

Exercise 1. (a) Softmax is invariant to constant offsets in the input, that is, for any input vector x and any constant c ,

$$\begin{aligned}
 \text{softmax}(x_i + c) &= \frac{e^{(x_i + c)}}{\sum_j e^{(x_j + c)}} \\
 &= \frac{e^x * e^c}{\sum_j e^{x_j} e^c} \\
 &= \frac{e^x * e^c}{e^c \sum_j e^{x_j}} \\
 &= \frac{e^{x_i}}{\sum_j e^{x_j}} \\
 &= \text{softmax}(x)_i
 \end{aligned}$$

Exercise 2. (a) Derive the gradients of the sigmoid function and show that it can be rewritten as a function of the function value

$$\begin{aligned}
 \frac{d}{dx} \sigma(x) &= \frac{d}{dx} \left[\frac{1}{1 + e^{-x}} \right] \\
 &= \frac{d}{dx} (1 + e^{-x})^{-1} \\
 &= -(1 + e^{-x})^{-2} (-e^{-x}) \\
 &= \frac{e^{-x}}{(1 + e^{-x})^2} \\
 &= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \\
 &= \frac{1}{1 + e^{-x}} \cdot \frac{(1 + e^{-x}) - 1}{1 + e^{-x}} \\
 &= \frac{1}{1 + e^{-x}} \cdot \left(1 - \frac{1}{1 + e^{-x}} \right) \\
 &= \sigma(x) \cdot (1 - \sigma(x))
 \end{aligned}$$

Exercise 2. (b) Derive the gradient with regard to the inputs of a softmax function when cross entropy loss is used for evaluation

$$\hat{y} - y$$

Exercise 2. (c) Derive the gradients with respect to the inputs x to an one-hidden-layer neural network

$$z_1 = xW_1 + b_1$$

$$z_2 = hW_2 + b_2$$

$$\delta_1 = \frac{\partial C E}{\partial x} = \hat{y} - y$$

$$\delta_2 = \frac{\partial C E}{\partial h} = \delta_1 \frac{\partial z_2}{\partial h} = \delta_1 W_2^T$$

$$\delta_3 = \frac{\partial C E}{\partial z_1} = \delta_2 \frac{\partial h}{\partial z_1}$$

$$\frac{\partial C E}{\partial x} = \delta_3 \frac{\partial z_1}{\partial x} = \delta_3 W_1^T$$

Exercise 2. (d) Parameters in this neural network is $(D \times 1)H + (H + 1)Dy$

Exercise 3. (a) Derive the gradients with respect to vc :

$$\frac{\partial J}{\partial v_c} = U^T (\hat{y} - y)$$

Exercise 3. (b) As in the previous part, derive gradients for the output word vectors uw

$$\frac{\partial J}{\partial U} = v_c (\hat{y} - y)^T$$

Exercise 3. (c) Describe with one sentence why this cost function is much more efficient to compute than the softmax-CE loss

$$\frac{\partial J}{\partial v_c} = (\sigma(u_o^T v_c) - 1) u_o - \sum_{k=1}^K (\sigma(-u_k^T v_c) - 1) u_k$$

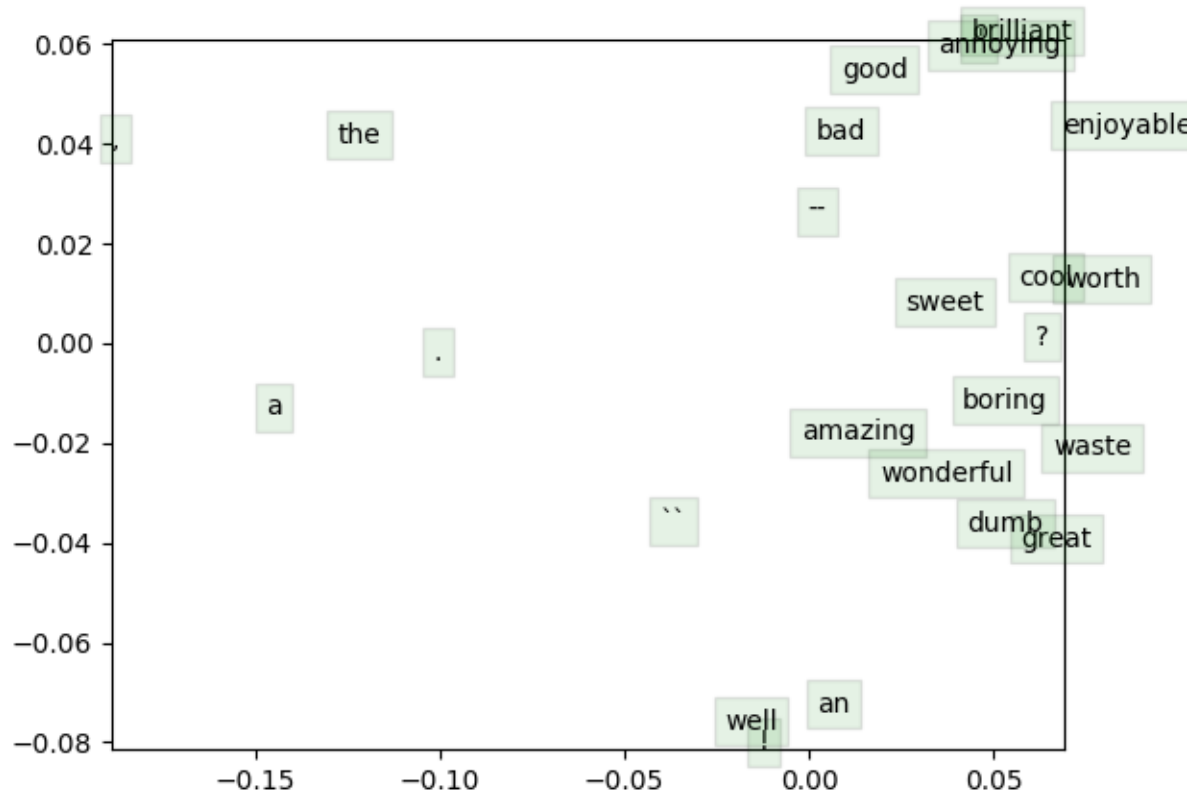
$$\frac{\partial J}{\partial u_o} = (\sigma(u_o^T v_c) - 1) v_c$$

$$\frac{\partial J}{\partial u_k} = -(\sigma(-u_k^T v_c) - 1) v_c$$

Exercise 3. (d) Derive gradients for all of the word vectors for skip-gram and CBOW given the previous parts and given a set of context words

$$\frac{\partial J_{\text{skip-gram}}(Word_{c-m, \dots, c+m})}{\partial U} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_{c+j}, v_c)}{\partial U}$$

$$\frac{\partial J_{\text{skip-gram}}(Word_{c-m, \dots, c+m})}{\partial v_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_{c+j}, v_c)}{\partial v_c}$$



$$\frac{\partial J_{\text{skip_gram}}(\text{Word}_{c-m, \dots, c+m})}{\partial v_j} = 0, j \neq c$$

$$\frac{\partial J_{CBOW}(\text{Word}_{c-m, \dots, c+m})}{\partial U} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_c, \hat{u})}{\partial U}$$

$$\frac{\partial J_{CBOW}(\text{Word}_{c-m, \dots, c+m})}{\partial u_j} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(w_c, \hat{u})}{\partial \hat{u}}, j \in (c-m, \dots, c+m)$$

$$\frac{\partial J_{\text{skip_gram}}(\text{Word}_{c-m, \dots, c+m})}{\partial u_j} = 0, j \notin (c-m, \dots, c+m)$$

Exercise 3. (g)