

This handout includes space for every question that requires a written response. Please feel free to use it to handwrite your solutions (legibly, please). If you choose to typeset your solutions, the `README.md` for this assignment includes instructions to regenerate this handout with your typeset \LaTeX solutions.

1.a

U is the outside word matrix composed of u_o outside vectors for each outside word.

V is the center word matrix composed of v_c center vectors for each center word.

Both U and V contains a vector for each word w in the vocabulary.

U : (embedsize X numtokens)

V : (embedsize X numtokens)

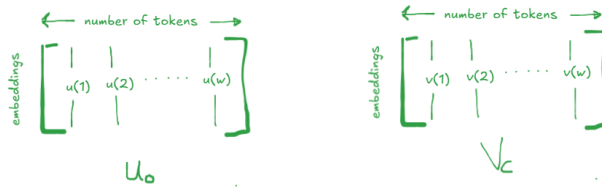


Figure 1: U V matrix

u_o is the column o of U

v_c is the column c of V

Cross-entropy "loss" between the true probability distribution y and predicted distribution \hat{y} is represented as:

$$-\sum_i y \log(\hat{y}) \dots (1)$$

where y and \hat{y} are vectors of length equal to the number of words in the vocabulary.

y : (num-tokens X 1)

\hat{y} : (num-tokens X 1)

Probability of an outside word o given a center word c is the softmax function:

$$\hat{y}_o = P(O = o | C = c) = \frac{\exp(u_o^T v_c)}{\sum_w \exp(u_w^T v_c)} \dots (2)$$

$$J_{\text{native-softmax}}(v_c, o, U) = -\sum_{w \in V_{\text{ocab}}} y_w \log(\hat{y}_o) \dots (3)$$

For a given center word v_c there is only one outside word u_o which is the k^{th} word in the vocabulary, where $y_k = 1$. Hence the $y_i = 0$ where i is not equal to k .

$$J_{\text{native-softmax}}(v_c, o, U) = -y_k \log(y_o(k)) \dots (4)$$

$$J_{\text{native-softmax}}(v_c, o, U) = -\log(\hat{y}_o) \dots (5)$$

Substituting (2) to (5)

$$J_{\text{native-softmax}}(v_c, o, U) = -\log\left(\frac{\exp(u_o^T v_c)}{\sum_w \exp(u_w^T v_c)}\right) \dots (6)$$

$$J_{native-softmax}(v_c, o, U) = -\log(\exp(u_o^T v_c)) + \log(\sum_w \exp(u_w^T v_c)) \dots (7)$$

$$J_{native-softmax}(v_c, o, U) = -u_o^T v_c + \log(\sum_w \exp(u_w^T v_c)) \dots (7)$$

Partial derivative of $J_{native-softmax}$ with respect to v_c

$$\frac{\partial J}{\partial v_c} = -u_o + \frac{\partial(\log(\sum_w \exp(u_w^T v_c)))}{\partial v_c} \dots (8)$$

Derivative of a nested functions (by chain rule):

$$\frac{df(g(x))}{dx} = \frac{d(f(g(x)))}{d(g(x))} \cdot \frac{d(g(x))}{dx}$$

$$\begin{aligned} f(g(v_c)) &= \log(g(v_c)) \\ g(v_c) &= \sum_w \exp(u_w^T v_c) \\ \frac{d(f(g(v_c)))}{d(g(v_c))} &= \frac{1}{g(v_c)} \end{aligned}$$

$$\frac{\partial J}{\partial v_c} = -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \frac{\partial(\sum_w \exp(u_w^T v_c))}{\partial v_c} \dots (9)$$

$$\frac{\partial J}{\partial v_c} = -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_w \exp(u_w^T v_c) \frac{\partial(u_w^T v_c)}{\partial v_c} \dots (10)$$

$$\frac{\partial J}{\partial v_c} = -u_o + \frac{1}{\sum_w \exp(u_w^T v_c)} \cdot \sum_w \exp(u_w^T v_c) u_w^T \dots (11)$$

$$\frac{\partial J}{\partial v_c} = -u_o + \sum_w \frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} \cdot u_w^T \dots (12)$$

Substituting (2)

$$\frac{\partial J}{\partial v_c} = -u_o + \sum_w^V \hat{y}_w \cdot u_w^T \dots (12)$$

$$\frac{\partial J}{\partial v_c} = -u_o + \sum_w^V \hat{y}_w \cdot u_w \dots (13)$$

V is the number of words in vocabulary.

1.b

$$J_{\text{native-softmax}}(v_c, o, U) = -u_o^T v_c + \log(\sum_w \exp(u_w^T v_c)) \dots (1)$$

$$\frac{\partial J}{\partial u_w} = -\frac{\partial u_o^T v_c}{\partial u_w} + \frac{\partial \log(\sum_w \exp(u_w^T v_c))}{\partial u_w} \dots (2)$$

$$\frac{\partial J}{\partial u_w} = -\frac{\partial u_o^T v_c}{\partial u_w} + \frac{1}{\sum_w \exp(u_w^T v_c)} \frac{\partial (\sum_w \exp(u_w^T v_c))}{\partial u_w} \dots (3)$$

$$\frac{\partial J}{\partial u_w} = -\frac{\partial u_o^T v_c}{\partial u_w} + \frac{\sum_w \exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} \frac{\partial (u_w^T v_c)}{\partial u_w} \dots (4)$$

$$\frac{\partial J}{\partial u_w} = -\frac{\partial u_o^T v_c}{\partial u_w} + \frac{\exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)} v_c \dots (5)$$

$$\frac{\partial J}{\partial u_w} = -\frac{\partial u_o^T v_c}{\partial u_w} + \hat{y}_w v_c \dots (6)$$

When $w = o$

$$\frac{\partial J}{\partial u_w} = -v_c + \hat{y}_w v_c = (\hat{y}_w - 1)v_c \dots (7)$$

When $w = \text{otherwise}$

$$\frac{\partial J}{\partial u_w} = \hat{y}_w v_c \dots (8)$$