

a) as it is mentioned before, y is a one-hot vector with a 1 for true outside word & 0 elsewhere so:

$$-\sum_{w \in \text{vocab}} y_w \log(\hat{y}_w) = -\left(\underbrace{y_{w_1}}_0 \log(\hat{y}_{w_1}) + \underbrace{y_{w_2}}_0 \log(\hat{y}_{w_2}) + \dots + \underbrace{y_o}_1 \log(\hat{y}_o) + \dots + \underbrace{y_{w_v}}_0 \log(\hat{y}_{w_v}) \right) =$$

$$= -y_o \log(\hat{y}_o) = -\log(\hat{y}_o)$$

$$b) J(v_c, o, u) = -\log P(o=o | C=c) = -\log \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} =$$

$$= \underbrace{-\log \exp(u_o^T v_c)}_A + \log \underbrace{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)}_B$$

$$\frac{\partial J}{\partial v_c} = -\frac{\frac{\partial A}{\partial v_c}}{A} + \frac{\frac{\partial B}{\partial v_c}}{B} = \frac{-u_o \exp(u_o^T v_c)}{\exp(u_o^T v_c)} + \frac{\sum_{w \in \text{vocab}} u_w^T \exp(u_w^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)} =$$

$$= -u_o + \sum_{w \in \text{vocab}} u_w^T \underbrace{\frac{\exp(u_w^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)}}_{\hat{y}_w} = -u_o + \sum_{w \in \text{vocab}} u_w^T \hat{y}_w = u(\hat{y} - y)$$

$$c) \text{ case 1 - } w=o : \frac{\partial J}{\partial u_w} = -\frac{\frac{\partial A}{\partial u_w}}{A} + \frac{\frac{\partial B}{\partial u_w}}{B} = \frac{-v_c \exp(u_o^T v_c)}{\exp(u_o^T v_c)} +$$

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

$$\text{case 2 - } w \neq o : \frac{\partial J}{\partial u_w} = \frac{\frac{\partial}{\partial u_w} \exp(u_o^T v_c)}{A} + v_c \hat{y}_w = v_c \hat{y}_w$$

$$d) \frac{\partial J}{\partial u} = \left[\frac{\partial J}{\partial u_1}, \frac{\partial J}{\partial u_2}, \dots, \frac{\partial J}{\partial u_{|\text{vocab}|}} \right]$$

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

$$= \sigma(x) \times (-) \sigma(x) \sigma(x) = \sigma(x)(1 - \sigma(x))$$

$$f) f(u) = \log(u), g(u) = \sigma(u) \quad (f(g(u)))' = f'(g(u)) \cdot g'(u)$$

$$\frac{\partial J}{\partial v_c} = - \frac{u_0 \sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c))}{\sigma(u_0^T v_c)} - \sum_{k=1}^K \frac{-u_k \sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c))}{\sigma(-u_k^T v_c)} =$$

$$= -u_0 (1 - \sigma(u_0^T v_c)) + \sum_{k=1}^K u_k (1 - \sigma(-u_k^T v_c))$$

$$\frac{\partial J}{\partial u_0} = - \frac{v_c \sigma(u_0^T v_c) (1 - \sigma(u_0^T v_c))}{\sigma(u_0^T v_c)} - \underbrace{0}_{0 \neq k} = -v_c (1 - \sigma(u_0^T v_c))$$

$$\frac{\partial J}{\partial u_k} = 0 - (-v_c) (1 - \sigma(-u_k^T v_c)) = v_c (1 - \sigma(-u_k^T v_c))$$

$$g) J = -\log \sigma(u_0^T v_c) - \sum_{\substack{\text{Count} \\ u_i = u_k}} \log \sigma(-u_i^T v_c) - \sum_{\substack{u_j \neq u_k}} \log \sigma(-u_j^T v_c)$$

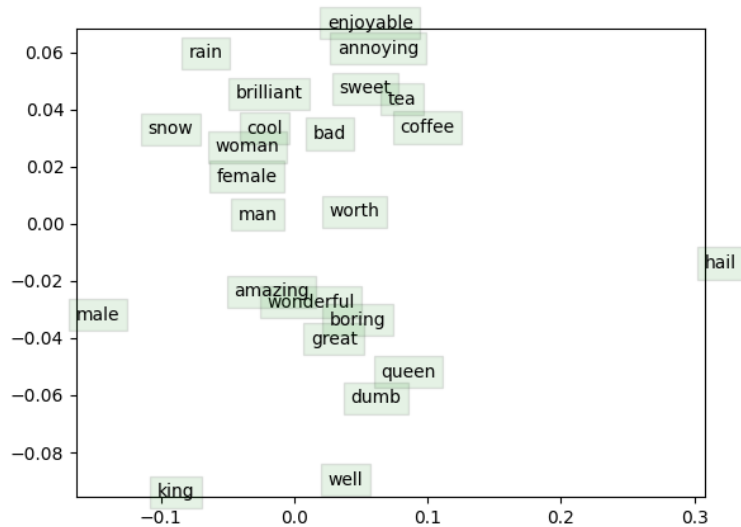
$$\frac{\partial J}{\partial u_k} = 0 - \text{Count} \frac{-v_c \sigma(-u_k^T v_c) (1 - \sigma(-u_k^T v_c))}{\sigma(-u_k^T v_c)} - 0 = \text{Count} \times v_c (1 - \sigma(-u_k^T v_c))$$

$$h) \frac{\partial J}{\partial u} = \sum_{m \leq j \leq m} \frac{\partial}{\partial u} J(v_c, w_{t+j}, u)$$

$$\frac{\partial J}{\partial v_c} = \sum_{m \leq j \leq m} \frac{\partial}{\partial v_c} J(v_c, w_{t+j}, u)$$

$$\frac{\partial J}{\partial v_w (w \neq c)} = 0$$

Coding part:



As it is seen, related words are almost clustered together: rain, snow/ tea, coffee/ amazing, wonderful, boring, great. Moreover, the semantic pattern is almost seen in this plot: male as king, female as queen. Some of the embeddings are not as good as expected (for example male should be clustered with woman, female, man); it is related to the corpus that we have trained on.