

Author: Benjamin Smidt
Created: October 3rd, 2022
Last Updated: October 3rd, 2022

CS 224N A2: Word Vectors

Note to the reader. This is my work for assignment one of Stanford's course CS 224N: Natural Language Processing with Deep Learning. You can find the lecture Winter 2021 lectures series on YouTube here. This document is meant to be used as a reference, explanation, and resource for the assignment, not necessarily a comprehensive overview of Word Vectors. If there's a typo or a correction needs to be made, feel free to email me at benjamin.smidt@utexas.edu so I can fix it. Thank you! I hope you find this document helpful :).

Contents

1	Written Understanding	1
1.1	Problem 1-A	1

1 Written Understanding

1.1 Problem 1-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

To start with, our naive-softmax loss is defined as

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

where

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

Let's define our variables. $\hat{\mathbf{y}}$ is our score vector. Note that the numerator is a vector of length V while the denominator is a scalar (this notation is a bit abusive but I think it actually makes things clearer). Thus, each index in the vector can be interpreted as the probability that the corresponding word (using the the index and one hot vector) is the center word.

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

\mathbf{y} is the one hot vector of the true center word. Then

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

Where the index o indicates the index containing the only 1 within $\hat{\mathbf{y}}$ (since it is a one-hot vector).

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

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

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

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

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

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 1-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