

CS 224n Assignment 3.

Abhishek Goswami.

January 24, 2019

1. Written : Understanding word2vec

(a)

$$\sum_{w \in Vocab} y_w \log(\hat{y}_w) = -\log(\hat{y}_o). \quad (1)$$

because

$$y_i = \begin{cases} 1, & \text{if } i = o \\ 0, & \text{otherwise} \end{cases}. \quad (2)$$

(b)

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = -\mathbf{u}_o + \sum_{w=1}^V \hat{y}_w \mathbf{u}_w. \quad (3)$$

also equivalent to

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}). \quad (4)$$

(c)

$$\frac{\partial J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{U}} = \mathbf{v}_c(\hat{\mathbf{y}} - \mathbf{y})^\top. \quad (5)$$

(d)

$$\sigma'(x) = \sigma(x)(1 - \sigma(x)) \quad (6)$$

(e)

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} = (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1)\mathbf{u}_o - \sum_{k=1}^K (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1)\mathbf{u}_k \quad (7)$$

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_o} = (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1)\mathbf{v}_c. \quad (8)$$

$$\frac{\partial J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{u}_k} = -(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1)\mathbf{v}_c. \quad (9)$$

$J_{\text{neg-sample}}$ is more efficient to compute than $J_{\text{naive-softmax}}$. In $J_{\text{naive-softmax}}$ we use softmax to compute the probability of outside word given the center word. In order to compute the probability, we need to normalize over all the words in the vocabulary. This is not efficient especially when the vocabulary size is large.

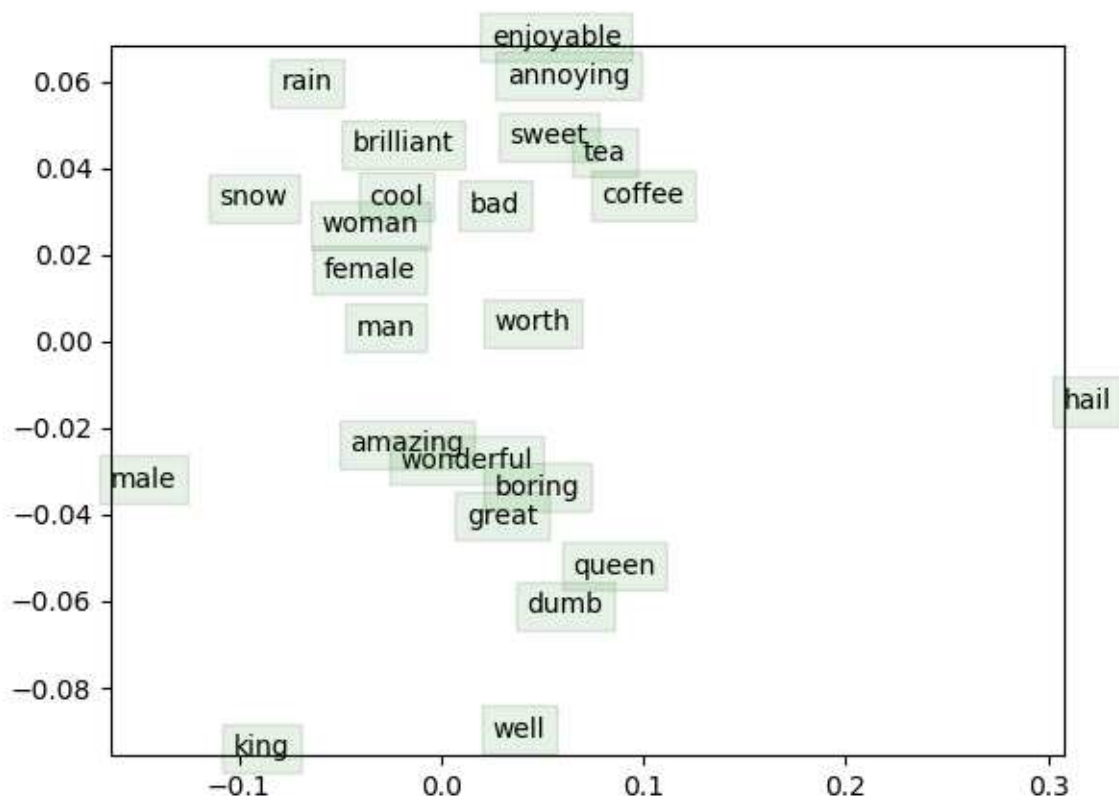


Figure 1: Word Vectors

(f) i.
$$\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots w_{t+m}, \mathbf{U})}{\partial \mathbf{U}} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}. \quad (10)$$

$$\text{ii.} \quad \frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c} \quad (11)$$

$$\text{iii.} \quad \frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_w} = 0, \text{ when } w \neq c \quad (12)$$

2. Coding

In Figure 1 we make the following observations :

- (a) Semantic similarity e.g. (king, male) (queen, female)
- (b) Syntactic structure e.g. (man, woman) (male, female)
- (c) Beverages e.g. (tea, coffee) appear together