

Solutions to assignment1 of CS224n

陈承勃

May 22, 2018

1 Softmax

(a) Omitted. (b) See `q1_softmax.py`

2 Neural Network Basics

(a) $\sigma'(x) = \sigma(x)\sigma(1-x)$

(b) Assume k is the correct class, then

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = -y_k \log \hat{y}_k = -\log \hat{y}_k = -\log \frac{\exp(\theta_k)}{\sum_i \exp(\theta_i)} = -\theta_k + \log \sum_i \exp(\theta_i).$$

$$\therefore \frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial \theta_k} = -1 + \frac{\exp(\theta_k)}{\sum_i \exp(\theta_i)} = \hat{y}_k - 1,$$

$$\frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial \theta_j} = \frac{\exp(\theta_j)}{\sum_i \exp(\theta_i)} = \hat{y}_j, \quad j \neq k.$$

$$\therefore \frac{\partial CE(\mathbf{y}, \hat{\mathbf{y}})}{\partial \boldsymbol{\theta}} = \hat{\mathbf{y}} - \mathbf{y}$$

(c) The forward propagation steps:

$$\mathbf{Z}_1 = \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1, \quad \mathbf{h} = \text{sigmoid}(\mathbf{Z}_1)$$

$$\mathbf{Z}_2 = \mathbf{h}\mathbf{W}_2 + \mathbf{b}_2, \quad \hat{\mathbf{y}} = \text{sigmoid}(\mathbf{Z}_2)$$

$$\mathbf{J} = CE(\mathbf{y}, \hat{\mathbf{y}})$$

The backward propagation:

$$\frac{\partial \mathbf{J}}{\partial \mathbf{Z}_2} = \hat{\mathbf{y}} - \mathbf{y} \triangleq \boldsymbol{\delta}_1, \quad \frac{\partial \mathbf{J}}{\partial \mathbf{h}} = \boldsymbol{\delta}_1 \mathbf{W}_2^T \triangleq \boldsymbol{\delta}_2$$

$$\frac{\partial \mathbf{J}}{\partial \mathbf{Z}_1} = \boldsymbol{\delta}_2 * \sigma'(\mathbf{Z}_1) \triangleq \boldsymbol{\delta}_3, \quad * \text{ denotes element-wise product.}$$

$$\frac{\partial \mathbf{J}}{\partial \mathbf{x}} = \boldsymbol{\delta}_3 \mathbf{W}_1^T$$

(d) $(1 + D_x) \times H + (1 + H) \times D_y$

(e) See `q2_sigmoid.py`

(f) See `q2_gradcheck.py`

(g) See `q2_neural.py`

3 word2vec

$$\begin{aligned} \text{(a)} \quad J_{\text{softmax-CE}}(o, \mathbf{v}_c, \mathbf{U}) &= CE(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_i y_i \log(\hat{y}_i) = -\log \hat{y}_o = -\log \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} \\ &= -\mathbf{u}_o^T \mathbf{v}_c + \log \sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c) \end{aligned}$$

$$\begin{aligned}\therefore \frac{\partial J}{\partial \mathbf{v}_c} &= -\mathbf{u}_o + \frac{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c) \mathbf{u}_w}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} = -\mathbf{u}_o + \sum_{w=1}^V \frac{\exp(\mathbf{u}_w^T \mathbf{v}_c)}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{u}_w = -\mathbf{u}_o + \sum_{w=1}^V \widehat{y}_w \mathbf{u}_w \\ \therefore \frac{\partial J}{\partial \mathbf{v}_c} &= \sum_{w=1}^V \widehat{y}_w \mathbf{u}_w - \sum_{w=1}^V y_w \mathbf{u}_o = \mathbf{U}(\widehat{\mathbf{y}} - \mathbf{y})\end{aligned}$$

$$(b) \quad \frac{\partial J}{\partial \mathbf{u}_o} = -\mathbf{v}_c + \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c) \mathbf{v}_c}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} = (\widehat{\mathbf{y}}_o - 1) \mathbf{v}_c$$

$$\frac{\partial J}{\partial \mathbf{u}_k} = \frac{\exp(\mathbf{u}_k^T \mathbf{v}_c) \mathbf{v}_c}{\sum_{w=1}^V \exp(\mathbf{u}_w^T \mathbf{v}_c)} = \widehat{\mathbf{y}}_k \mathbf{v}_c, \text{ for } k \neq o.$$

$$\therefore \frac{\partial J}{\partial \mathbf{U}} = \mathbf{v}_c (\widehat{\mathbf{y}} - \mathbf{y})^T, \text{ or } \frac{\partial J}{\partial \mathbf{u}_k} = \begin{cases} (\widehat{\mathbf{y}}_o - 1) \mathbf{v}_c & k = o \\ \widehat{\mathbf{y}}_k \mathbf{v}_c & k \neq o \end{cases}$$

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

(d) Let \mathbf{U} be the collection of all output vectors for all words in the vocabulary.

For Skip-Gram model,

$$\frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{v}_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial \mathbf{F}(w_{t+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c}, \quad \frac{\partial J_{\text{skip-gram}}}{\partial \mathbf{U}} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial \mathbf{F}(w_{t+j}, \mathbf{v}_c)}{\partial \mathbf{U}}$$

For CBOW model,

$$\begin{aligned}\frac{\partial J_{\text{CBOW}}}{\partial \mathbf{v}_{w_{t+j}}} &= \frac{\partial \mathbf{F}(w_t, \widehat{\mathbf{v}})}{\partial \widehat{\mathbf{v}}}, \text{ for all } j \in \{-m, \dots, -1, 1, \dots, m\} \\ \frac{\partial J_{\text{CBOW}}}{\partial \mathbf{v}_{w_{t+j}}} &= 0, \text{ for all } j \notin \{-m, \dots, -1, 1, \dots, m\} \\ \frac{\partial J_{\text{CBOW}}}{\partial \mathbf{U}} &= \frac{\partial \mathbf{F}(w_t, \widehat{\mathbf{v}})}{\partial \mathbf{U}}\end{aligned}$$