

CS 224n Assignment #2: word2vec

Intae Jun

1 Written: Understanding word2vec (31 points)

- (a) (2 points) Prove that the naive-softmax loss (Eq. 2) is the same as the cross-entropy loss between \mathbf{y} and $\hat{\mathbf{y}}$.

Since the true label \mathbf{y}_w is a one-hot vector, all the terms where $w \neq o$ are reduced and only the term where $w = o$ remains.

- (b) (7 points)

- (i) Compute the partial derivative of $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{v}_c .

Please write your answer in terms of $\mathbf{y}, \hat{\mathbf{y}}, \mathbf{U}$, and show your work to receive full credit.

$$\begin{aligned} \frac{\partial \mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})}{\partial \mathbf{v}_c} &= - \frac{\partial \log P(O = o | C = c)}{\partial \mathbf{v}_c} \\ &= - \frac{1}{P(O = o | C = c)} \frac{\partial P(O = o | C = c)}{\partial \mathbf{v}_c} \\ &= - \frac{\mathbf{u}_o \exp(\mathbf{u}_o^T \mathbf{v}_c) \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) - \exp(\mathbf{u}_o^T \mathbf{v}_c) \sum_{w \in \text{Vocab}} \mathbf{u}_w \exp(\mathbf{u}_w^T \mathbf{v}_c)}{P(O = o | C = c) (\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c))^2} \\ &= - \frac{\mathbf{u}_o \exp(\mathbf{u}_o^T \mathbf{v}_c) \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) - \exp(\mathbf{u}_o^T \mathbf{v}_c) \sum_{w \in \text{Vocab}} \mathbf{u}_w \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\exp(\mathbf{u}_o^T \mathbf{v}_c) (\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c))} \\ &= - \frac{\mathbf{u}_o \sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c) - \sum_{w \in \text{Vocab}} \mathbf{u}_w \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{u}_o + \sum_{w \in \text{Vocab}} P(O = w | C = c) \mathbf{u}_w \\ &= \mathbf{U}^T (\hat{\mathbf{y}} - \mathbf{y}). \end{aligned}$$

- (ii) When $\hat{\mathbf{y}} - \mathbf{y} = 0$, the gradient becomes zero.

- (iii) As the gradient is subtracted from the word vector \mathbf{v}_c , the vector is updated as follows:

$$\mathbf{v}_c^{(i+1)} \leftarrow \mathbf{v}_c^{(i)} - \eta \mathbf{U}^T (\hat{\mathbf{y}} - \mathbf{y}),$$

where i is the index for the iteration number and η is the learning rate for the gradient descent method. By subtracting this gradient, we are moving the center word vector \mathbf{v}_c towards the correct output vector \mathbf{u}_o and away from the incorrect predictions $\mathbf{U}^T \hat{\mathbf{y}}$.

(iv) For two different words $x \neq y$, if, for some scalar α , they have a relationship such as $\mathbf{u}_x = \alpha \mathbf{u}_y$, the L2 normalized vectors \mathbf{u}'_x and \mathbf{u}'_y are not discriminated since L2 normalization makes every vector have the same magnitude. However the primary information contained in the vectors, such as representation of each word does not disappear. Also, keeping consistency in vector scale are recommended for some models that are sensitive to the scale of input features.

(c) (5 points) Compute the partial derivatives of $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ 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 . In this subpart, you may use specific elements within these terms as well (such as $\mathbf{y}_1, \mathbf{y}_2, \dots$).

$$\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c) = -\log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)} = -\mathbf{u}_o^T \mathbf{v}_c + \log \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c).$$

If $w = o$ (i.e. the outside word is the answer):

$$\begin{aligned} \frac{\partial \mathbf{J}}{\partial \mathbf{u}_w} &= -\mathbf{v}_c + \frac{\partial}{\partial \mathbf{u}_w} (\log \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)) \\ &= -\mathbf{v}_c + \frac{\mathbf{v}_c \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\log \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{v}_c + P(O = o | C = c) \mathbf{v}_c \\ &= (\hat{\mathbf{y}} - \mathbf{y})^T \mathbf{v}_c. \end{aligned}$$

If $w \neq o$ (i.e. the outside word is not the answer):

$$\begin{aligned} \frac{\partial \mathbf{J}}{\partial \mathbf{u}_w} &= -\frac{\partial}{\partial \mathbf{u}_w} (\mathbf{u}_o^T \mathbf{v}_c) + \frac{\partial}{\partial \mathbf{u}_w} (\log \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)) \\ &= 0 + \frac{\mathbf{v}_c \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)}{\log \sum_w \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= P(O = w | C = c) \mathbf{v}_c \\ &= \hat{\mathbf{y}}^T \mathbf{v}_c \\ &= (\hat{\mathbf{y}} - \mathbf{y})^T \mathbf{v}_c. \end{aligned}$$

(d) (1 point) Write down the partial derivative of $\mathbf{J}_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ with respect to \mathbf{U} .

(e) (2 points) Compute the derivate of $f(x)$ with respect to x .

If $x > 0$, $f(x) = x$ and $f'(x) = 1$. If $x < 0$, then $f(x) = \alpha x$ and $f'(x) = \alpha$.

In summary, we have that

$$\frac{d}{dx}f(x) = \begin{cases} x & (x > 0), \\ \alpha & (x < 0). \end{cases}$$

(f) (3 points) Compute the derivative of $\sigma(x)$ with respect to x , where x is a scalar.

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

(g) (6 points)

(i) .

(h) (2 points)

(i) .