

Advanced Methods in NLP - Assignment 1: Word Vectors**Question 1(a)**

We remember that:

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

We note that:

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

And so we get:

$$\text{softmax}(\mathbf{x}) = \text{softmax}(\mathbf{x} + c)$$

We can use this property to optimize *softmax* by choosing $c = -\max_i x_i$.

Question 1(c)

The sigmoid function is: $\sigma(x) = \frac{1}{1 + e^{-x}}$

Let's consider: $f(x) = \frac{1}{\sigma(x)} = 1 + e^{-x}$.

On the one hand, $\frac{d}{dx}f(x) = \frac{d}{dx}\sigma(x)^{-1} = -\frac{d}{dx}\sigma(x) \cdot \sigma(x)^{-2}$

On the other hand, $\frac{d}{dx}f(x) = \frac{d}{dx}(1 + e^{-x}) = -e^{-x} = 1 - f(x) = 1 - \frac{1}{\sigma(x)} = \frac{\sigma(x) - 1}{\sigma(x)}$

So we get that:

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

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

So the gradients of the sigmoid function can be written as a function of the function value:

$$\frac{d}{dx}\sigma(x) = \sigma(x) \cdot (1 - \sigma(x))$$

Question 2(a)

$$\hat{y}_o = p(o | c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w=1}^W \exp(\mathbf{u}_w^\top \mathbf{v}_c)}$$

$$J_{softmax-CE}(o, \mathbf{v}_c | U) = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \cdot \log(\hat{y}_i)$$

$$\frac{\partial J}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_{w=1}^W \hat{y}_w \cdot \mathbf{u}_w$$

Question 2(b)

$$\frac{\partial J}{\partial \mathbf{u}_w} = \begin{cases} (\hat{y}_w - 1) \mathbf{v}_c, & w = o \\ \hat{y}_w \mathbf{v}_c, & \text{otherwise} \end{cases}$$

Question 2(c)

The sigmoid function is: $\sigma(x) = \frac{1}{1 + e^{-x}}$

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

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

$$\frac{\partial J}{\partial \mathbf{u}_o} = (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{v}_c$$

$$\text{for } k \in [K]: \frac{\partial J}{\partial \mathbf{u}_k} = -(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1) \mathbf{v}_c$$

Question 2(d)

If we use $F(o, \mathbf{v}_c)$ for $J_{softmax-CE}(o, \mathbf{v}_c | \dots)$ or $J_{neg-sample}(o, \mathbf{v}_c | \dots)$, the cost for a context centered around c is:

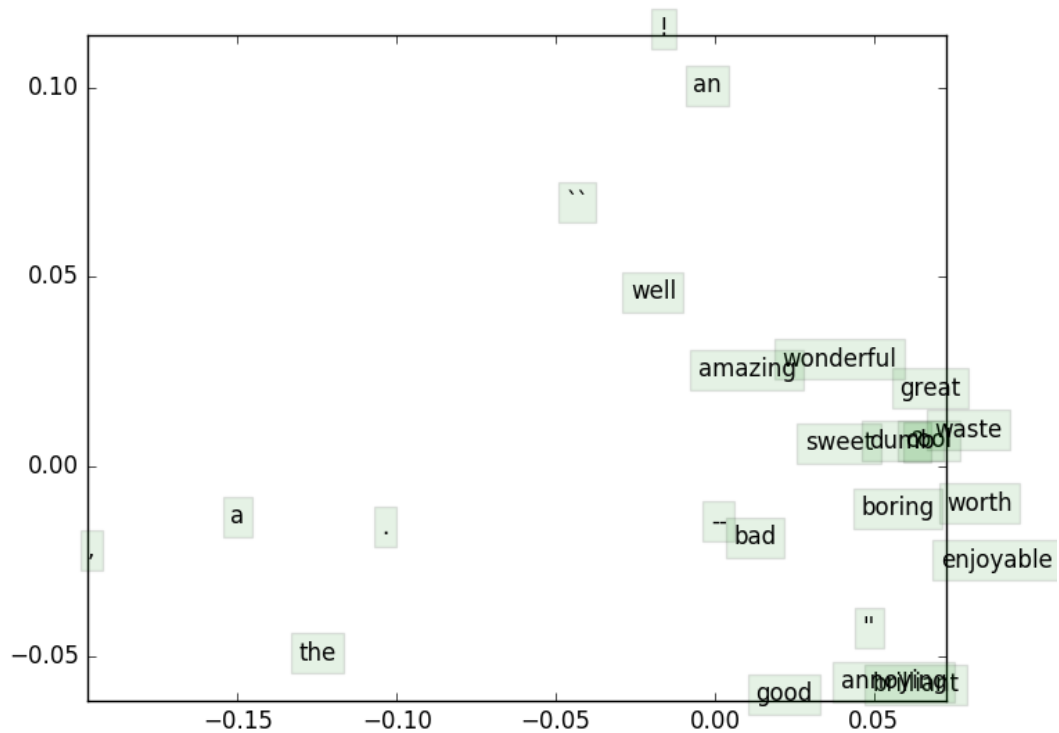
$$J = \sum_{-m \leq j \leq m, j \neq 0} F(\mathbf{w}_{c+j}, \mathbf{v}_c)$$

We denote the set of all output vectors \mathbf{u}_k as \mathbf{U} .

$$\frac{\partial J}{\partial \mathbf{U}} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{U}}$$

and $\frac{\partial J}{\partial v_j} = 0$ if $j \neq c$.

Question 2(h)



KNN:

Words related to "the":

['comedy\\thriller', 'a', 'if', 'that', 'or', 'the', '.', 'is', 'derek', 'bolt', 'decide']

Words related to "unique":

['unique', '1979', 'dares', 'realized', 'succumb', 'chabrolian', 'ba', 'puns', 'regardless',
'imaginative', 'lunar']

Words related to "superb":

['best', 'pool', 'industry', 'mine', 'gold', 'ghoulish', 'zingers', 'superb', 'moppets', 'roussillon', 'transporter']

Words related to "comedy":

['cute', 'longest', 'comedy', 'sensation', 'fast', 'cleaving', 'mature', 'first-timer', 'singing', 'observation', 'often-funny']

Words related to "surprisingly":

['hundred', 'surprisingly', 'philandering', 'either', 'protective', '20-car', 'unusually', 'bollywood', 'dogs', 'thinking', 'soderbergh']

- Notice that we implemented knn without full sorting (using argpartition for performance) therefore the related words list order is not by closeness. We can see that a lot of syntactic related words in the related words list and sometimes we see even semantic related words ('comedy' – 'often funny', 'superb' – 'best' and so on).
- We can see that in the image there are some close words with similar meaning such as 'wonderful' and 'great', and close words with similar connotation such as 'bad', 'boring' and 'worth' (negative connotation). Still we can observe that words which are not semantically related such as 'waste' and 'great' are still close in the image, but they are indeed syntactically close (both can be interpreted as adjectives). So, bottom line word2vec words embedding shows syntactic and semantic closeness for some words, but not for all of them (might perform better with larger corpus).