

Stanford CS 224n Assignment 2

Kuei-Da Liao (lkueida@eng.ucsd.edu)

October 3, 2020

1. Written: Understanding word2vec

(a)

$$- \sum_{w \in Vocab} y_w \log(\hat{y}_w) = - \sum_{w \in Vocab} \mathbf{1}(w = o) \log(\hat{y}_w) = -\log(\hat{y}_o)$$

(b)

$$-\log P(O = o | C = c) = -\mathbf{u}_o^\top \mathbf{v}_c + \log \sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)$$

$$\begin{aligned} \frac{\partial J_{\text{naive-softmax}}}{\partial \mathbf{v}_c} &= -\mathbf{u}_o^\top + \frac{1}{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c)} \frac{\partial \sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\partial \mathbf{v}_c} \\ &= -\mathbf{u}_o^\top + \sum_{w \in Vocab} \mathbf{u}_w^\top \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c)} \\ &= -\mathbf{u}_o^\top + \sum_{w \in Vocab} \mathbf{u}_w^\top \hat{y}_w \\ &= \sum_{w \in Vocab} (\hat{y}_w - \mathbf{1}(w = o)) \mathbf{u}_w^\top \end{aligned}$$

(c)

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_o} &= -\mathbf{v}_c^\top + \frac{1}{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c)} \frac{\partial \sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\partial \mathbf{u}_o} \\ &= -\mathbf{v}_c^\top + \sum_{w \in Vocab} \mathbf{v}_c^\top \hat{y}_w \mathbf{1}(w = o) \\ &= -\mathbf{v}_c^\top + \mathbf{v}_c^\top \hat{y}_o \\ &= (\hat{y}_o - 1) \mathbf{v}_c^\top \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{u}_{w, w \neq o}} &= \frac{1}{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c)} \frac{\partial \sum_{w \in Vocab} \exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\partial \mathbf{u}_{w, w \neq o}} \\ &= \sum_{w \in Vocab, w \neq o} \mathbf{v}_c^\top \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\sum_{w' \in Vocab} \exp(\mathbf{u}_{w'}^\top \mathbf{v}_c)} \\ &= \sum_{w \in Vocab, w \neq o} \hat{y}_w \mathbf{v}_c^\top \end{aligned}$$

$$\frac{\partial J}{\partial \mathbf{U}} = (\hat{\mathbf{y}} - \mathbf{e}_o^\top) \mathbf{v}_c^\top$$

(d)

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

(e)

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{neg-sample}}}{\partial \mathbf{v}_c} &= -\frac{1}{\cancel{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)}} \cancel{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \sigma(-\mathbf{u}_o^\top \mathbf{v}_c) \mathbf{u}_o^\top - \sum_{k=1}^K -\frac{1}{\cancel{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)}} \cancel{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{u}_k^\top \\ &= (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{u}_o^\top + \sum_{k=1}^K \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{u}_k^\top \\ \frac{\partial \mathbf{J}}{\partial \mathbf{u}_o} &= -\frac{1}{\cancel{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)}} \cancel{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \sigma(-\mathbf{u}_o^\top \mathbf{v}_c) \mathbf{v}_c = (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{v}_c \\ \frac{\partial \mathbf{J}}{\partial \mathbf{u}_k} &= -\frac{1}{\cancel{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)}} \cancel{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \sigma(\mathbf{u}_k^\top \mathbf{v}_c) (-\mathbf{v}_c) = \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{v}_c \end{aligned}$$

(f) (i)

$$\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{U}} = \sum_{\substack{-m \leq j < m \\ j \neq 0}} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}}$$

(ii)

$$\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_c} = \sum_{\substack{-m \leq j < m \\ j \neq 0}} \frac{\partial \mathbf{J}(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c}$$

(iii)

$$\frac{\partial \mathbf{J}_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U})}{\partial \mathbf{v}_{w, w \neq c}} = 0$$