**1|**

**a.**

$$CE(y, \hat{y}) = -\sum_w y_w \log(\hat{y}_w)$$

$$y_w = [0, \dots, \overset{o}{\underset{\downarrow}{1}}, \dots, 0]$$

$$= -1 \cdot \log(\hat{y}_0) - \sum_{w \neq 0} 0 \cdot \log(\hat{y}_w)$$

$$= -\log(\hat{y}_0)$$

**b.**

$$\frac{\partial J}{\partial v_c} = \frac{\partial}{\partial v_c} \left[ -\log(\hat{y}_0) \right] = -\frac{\partial}{\partial v_c} \left[ \log\left(\exp(u_0^T v_c)\right) - \log\left(\sum_w \exp(u_w^T v_c)\right) \right]$$

$$= -u_0 + \frac{\partial}{\partial v_c} \log\left(\sum_w \exp(u_w^T v_c)\right)$$

$$= -u_0 + \frac{\sum_w u_w \exp(u_w^T v_c)}{\sum_w \exp(u_w^T v_c)}$$

$$= -u_0 + \sum_x \frac{\exp(u_x^T v_c)}{\sum_w \exp(u_w^T v_c)} u_x$$

$$= \boxed{-u_0 + \sum_w \overset{\vee}{\hat{y}_w} u_w}$$

$u_w : d \times 1 \qquad \frac{\partial}{\partial u_w} : d \times 1$

$v_c : d \times 1$

$\frac{\partial}{\partial u_w} u_w^T v_c = v_c \qquad \frac{\partial}{\partial v_c} u_w^T v_c = u_w$

**c.**

$$\frac{\partial J}{\partial u_w} = -\frac{\partial}{\partial u_w} \left[ \underbrace{\log\left(\exp(u_0^T v_c)\right)}_{\substack{v_c \text{ if } w=0 \\ 0 \quad \text{otherwise}}} - \underbrace{\log\left(\sum_i \exp(u_i^T v_c)\right)}_{\frac{\exp(u_w^T v_c) v_c}{\sum_i \exp(u_i^T v_c)} = \hat{y}_w v_c} \right]$$

$$= \begin{cases} v_c (\hat{y}_w - 1) & \text{if } w = 0 \\ v_c \, \hat{y}_w & \text{otherwise} \end{cases}$$

d.

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

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

$$= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \frac{1}{(1 + e^{-x})^2}$$

$$= \sigma(x) - \sigma^2(x) = \boxed{\sigma(x)(1 - \sigma(x))}$$

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

e.

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

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

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

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

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

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

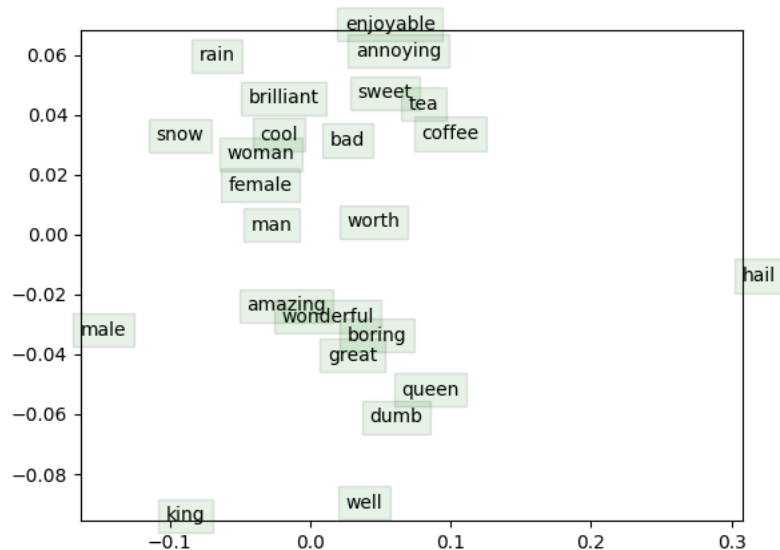Looping over $K \ll V$ much faster than over entire vocabulary

f.

i) $\dfrac{\partial J_{skip}(v_c, w_{t-m}, ..., w_{t+m}, U)}{\partial U} = \displaystyle\sum_{\substack{-m \le j \le m \\ j \ne 0}} \dfrac{\partial J(v_c, w_{t+j}, U)}{\partial U}$

ii) $\dfrac{\partial J_{skip}}{\partial v_c} = \displaystyle\sum_j \dfrac{\partial J(v_c, w_{t+j}, U)}{\partial v_c}$

iii) $\dfrac{\partial J_{skip}}{\partial v_w} = 0$

$w \ne c$

## 2)

c.



- Interchangeable words are grouped [amazing, wonderful]
- Grouped words do not necessarily have similar meaning:
  [boring, great]   [enjoyable, annoying]
- Consistent vector from good to bad:
  bad − brilliant  $\approx$  boring − amazing  $\approx$  dumb − great
  $\approx$ annoying − enjoyable