

cs224n 2019 Assignment 1

Zhou Xiaorui - jarvan000@gmail.com

March 20, 2019

Notice: The answer is done by myself, and I'm not a Stanford student.

(a) only the o 'th position of \mathbf{y} is not zero, so the sum can be reduced to only one part:

$$- \sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -y_o \log(\hat{y}_o) = -\log(\hat{y}_o) \quad (1)$$

(b) since

$$\begin{aligned} \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) &= -\log \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \\ &= -\mathbf{u}_o^\top \mathbf{v}_c + \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c) \end{aligned} \quad (2)$$

so we have

$$\begin{aligned} \frac{\partial \mathbf{J}}{\partial \mathbf{v}_c} &= -\mathbf{u}_o + \frac{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^\top \mathbf{v}_c) \mathbf{u}_x}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \\ &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} \frac{\exp(\mathbf{u}_x^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \mathbf{u}_x \\ &= -\mathbf{u}_o + \sum_{x \in \text{Vocab}} \hat{y}_x \mathbf{u}_x \\ &= \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}) \end{aligned}$$

According to the result, we need to compute the matrix multiplication over the whole vocabulary, which may contain millions of words, and that is time-consuming and unnecessary.

(c) when $\mathbf{u}_w = \mathbf{u}_o$

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

when $\mathbf{u}_w \neq \mathbf{u}_o$

$$\begin{aligned}\frac{\partial \mathbf{J}}{\partial \mathbf{u}_w} &= \frac{\exp(\mathbf{u}_w^\top \mathbf{v}_c)}{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^\top \mathbf{v}_c)} \mathbf{v}_c \\ &= \hat{y}_w \mathbf{v}_c\end{aligned}$$

so from above two equation, we can get

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

(d) very basic partial derivative computation.

$$\begin{aligned}\frac{d\sigma(\mathbf{x})}{d\mathbf{x}} &= \frac{e^{-\mathbf{x}}}{(1 + e^{-\mathbf{x}})^2} \\ &= \sigma(\mathbf{x})(1 - \sigma(\mathbf{x}))\end{aligned}$$

(e) basic partial derivative computation, using chain rule.

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

for \mathbf{u}_o and \mathbf{u}_w , we apply the same technique.

$$\begin{aligned}\frac{\partial \mathbf{J}}{\partial \mathbf{u}_o} &= -\frac{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))\mathbf{v}_c}{\sigma(\mathbf{u}_o^\top \mathbf{v}_c)} \\ &= -(1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c))\mathbf{v}_c \\ \frac{\partial \mathbf{J}}{\partial \mathbf{u}_k} &= -\frac{(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))(1 - (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)))(-\mathbf{v}_c)}{\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)} \\ &= (1 - (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c)))(-\mathbf{v}_c)\end{aligned}$$

using Negative Sampling loss, we only need to compute K+1 parameter, which is much more efficient than if we need to compute over the whole vocabulary.

(f) just add the loss of every context word, we can get the answer.

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

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

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

Useful reference: <http://www.amendgit.com/post/cs224n/cs224n-assignment-1/>