

- 1-(a)
- 1-(b)
- 1-(c)
- 1-(d)
- 1-(e)
- 1-(f)

定义如下符号

- d , 词向量的维度
- n , 词汇数量
- $\mathbf{U} \in \mathbb{R}^{d \times n}$, outside词向量矩阵, 每一列是一个词向量, $\mathbf{u}_w \in \mathbb{R}^{d \times 1}$
- $\mathbf{V} \in \mathbb{R}^{d \times n}$, center词向量矩阵, 每一列是一个词向量, $\mathbf{v}_w \in \mathbb{R}^{d \times 1}$
- \mathbf{y} , 真实值, one-hot, 表示是哪个词, $\mathbf{y} \in \mathbb{R}^{n \times 1}$
- $\hat{\mathbf{y}}$, 预测值, 表示属于某个词的概率, $\hat{\mathbf{y}} \in \mathbb{R}^{n \times 1}$

1-(a)

$$J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c). \quad (2)$$

(a) (3 points) Show that the naive-softmax loss given in Equation (2) is the same as the cross-entropy loss between \mathbf{y} and $\hat{\mathbf{y}}$; i.e., show that

¹Assume that every word in our vocabulary is matched to an integer number k . \mathbf{u}_k is both the k^{th} column of \mathbf{U} and the ‘outside’ word vector for the word indexed by k . \mathbf{v}_k is both the k^{th} column of \mathbf{V} and the ‘center’ word vector for the word indexed by k . **In order to simplify notation we shall interchangeably use k to refer to the word and the index-of-the-word.**

²The Cross Entropy Loss between the true (discrete) probability distribution p and another distribution q is $-\sum_i p_i \log(q_i)$.

$$-\sum_{w \in \text{Vocab}} y_w \log(\hat{y}_w) = -\log(\hat{y}_o). \quad (3)$$

Your answer should be one line.

A:因为one-hot向量只有标记的那一维是1, 其他维度都是0, 所以展开左边的 \sum 自然就是右边的代数式。

1-(b)

(b) (5 points) Compute the partial derivative of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{v}_c . Please write your answer in terms of \mathbf{y} , $\hat{\mathbf{y}}$, and \mathbf{U} .

直接对 $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c)$ 求导,

$$\begin{aligned} & \frac{\partial}{\partial \mathbf{v}_c} \left[- \left(\mathbf{u}_o^T \mathbf{v}_c - \log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) \right) \right] \\ &= \frac{\partial}{\partial \mathbf{v}_c} \left(\log \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) \right) - \mathbf{u}_o \\ &= \frac{\sum_{x \in \text{Vocab}} \exp(\mathbf{u}_x^T \mathbf{v}_c) \mathbf{u}_x}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} - \mathbf{u}_o \\ &= \sum_{x \in \text{Vocab}} \frac{\exp(\mathbf{u}_x^T \mathbf{v}_c) \mathbf{u}_x}{\exp(\mathbf{u}_w^T \mathbf{v}_c)} - \mathbf{u}_o \\ &= \sum_{x \in \text{Vocab}} P(O = x | C = c) \mathbf{u}_x - \mathbf{u}_o \end{aligned}$$

该题要求使用 $\mathbf{y}, \hat{\mathbf{y}}, \mathbf{U}$ 表示出结果，我们整理一下推导结果即可。推导的结果表示损失函数对于 \mathbf{v}_c 的偏导数是“outside词向量的期望减去当前outside词向量”，根据这个意思即可以整理答案为

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

1-(c)

(c) (5 points) Compute the partial derivatives of $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to each of the ‘outside’ word vectors, \mathbf{u}_w ’s. There will be two cases: when $w = o$, the true ‘outside’ word vector, and $w \neq o$, for all other words. Please write your answer in terms of $\mathbf{y}, \hat{\mathbf{y}}$, and \mathbf{v}_c .

该题求的是 J 对 \mathbf{U} 的偏导数，可以先求 $\frac{\partial J}{\partial \mathbf{u}_w}$ ，这个推导和1-(b)的推导过程非常的相似，我们此处直接写出答案，就不写详细的推导了。

$$\frac{\partial J}{\partial \mathbf{u}_w} = \begin{cases} (\hat{y}_w - 1) \mathbf{v}_c, & w = o \\ \hat{y}_w \mathbf{v}_c, & w \neq o \end{cases}$$

根据标量 f 对矩阵 \mathbf{X} 的求导规律： $\frac{\partial f}{\partial \mathbf{X}} = \left[\frac{\partial f}{\partial X_{ij}} \right]$ ，可以得到答案

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

小技巧：如果对于向量或者矩阵的求导犯迷糊了，可以多从维度的角度去check。

1-(d)

(d) (3 Points) The sigmoid function is given by Equation 4:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (4)$$

Please compute the derivative of $\sigma(x)$ with respect to \mathbf{x} , where \mathbf{x} is a vector.

标量 f 对向量 \mathbf{x} 求导, 结果是 $\frac{\partial f}{\partial \mathbf{x}} = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_m} \right]^T$, 又 $\sigma'(x) = \sigma(x)(1 - \sigma(x))$, 所以有

$$\sigma'(\mathbf{x}) = \begin{bmatrix} \sigma(x_1)(1 - \sigma(x_1)) \\ \sigma(x_2)(1 - \sigma(x_2)) \\ \dots \\ \sigma(x_m)(1 - \sigma(x_m)) \end{bmatrix}$$

1-(e)

(e) (4 points) Now we shall consider the Negative Sampling loss, which is an alternative to the Naive Softmax loss. Assume that K negative samples (words) are drawn from the vocabulary. For simplicity of notation we shall refer to them as w_1, w_2, \dots, w_K and their outside vectors as $\mathbf{u}_1, \dots, \mathbf{u}_K$. Note that $o \notin \{w_1, \dots, w_K\}$. For a center word c and an outside word o , the negative sampling loss function is given by:

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^T \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^T \mathbf{v}_c)) \quad (5)$$

for a sample w_1, \dots, w_K , where $\sigma(\cdot)$ is the sigmoid function.³

Please repeat parts (b) and (c), computing the partial derivatives of $J_{\text{neg-sample}}$ with respect to \mathbf{v}_c , with respect to \mathbf{u}_o , and with respect to a negative sample \mathbf{u}_k . Please write your answers in terms of the vectors \mathbf{u}_o , \mathbf{v}_c , and \mathbf{u}_k , where $k \in [1, K]$. After you've done this, describe with one sentence why this loss function is much more efficient to compute than the naive-softmax loss. Note, you should be able to use your solution to part (d) to help compute the necessary gradients here.

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

类似的有

$$\frac{\partial J_{\text{neg-sample}}}{\partial \mathbf{u}_o} = (\sigma(\mathbf{u}_o^T \mathbf{v}_c) - 1)\mathbf{v}_c$$

类似的有

$$\frac{\partial J_{\text{neg-sample}}}{\partial \mathbf{u}_k} = (1 - \sigma(-\mathbf{u}_k^T \mathbf{v}_c))\mathbf{v}_c$$

$J_{\text{neg-sample}}$ 比 $J_{\text{naive-softmax}}$ 训练效率高的主要原因是前者只需要和相关的词向量进行内积, 而后者的偏导数里面有 \hat{y} 这一项, 该项使用 Softmax 公式得出, 需要计算 \mathbf{v}_c 和所有的词向量内积, 这个步骤会消耗较多的计算资源。

1-(f)

- (f) (3 points) Suppose the center word is $c = w_t$ and the context window is $[w_{t-m}, \dots, w_{t-1}, w_t, w_{t+1}, \dots, w_{t+m}]$, where m is the context window size. Recall that for the skip-gram version of word2vec, the total loss for the context window is:

$$J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) \quad (6)$$

Here, $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ represents an arbitrary loss term for the center word $c = w_t$ and outside word w_{t+j} . $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ could be $J_{\text{naive-softmax}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ or $J_{\text{neg-sample}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$, depending on your implementation.

Write down three partial derivatives:

- (i) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{U}$
- (ii) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_c$

³Note: the loss function here is the negative of what Mikolov et al. had in their original paper, because we are doing a minimization instead of maximization in our assignment code. Ultimately, this is the same objective function.

CS 224n Assignment #2: word2vec (43 Points)

- (iii) $\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) / \partial \mathbf{v}_w$ when $w \neq c$

Write your answers in terms of $\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{U}$ and $\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U}) / \partial \mathbf{v}_c$. This is very simple – each solution should be one line.

损失函数就是窗口滑动的过程中，每个窗口的损失累加，所以偏导数也容易得出（“多项和的导数”等于“多项导数的和”）。

$$\begin{aligned} \frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{U}} &= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{U}} \\ \frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{v}_c} &= \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial J(\mathbf{v}_c, w_{t+j}, \mathbf{U})}{\partial \mathbf{v}_c} \\ \frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{v}_w} &= 0, \text{ when } w \neq c \end{aligned}$$