

# CS 224n Problem Set #2 Solutions: word2vec

LA DUC CHINH

---

**Due Wednesday, Nov 14 at 11:59 pm on Gradescope.**

**Written: Understanding word2vec**

(a) We have

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

(b) We have

$$\frac{\partial J_{naive-softmax}(v_c, o, U)}{\partial v_c} = -u_o + \frac{1}{\sum_{w \in Vocab} \exp(u_w^T v_c)} \sum_{w \in Vocab} \exp(u_w^T v_c) u_w \quad (2)$$

$$= -u_o + \sum_{w \in Vocab} \frac{\exp(u_w^T v_c)}{\sum_{w \in Vocab} \exp(u_w^T v_c)} u_w \quad (3)$$

$$= -u_o + \sum_{w \in Vocab} P(O = w | C = c) u_w \quad (4)$$

$$= -u_o + \sum_{w \in Vocab} \hat{y}_w u_w \quad (5)$$

Let  $\hat{y}_w = y_w + (\hat{y}_w - y_w)$  and we have  $\sum_{w \in Vocab} y_w u_w = \sum_{w \in Vocab} 1\{w = o\} u_w u_o$

$$\frac{\partial J_{naive-softmax}(v_c, o, U)}{\partial v_c} = -u_o + \sum_{w \in Vocab} \hat{y}_w u_w \quad (6)$$

$$= -u_o + \sum_{w \in Vocab} [y_w + (\hat{y}_w - y_w)] u_w \quad (7)$$

$$= -u_o + \sum_{w \in Vocab} y_w u_w + \sum_{w \in Vocab} (\hat{y}_w - y_w) u_w \quad (8)$$

$$= -u_o + u_o + \sum_{w \in Vocab} (\hat{y}_w - y_w) u_w \quad (9)$$

$$= \sum_{w \in Vocab} (\hat{y}_w - y_w) u_w \quad (10)$$

$$= \sum_{j=1}^{|V|} (\hat{y}_j - y_j) u_j = (\hat{y}_1 - y_1) u_1 + (\hat{y}_2 - y_2) u_2 + \dots + (\hat{y}_{|V|} - y_{|V|}) u_{|V|} \quad (11)$$

$$= (\hat{y}_1 - y_1 \quad \hat{y}_2 - y_2 \quad \dots \quad \hat{y}_{|V|} - y_{|V|}) \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_{|V|} \end{pmatrix} \quad (12)$$

$$= (\hat{y} - y) U \quad (13)$$

(c) partial derivatives outside word when  $w = o$   
 Similar to (b) we have

$$\frac{\partial J}{\partial \mathbf{u}_o} = -(1 - \hat{y}) \mathbf{v}_c = (\hat{y}_{w=o} - y_{y=o}) \mathbf{v}_c \quad (14)$$

When  $w \neq o$

$$\frac{\partial J}{\partial \mathbf{u}_w} = -\frac{\partial \log \hat{y}}{\partial \mathbf{u}_w} \quad (15)$$

$$= -\frac{\partial \log \hat{y}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{u}_w^T \mathbf{v}_c} \frac{\partial \mathbf{u}_w^T \mathbf{v}_c}{\partial \mathbf{u}_w} \quad (16)$$

$$= -\frac{1 - \exp(\mathbf{u}_o^T \mathbf{v}_c) \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\hat{y} (\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c))^2} \mathbf{v}_c \quad (17)$$

$$= \frac{1}{\hat{y}} \hat{y}_{w, w \neq o} \mathbf{v}_c \quad (18)$$

$$= \hat{y}_{w, w \neq o} \mathbf{v}_c \quad (19)$$

$$= (\hat{y}_{w, w \neq o} - y_{w, w \neq o}) \mathbf{v}_c \quad (20)$$

Therefore we have

$$\frac{\partial J}{\partial \mathbf{u}_w} = (\hat{y}_w - y_w) \mathbf{v}_c \quad (21)$$

(d) Sigmoid

$$\frac{d\sigma}{dx} = \sigma(1 - \sigma) \quad (22)$$

(e) Negative Sampling loss

$$\mathbf{J}_{neg-sample}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \quad (23)$$

Derivative  $\mathbf{v}_c$ 

$$\frac{\partial \mathbf{J}}{\partial \mathbf{v}_c} = -\frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \sigma(\mathbf{u}_o^T \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o - \sum_{k=1}^K \frac{1}{\sigma(-\mathbf{u}_k^T \mathbf{v}_c)} \sigma(-\mathbf{u}_k^T \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c))(-\mathbf{u}_k) \quad (24)$$

$$= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{u}_o + \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{u}_k \quad (25)$$

Derivative  $\mathbf{u}_o$ 

$$\frac{\partial \mathbf{J}}{\partial \mathbf{u}_o} = -\frac{1}{\sigma(\mathbf{u}_o^T \mathbf{v}_c)} \sigma(\mathbf{u}_o^T \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{v}_c \quad (26)$$

$$= -(1 - \sigma(\mathbf{u}_o^T \mathbf{v}_c)) \mathbf{v}_c \quad (27)$$

Derivative  $\mathbf{u}_k$ 

$$\frac{\partial \mathbf{J}}{\partial \mathbf{u}_k} = \sum_{k=1}^K -\frac{1}{\sigma(-\mathbf{u}_k^T \mathbf{v}_c)} \sigma(-\mathbf{u}_k^T \mathbf{v}_c)(1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c))(-\mathbf{v}_c) \quad (28)$$

$$= \sum_{k=1}^K (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \mathbf{v}_c \quad (29)$$

This loss function is more efficient because sigmoid function has less computation cost than softmax function.

(f) Skip-gram

$$\frac{\partial J_{\text{skip-gram}}}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U} \quad (30)$$

$$\frac{\partial J_{\text{skip-gram}}}{\partial v_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial v_c} \quad (31)$$

$$\frac{\partial J_{\text{skip-gram}}}{\partial v_w} = 0 \quad (32)$$