

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

This happens because the true empirical distribution y is a one-hot vector with a 1 for the true outside word o , and 0 (zero) everywhere else.

$$\begin{aligned} \text{b) } \frac{\partial \text{J}_{\text{naive-softmax}}(v_c, o, U)}{\partial v_c} &= \frac{\partial}{\partial v_c} \left(-\log \left(\frac{\exp(u_o^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \right) \right) = \\ &= \frac{\partial}{\partial v_c} \left(-\log(\exp(u_o^T v_c)) + \log \left(\sum_{w \in \text{Vocab}} \exp(u_w^T v_c) \right) \right) = \\ &= \frac{\partial}{\partial v_c} \left(-u_o^T v_c + \log \left(\sum_{w \in \text{Vocab}} \exp(u_w^T v_c) \right) \right) = -u_o + \\ &+ \frac{1}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \times \frac{\partial}{\partial v_c} \left(\sum_{w \in \text{Vocab}} \exp(u_w^T v_c) \right) = \\ &= -u_o + \frac{1}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \times \left(\sum_{w \in \text{Vocab}} \exp(u_w^T v_c) u_w \right) = \\ &= \sum_{w \in \text{Vocab}} \frac{\exp(u_w^T v_c) u_w}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} - u_o \quad \ominus \text{ since we know that} \end{aligned}$$

$$\begin{aligned} \hat{y}_o &= \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{Vocab}} \exp(u_w^T v_c)} \quad \text{and} \quad U y_o = u_o \quad \ominus \left(\sum_{w \in \text{Vocab}} \hat{y}_w u_w \right) - U y_o = \\ &= U \hat{y} - U y = U(\hat{y} - y) \quad [\text{assuming that the } y \text{ is for the outside word } o] \end{aligned}$$

1) The gradient is zero when the model perfectly classifies the context word, meaning that it assigns probability of 1 to it.

where o is any word in the Vocab

2) Subtracting the gradient of the loss function with respect to the v_c vector from the v_c vector moves that vector in the direction of the negative gradient, which is the direction of the steepest descent, and thus maximally minimizes the loss function.

$$c) \frac{\partial J_{\text{neural-softmax}}(v_c, 0, U)}{\partial u_w} = \frac{\partial}{\partial u_w} \left(-\log \left(\frac{\exp(u_0^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \right) \right) =$$

$$\Leftrightarrow \text{based on b)} \Leftrightarrow \frac{\partial}{\partial u_w} \left(-u_0^T v_c + \log \left(\sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right) \right)$$

when $w=0$,

$$\frac{\partial}{\partial u_0} \left(-u_0^T v_c + \log \left(\sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right) \right) = \left(-v_c + \frac{1}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \right) \times$$

$$\times \left(\exp(u_0^T v_c) v_c \right) = \frac{\exp(u_0^T v_c) v_c}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} - v_c =$$

$$= \hat{y}_0 v_c - v_c = (\hat{y}_0 - 1) v_c$$

when $w \neq 0$,

$$\frac{\partial}{\partial u_w} \left(-u_0^T v_c + \log \left(\sum_{w \in \text{vocab}} \exp(u_w^T v_c) \right) \right) = \frac{1}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} \times$$

$$\times \left(\exp(u_w^T v_c) v_c \right) = \frac{\exp(u_w^T v_c) v_c}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} = \hat{y}_w v_c$$

d) In c), we calculated the derivatives of $J_{\text{naive-softmax}}(v_c, o, U)$ with respect to each of the 'outside' word vectors, u_w 's.

Since those u_w vectors represent the columns of the U matrix, we can just reuse our answers from c) here.

$$\text{So, } \frac{\partial J_{\text{naive-softmax}}(v_c, o, U)}{\partial U} = [\hat{y}_1 v_c \quad \hat{y}_2 v_c \quad \dots \quad (\hat{y}_o - 1) v_c \quad \dots \quad \hat{y}_{|\text{vocab}|} v_c]$$

$$\oplus v_c(\hat{y} - y)^T$$

$$\in \mathbb{R}^{d \times |\text{vocab}|}$$

where d is the dimension of the embeddings.

$$e) f(x) = \max(0, x)$$

$$\text{when } x < 0, f'(x) = 0$$

$$\text{when } x > 0, f'(x) = 1$$

$$f) \sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}, \quad \sigma'(x) = \frac{e^x(e^x+1) - e^{2x}}{(e^x+1)^2} = \frac{e^x}{(e^x+1)^2} =$$

$$= \frac{e^x}{(e^x+1)^2} \times \frac{1}{e^x+1} = \sigma(x) \times \frac{e^x+1 - e^x}{e^x+1} = \sigma(x) \left(1 - \frac{e^x}{e^x+1}\right) = \sigma(x)(1 - \sigma(x))$$

$$g) \text{ (i) } \frac{\partial J_{\text{neg-sample}}(v_c, o, U)}{\partial v_c} \quad \text{with respect to } v_c$$

$$= \frac{\partial}{\partial v_c} \left(-\log(\sigma(u_o^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c)) \right) =$$

$$= -\frac{1}{\sigma(u_o^T v_c)} \times (1 - \sigma(u_o^T v_c)) \sigma(u_o^T v_c) u_o + \sum_{s=1}^K \frac{(1 - \sigma(-u_{w_s}^T v_c)) \sigma(-u_{w_s}^T v_c) u_{w_s}}{\sigma(-u_{w_s}^T v_c)}$$

$$= -(1 - \sigma(u_o^T v_c)) u_o + \sum_{s=1}^K (1 - \sigma(-u_{w_s}^T v_c)) u_{w_s}$$

$$\text{with respect to } u_o$$

$$\frac{\partial J_{\text{neg-sample}}(v_c, o, U)}{\partial u_o} = \frac{\partial}{\partial u_o} \left(-\log(\sigma(u_o^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c)) \right) =$$

$$= -\frac{1}{\sigma(u_o^T v_c)} \times (1 - \sigma(u_o^T v_c)) \sigma(u_o^T v_c) v_c = -(1 - \sigma(u_o^T v_c)) v_c$$

With respect to u_{w_s}

$$\begin{aligned} \frac{\partial J_{\text{neg-sample}}(v_c, o, u)}{\partial u_{w_s}} &= \frac{\partial}{\partial u_{w_s}} \left(-\log(\sigma(u_o^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c)) \right) = \\ &= -\frac{1}{\sigma(-u_{w_s}^T v_c)} \times \left((1 - \sigma(-u_{w_s}^T v_c)) \sigma(-u_{w_s}^T v_c) \right) \times (-v_c) = \\ &= (1 - \sigma(-u_{w_s}^T v_c)) v_c \end{aligned}$$

(ii) To minimize duplication, ~~we can reuse~~, we can reuse the quantities $(1 - \sigma(u_o^T v_c))$ and/or $(1 - \sigma(-u_{w_s}^T v_c))$ [$s \in [1, \dots, K]$]

We calculated these quantities to find the value of $\partial J / \partial v_c$. Now,

we can express $\partial J / \partial v_c$ and $\partial J / \partial u$ in the following way.

$$\partial J / \partial v_c = -U_{o, \{w_1, \dots, w_K\}} (1 - \sigma(U_{o, \{w_1, \dots, w_K\}}^T v_c)) \in \mathbb{R}^{d \times 1}$$

$$\partial J / \partial u = -\text{sing}(U_{o, \{w_1, \dots, w_K\}}) \otimes v_c (1 - \sigma(U_{o, \{w_1, \dots, w_K\}}^T v_c))^T \in \mathbb{R}^{d \times (K+1)}$$

sing function

this is done to resolve the discrepancy

element-by-element multiplication

between the gradients with respect to u_o and u_{w_s} , $w \in [1, \dots, K]$

(iii) What makes the Negative Sampling loss efficient is that it only involves a fixed K number (much smaller than the vocab size) of words in its computation in addition with the true outsize word, unlike soft max loss.

$$h) \frac{\partial J_{\text{neg-sample}}(v_c, o, u)}{\partial u_{w_s}} = \frac{\partial}{\partial u_{w_s}} \left(-\log(\sigma(u_o^T v_c)) - \sum_{s=1}^K \log(\sigma(-u_{w_s}^T v_c)) \right)$$

$$= -\frac{\partial}{\partial u_{w_s}} \left(\sum_{\substack{a \in [1, K] \\ a \neq s}} \log(\sigma(-u_{w_a}^T v_c)) + \sum_{\substack{a \in [1, K] \\ a = s}} \log(\sigma(-u_{w_a}^T v_c)) \right) =$$

$$= \sum_{\substack{a \in [1, K] \\ a = s}} (1 - \sigma(-u_{w_a}^T v_c)) v_c$$

$$(i) \frac{\partial j_{\text{skip-gram}}(v_c, w_{+m}, \dots, w_t, \dots, w_{+m}, U)}{\partial U} =$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial j(v_c, w_{+j}, U)}{\partial U}$$

$$(ii) \frac{\partial j_{\text{skip-gram}}(v_c, w_{+m}, \dots, w_t, \dots, w_{+m}, U)}{\partial v_c} =$$

$$= \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial j(v_c, w_{+j}, U)}{\partial v_c}$$

$$(iii) \frac{\partial j_{\text{skip-gram}}(v_c, w_{+m}, \dots, w_t, \dots, w_{+m}, U)}{\partial w_t} =$$

$$= 0 \text{ (when } w \neq c)$$

- ② c) In the graph, we can see that a number of words, which have similar meanings or have closeness in terms of general topic, ended up being close to each other. Clear examples are (1) woman, female and man; (2) amazing, wonderful and great; (3) tea, coffee, sweet; (4) rain and snow. However, the word "male" did not end up being close to the cluster (1) and the word "hail" was also too far from the words "rain" and "snow", which are expected to be close.