

$$J_{\text{naive-softmax}}(V_c, 0, U) \Leftrightarrow J_{\text{ns}}(V_c, 0, U) \quad (1b)$$

$$\begin{aligned} \frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial V_c} &= - \frac{\partial \log P(O=0 | C=c)}{\partial V_c} = - \frac{\partial}{\partial V_c} \log \frac{e^{u_0^T V_c}}{\sum_{w \in V} e^{u_w^T V_c}} \\ &= - \frac{\partial}{\partial V_c} \left[ \underbrace{\log e^{u_0^T V_c}}_{u_0^T V_c} - \log \sum_{w \in V} e^{u_w^T V_c} \right] = - \left( u_0 - \sum_{w \in V} \underbrace{\frac{u_w e^{u_w^T V_c}}{\sum_{w \in V} e^{u_w^T V_c}}}_{P(O=w | C=c) u_w} \right) \quad (\star) \\ &= -u_0 + \sum_{w \in V} u_w \underbrace{P(O=w | C=c)}_{\text{predicted distro } \hat{y}} = \sum_{w \in V} u_w \hat{y} - u_0 = U(\hat{y} - y) \end{aligned}$$

show that:  $-\sum_{w \in V} y_w \log(\hat{y}_w) = -\log(\hat{y}_0)$  (1a)

$$-\sum_{w \in V} y_w \log(\hat{y}_w) = - \left( y_0 \log(\hat{y}_0) + \sum_{w \in V, w \neq 0} y_w \log(\hat{y}_w) \right) = -y_0 \log(\hat{y}_0)$$

$$\forall w \in V \quad \frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial u_w} \rightarrow \text{from } \star \rightarrow \forall w \in V \quad \frac{-\partial}{\partial u_w} (u_0^T V_c - \log \sum_{w \in V} e^{u_w^T V_c})$$

$\rightarrow w \neq 0 \rightarrow \text{zero} + \frac{V_c}{P(O=w | C=c)} = V_c \hat{y}_w$  from  $\star$   
 $\rightarrow w = 0 \rightarrow -V_c + V_c P(O=0 | C=c) = -V_c + V_c \hat{y}_0 = V_c (-1 + \hat{y}_0)$

$$\frac{\partial J_{\text{ns}}(V_c, 0, U)}{\partial U} = V_c (\hat{y} - y)^T \quad (1d) \quad (1c)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = (1 + e^{-x})^{-1} \rightarrow \sigma(x) = \frac{e^x}{1 + e^x}$$

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

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

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

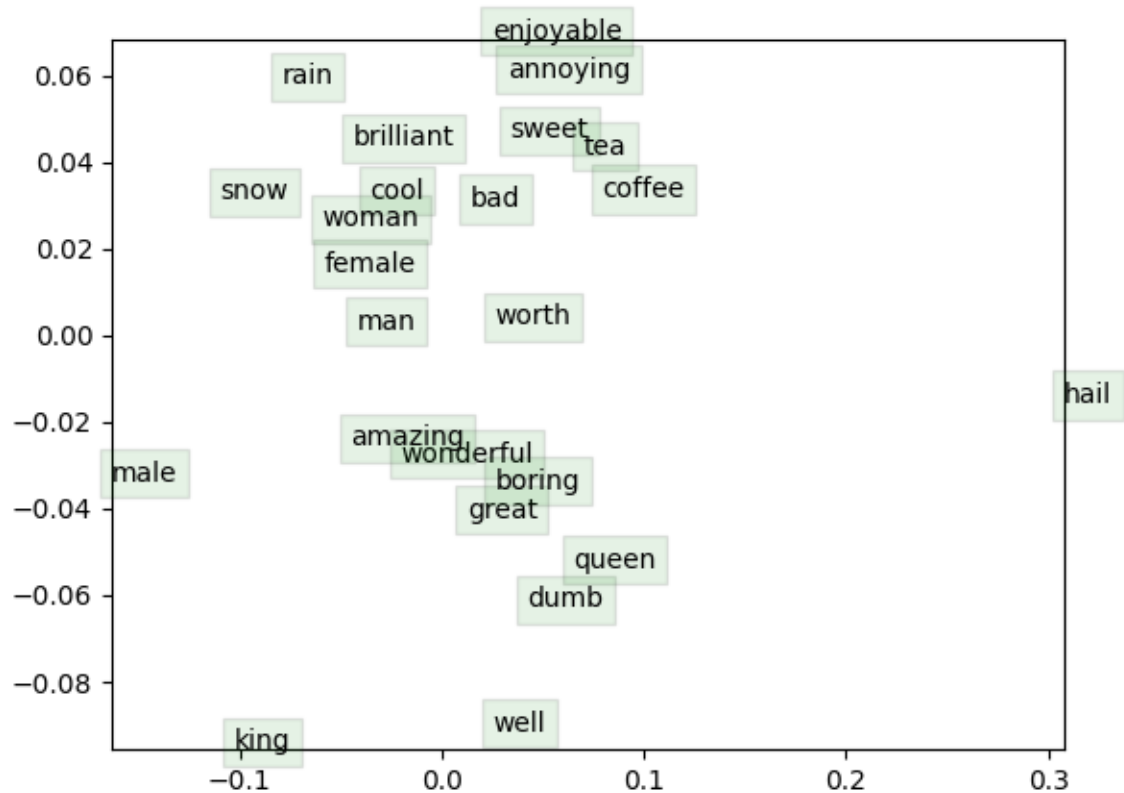
$$= \sum_{k=1}^K \sigma(u_k^T v_c) u_k + u_0 (\sigma(u_0^T) - 1)$$

۱۹

$$\frac{\partial J_{res}}{\partial u_0} = v_c (\sigma(u_0^T v_c) - 1)$$

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

در  $softmax$  پیچیدگی محاسبات با سایر  $V$  (تعداد کمات) رابطه دارد اما در  $res$  این رابطه برقرار نیست و می توان الگوریتم را در زمان محقول روی دایره لغات گسترده تر اجرا کرد.



In this plot, we can see that similar words are mapped near together. (Their vectors are similar)

Woman, female, man

Enjoyable, annoying

Sweet, tea, coffee

Amazing, wonderful, boring, great

We can see, the words in each of these groups are having similar contexts.

Also this is not perfect, since for example king and queen are far from each other, or male is not near man.