

Calculating gradient

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t)$$

$$P(o|c) = \frac{e^{u_o^T v_c}}{\sum_{w=1}^V e^{u_w^T v_c}}$$

$$\frac{\partial}{\partial v_c} \log P(o|c) = \frac{\partial}{\partial v_c} \log e^{u_o^T v_c} - \frac{\partial}{\partial v_c} \log \sum_{w=1}^V e^{u_w^T v_c}$$

$$= \frac{\partial}{\partial v_c} u_o^T v_c - \dots$$

$$= u_o - \frac{\partial}{\partial v_c} \underbrace{\log}_{f} \underbrace{\sum_{w=1}^V e^{u_w^T v_c}}_g \quad (f \cdot g)$$

$$= u_o - \frac{1}{\sum_{w=1}^V e^{u_w^T v_c}} \cdot \frac{\partial}{\partial v_c} \sum_{n=1}^V e^{u_n^T v_c}$$

$$= u_o - \dots \cdot \sum_{n=1}^V \frac{\partial}{\partial v_c} e^{u_n^T v_c}$$

$$= u_o - \dots \cdot \sum_{n=1}^V e^{u_n^T v_c} \cdot \frac{\partial}{\partial v_c} u_n^T v_c$$

$$= u_o - \dots \cdot \sum_{n=1}^V e^{u_n^T v_c} \cdot u_n$$

$$= u_o - \frac{\sum_{n=1}^V e^{(u_n^T v_c)} \cdot u_n}{\sum_{w=1}^V e^{(u_w^T v_c)}}$$

$$= u_o - \sum_{n=1}^V \underbrace{\left(\frac{e^{(u_n^T v_c)}}{\sum_{w=1}^V e^{u_w^T v_c}} \right)}_{\text{softmax}} u_n$$

$$= u_o - \underbrace{\sum_{x=1}^V p(x|c) u_x}_{\text{expectation}}$$

(average over all context vectors weighted by P over model)

$$= \text{observed} - \text{expected}$$