{j=1}^K \log(\sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c)) \\ \frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} &= -\frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{j=1, w_s \neq w_j}^K \log(\sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c)) - \frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{j=1, w_s = w_j}^K \log(\sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c)) \\ \frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} &= -\frac{\partial}{\partial \mathbf{u}_{w_s}} \sum_{j=1, w_s = w_j}^K \log(\sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c)) \\ \frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} &= -\sum_{j=1, w_s = w_j}^K \frac{\sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c)}{\sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c)} [1 - \sigma(-\mathbf{u}_{w_j}^\top \mathbf{v}_c)] \mathbf{v}_c \\ \frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} &= \sum_{j=1, w_s = w_j}^K \sigma(\mathbf{u}_{w_j}^\top \mathbf{v}_c) \mathbf{v}_c \end{aligned}$$

## 1.9 Problem I

### Instructions

Suppose the center word is  $c = w_t$  and the context window is  $[w_{t-m}, \dots, w_{t-1}, w_t, w_{t+1}, \dots, w_{t+m}]$ , where  $m$  is the context window size. Recall that for the skip-gram version of word2vec, the total loss for the context window is:

$$\mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U}) \quad (78)$$

Here,  $\mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$  represents an arbitrary loss term for the center word  $c = w_t$  and outside word  $w_{t+j}$ .  $\mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$  could be  $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$  or  $\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ , depending on your implementation.

Write down three partial derivatives:

- (i)  $\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{U}}$
- (ii)  $\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_c}$
- (iii)  $\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_w}$  when  $w \neq c$

Write your answers in terms of  $\frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$  and  $\frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$ . This is very simple – each solution should be one line.

### Solution

At first I thought these solutions were supposed to be substituted and simplified but I'm pretty sure this is literally all the question is asking.

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{U}} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}} \quad (79)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c} \quad (80)$$

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_w} = 0 \quad (81)$$

This last one should make sense because  $\mathbf{v}_w$  isn't part of the loss function at all.

## 2 Programming

If you don't know how to use NumPy (or Python more generally) I'd suggest going through this tutorial to learn the basics

### 2.1 Sigmoid

This one is a gimme if you know how to use NumPy at all. Literally just write down the sigmoid function using `np.exp()` and you're done.

## 2.2 Naive Softmax Loss and Gradient

Here is where all the math we've done really comes in handy. Literally 90% of this assignment is understanding and deriving the equations while the other 10% is just writing down our mathematical results for the computer to calculate. Now, let's get to programming. As we saw in Problem A, our loss function is

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log\left(\frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}\right)$$

Let's rearrange this a bit and make it more numerically stable. These exponentials can easily blow up to numbers much larger than we can store in a computer. We first add a constant  $\mathbf{C}$  to both the numerator and the denominator, which you can see doesn't affect the equation at all.

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log\left(\frac{\mathbf{C} \exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \mathbf{C} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}\right)$$

Then we move the constant inside the exponential.

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log\left(\frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c + \log \mathbf{C})}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c + \log \mathbf{C})}\right)$$

We can then set  $\log \mathbf{C} = -\max(\mathbf{u}_w^\top \mathbf{v}_c)$  so that the maximum value of our exponential will be 1. This prevents the exponential from blowing up the loss such that it makes our code numerically unstable. For more info see these notes from CS231N.

In my implementation, I first found the softmax probabilities ( $\hat{\mathbf{y}}$ ) for every vector in  $\mathbf{U}$  since we'll need it later for the gradient. Speaking of which, here are our gradients computed in problems B, C, and D.

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = \mathbf{U}^\top (\hat{\mathbf{y}} - \mathbf{y})$$

$$\frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{v}_c^\top$$

I chose not to use a one-hot vector  $\mathbf{y}$  in my code because it's quite simple to just subtract 1 from the value  $\hat{y}_o$  in our  $\hat{\mathbf{y}}$  vector. A one-hot vector would just be a waste of space and excess computation. Hopefully my code is easy to follow. It's labeled exactly the same as the math so if you know how to use NumPy it should be easy to follow along.

## 2.3 Negative Sampling Loss and Gradient

Onto the more complicated negative sampling loss and gradient. Our loss from Problem G is

$$\mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))$$

Our first step is to get our negative samples. This is implemented for us already (yay!? I'm always curious as to how things work but I'll leave my curiosity for more important things). Then, we get our matrix  $\mathbf{U}_{o, \{w_1, \dots, w_K\}} = [\mathbf{u}_o, -\mathbf{u}_{w_1}, \dots, -\mathbf{u}_{w_K}]$ . We matrix multiply  $\mathbf{U}_o \mathbf{v}_c$ ,

pass the result through our sigmoid, take the negative log of every term, and sum everything to find our loss. And boom, loss calculated.

For the gradient with respect to  $\mathbf{v}_c$  we have

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] \mathbf{u}_{w_s}$$

Note that this isn't the final form I placed the gradient in but this is actually the easiest one to work with. We can rearrange this a bit to find

$$\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{u}_o + \sum_{s=1}^K [\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c - 1)] (-\mathbf{u}_{w_s})$$

We can easily reuse  $\mathbf{U}_o$  and the sigmoid we calculated before using this simple line of code

```
gradCenterVec = (sigmoid_vec - 1) @ Uo
```

I'll leave it to you to figure out why this works. Onto our final gradient, the gradient with respect to  $\mathbf{u}_o$  and  $\mathbf{u}_{w_s}$ . From problem G part (ii) we know

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} &= [\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1] \mathbf{v}_c \\ \frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_{w_s}} &= [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] \mathbf{v}_c \end{aligned}$$

We use the following lines of code to find  $\frac{\partial \mathbf{J}_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}_o}$

```
sigmoidVec[0] -= 1
sigmoidVec[1:] = 1 - sigmoidVec[1:]
gradUo = np.outer(sigmoidVec, centerWordVec)
```

Finally, we iterate through  $\mathbf{U}_o$  and add the vector to its proper place in our gradient matrix. I racked my brain for a long time to try and do this without using a for loop but I could not find one. Please let me know if you did, I would thoroughly enjoy knowing.

## 2.4 Skipgram

This function is pretty simple. All you need to do is iterate over all the words, add all the losses, add all the gradients, and spit back the output. Nothing much going on here to be honest. Although I will note that in the gradient checker, if you switch the order in which you check the two different functions (negative sampling vs. softmax), you actually get a different loss. I'm not totally sure what that's about but it was something worth mentioning.

## 2.5 Stochastic Gradient Descent

The update rule is pretty simple here. We just use the following code per usual

```
loss, gradient = f(x)
x -= step * gradient
```

Anyways, that's all for this assignment! My implementation ran in 30 minutes on Macbook with M1 chip for reference.

### 3 References

1. Python Tutorial
2. Softmax