

3: word2vec

(a) Problem: Given a predicted word vector v_c for center word c for skip-gram and softmax prediction, find the derivative of the cross entropy cost wrt v_c

Answer: The soft-max prediction gives:

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

Cross entropy loss is given by

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = J = - \sum_{i=1}^V y_i \log \hat{y}_i \quad (2)$$

where $\hat{\mathbf{y}}$ is a $1 \times V$ matrix/vector whose components are the softmax results, and \mathbf{y} is a $1 \times V$ matrix/vector that is 1-hot encoded with 1 at the o th entry. So, we can write $\hat{\mathbf{y}}$ as

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z}) \quad (3)$$

where \mathbf{z} comprises $z_i = \mathbf{u}_i^T \mathbf{v}_c$.

We want to calculate, by chain rule

$$\frac{\partial J}{\partial \mathbf{v}_c} = \frac{\partial J}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{v}_c} \quad (4)$$

where we are multiplying a $1 \times V$ matrix with a $V \times V$ matrix.

From previous results we have

$$\frac{\partial J}{\partial \mathbf{z}} = \hat{\mathbf{y}} - \mathbf{y} \quad (5)$$

If we consider \mathbf{z} a column vector, and adopt numerator layout for derivatives we get

$$\frac{\partial \mathbf{z}}{\partial \mathbf{v}_c} = [\mathbf{u}_1 \quad \dots \quad \mathbf{u}_V] = \mathbf{U} \quad (6)$$

This gives

$$\frac{\partial J}{\partial \mathbf{v}_c} = (\hat{\mathbf{y}} - \mathbf{y})\mathbf{U} \quad (7)$$

(b) Problem: Find the partial derivative of the above cost function wrt \mathbf{u}_k

Answer: We need to calculate

$$\frac{\partial J}{\partial \mathbf{u}_k} = \frac{\partial J}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{u}_k} \quad (8)$$

We use results from the preceding sub-problem for $\frac{\partial J}{\partial \mathbf{z}}$ and note that

$$\frac{\partial z_i}{\partial \mathbf{u}_k} = \mathbf{v}_c \delta_{ik} \quad (9)$$

Because of the Kronecker delta we can write

$$\frac{\partial J}{\partial \mathbf{u}_k} = \frac{\partial J}{\partial z_i} \frac{\partial z_i}{\partial \mathbf{u}_k} \quad (10)$$

This gives

$$\frac{\partial J}{\partial \mathbf{u}_k} = (\hat{\mathbf{y}} - \mathbf{y})_k \mathbf{v}_c \quad (11)$$

(c) Problem: Repeating **(a)** for a negative sampling cost function

Answer: Given the cost function

$$J = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \quad (12)$$

we use the chain rule

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{v}_c} &= -\frac{\mathbf{u}_o^T \sigma(\mathbf{u}_o^T \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c))}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} - \sum_{k=1}^K \frac{\mathbf{u}_k^T \sigma(\mathbf{u}_k^T \mathbf{v}_c)(1 - \sigma(\mathbf{u}_k^T \mathbf{v}_c))}{\sigma(\mathbf{u}_k^T \mathbf{v}_c)} \\ &= -\mathbf{u}_o^T (1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \mathbf{u}_k^T (1 - \sigma(\mathbf{u}_k^T \mathbf{v}_c)) \end{aligned} \quad (13)$$

Similarly

$$\frac{\partial J}{\partial \mathbf{u}_o} = (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1) \mathbf{v}_c^T \quad (14)$$

and

$$\frac{\partial J}{\partial \mathbf{u}_k} = (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{v}_c^T \quad (15)$$

for $o \neq k$