

Solution for Assignment 1

1. Softmax

$$(a). \text{softmax}(x_i + c) = \frac{e^{x_i+c}}{\sum_j e^{x_j+c}} = \frac{e^c * e^{x_i}}{e^c * \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(x_i)$$

2. Neural Network Basics

$$(a). \sigma(x) = \frac{1}{1+e^{-x}} \quad \sigma'(x) = \frac{1}{(1+e^{-x})^2} e^{-x} = \sigma^2(x) \left(\frac{1}{\sigma(x)} - 1 \right) = \sigma(x)(1 - \sigma(x))$$

$$(b). CE(y, \hat{y}) = -\sum_i y_i \log(\hat{y}_i)$$

Assume: $y_i = 1, i = k; y_i = 0, i \neq k$

Thus: $CE(y, \hat{y}_i) = -\log(\text{softmax}(\theta_k))$

$$CE(y, \hat{y}) = -\frac{\alpha y_k \log \frac{e^{\theta_k}}{\sum_j e^{\theta_j}}}{\alpha \theta} = \hat{y} - y$$

$$(c). \text{Assume: } \theta_i^{(1)} = \sum_j x_j W_{ji}^{(1)} + b_1$$

$$h_i = \sigma(\theta_i^{(1)})$$

$$\theta_i^{(2)} = \sum_j h_j W_{ji}^{(2)}$$

$$\hat{y}_i = \text{softmax}(\theta_i^{(2)})$$

Then:

$$\frac{\alpha}{\alpha h_i} J = \sum_j \frac{\alpha J}{\alpha \theta_j^{(2)}} \frac{\alpha \theta_j^{(2)}}{\alpha h_i} = \sum_j (\hat{y}_j - y_j) * W_{ij}^{(2)} = (\hat{y} - y) W_{i*}^{(2)T}$$

$$\frac{\alpha}{\alpha \theta_i^{(1)}} J = \sum_j \frac{\alpha J}{\alpha h_j} \frac{\alpha h_j}{\alpha \theta_i^{(1)}} = \sum_j (\hat{y} - y) W^{(2)T} \frac{\alpha \sigma(\theta_j^{(1)})}{\alpha \theta_i^{(1)}} = (\hat{y} - y) W_{i*}^{(2)T} \sigma'(\theta_i^{(1)})$$

Thus:

$$\frac{\alpha}{\alpha x_i} J = \sum_j \frac{\alpha J}{\alpha \theta_j^{(1)}} \frac{\alpha \theta_j^{(1)}}{\alpha x_i} = \sum_j (\hat{y} - y) W_{j*}^{(2)T} \frac{\alpha \sum_i W_{ij}^{(1)} + b_i}{\alpha x_i} = \sum_j (\hat{y} - y) W_{j*}^{(2)T} \sigma'(\theta_j) W_{ij} = [(\hat{y} - y)(W^{(2)T} \otimes \sigma'(\theta^{(1)}))] W_{i*}^{(1)T}$$

, Where \otimes denotes broadcasting and element-wise multiply.

(d). The first Parameter Matrix is of $(D_x + 1) * H$ size, The second Parameter Matrix is of $(H + 1) * D_y$ size (1 denotes bias); Thus total parameter number is

$$(D_x + 1) * H + (H + 1) * D_y = H(D_x + D_y) + H + D_y$$

(g). For $W^{(2)}$:

$$\frac{\alpha}{\alpha W_{ij}^{(2)}} J = \sum_j \frac{\alpha J}{\alpha \theta_j^{(2)}} \frac{\alpha \theta_j^{(2)}}{\alpha W_{ij}^{(2)}} = (\hat{y}_j - y_j) h_i$$

$$\text{Thus: } \frac{\alpha}{\alpha W^{(2)}} J = h^T (\hat{y} - y)$$

For $W^{(1)}$:

$$\frac{\alpha}{\alpha W_{ij}^{(1)}} J = \sum_j \frac{\alpha J}{\alpha \theta_j^{(1)}} \frac{\alpha \theta_j^{(1)}}{\alpha W_{ij}^{(1)}} = (\hat{y} - y) W_{i*}^{(2)T} \sigma'(\theta_i^{(1)}) x_i$$

$$\text{Thus: } \frac{\alpha}{\alpha W^{(1)}} J = X^T \frac{\alpha J}{\alpha \theta^{(1)}}$$

3. word2vec

(a). assume: all vectors below are coloumn vectors. Let $\theta = U^T v_c$,

$$\text{Then } \frac{\alpha CE(y, \hat{y}_o)}{\alpha v_c} = \frac{\alpha CE(y, \hat{y}_o)}{\alpha \theta} \frac{\alpha \theta}{\alpha v_c} = U(\hat{y} - o)$$

$$(b). \frac{\alpha CE(y, \hat{y}_o)}{\alpha U_{ij}} = \sum_k \frac{\alpha CE(y, \hat{y}_o)}{\alpha \theta_k} \frac{\alpha \theta_k}{\alpha U_{ij}} = \sum_k (\hat{y}_k - o_k) \frac{\alpha \sum_i U_{ki}^T (V_c)_i}{\alpha U_{ij}} = (\hat{y}_j - o_j)(v_c)_i$$

$$\text{Thus: } \frac{\alpha CE(y, \hat{y}_o)}{\alpha U} = v_c(\hat{y} - o)^T$$

$$(c). \frac{\alpha J_{neg}}{\alpha v_c} = -\frac{\alpha \log(\sigma(u_o^T v_c))}{\alpha v_c} - \sum_{k=1}^K \frac{\alpha \log(\sigma(-u_k^T v_c))}{\alpha v_c}$$

$$= -\frac{1}{\sigma(u_o^T v_c)} \sigma(u_o^T v_c)(1 - \sigma(u_o^T v_c))u_o - \sum_{k=1}^K \frac{1}{\sigma(-u_k^T v_c)} \sigma(-u_k^T v_c)(1 - \sigma(-u_k^T v_c))(-u_k)$$

$$= (\sigma(u_o^T v_c) - 1)u_o - \sum_{k=1}^K (\sigma(-u_k^T v_c) - 1)u_k$$

$$\frac{\alpha J_{neg}}{\alpha u_w} = -\frac{\alpha \log(\sigma(u_o^T v_c))}{\alpha u_w} - \sum_{k=1}^K \frac{\alpha \log(\sigma(1 - u_k^T v_c))}{\alpha u_w},$$

$$\text{When } w = 0, \frac{\alpha J_{neg}}{\alpha u_w} = [\sigma(u_o^T v_c) - 1]v_c,$$

$$\text{When } w \in S, \frac{\alpha J_{neg}}{\alpha u_w} = [1 - \sigma(-u_w^T v_c)]v_c$$

$$\text{When } w \notin S, \frac{\alpha J_{neg}}{\alpha u_w} = 0$$

Reason for why this cost function can be more effienct: $CE(y, \hat{y})$ contains lots of exponents calculation while this cost function only contains K+1 exponents calculation. And exponents calculation requires lots of calculating time.

(d). For Skip-Gram:

$$\frac{\alpha}{\alpha v_k} J_{skip-gram} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\alpha F(w_{c+j}, v_c)}{\alpha v_k},$$

$$\frac{\alpha}{\alpha u_k} J_{skip-gram} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\alpha F(w_{c+j}, v_c)}{\alpha u_k}$$

For CBOW:

$$\frac{\alpha}{\alpha v_k} J_{CBOW} = \frac{\alpha F(w_{c+w}, \hat{v})}{\alpha v_k} = \frac{\alpha F(w_{c+w}, \hat{v})}{\alpha \hat{v}} \frac{\alpha \hat{v}}{\alpha v_k}$$

$$\text{when } k \in [c - m, c + m], k \neq c, \frac{\alpha}{\alpha v_k} J_{CBOW} = \frac{\alpha F(w_{c+w}, \hat{v})}{\alpha \hat{v}}$$

$$\text{else, } \frac{\alpha}{\alpha v_k} J_{CBOW} = 0$$

$$\frac{\alpha}{\alpha u_k} J_{CBOW} = \frac{\alpha F(w_c, \hat{v})}{\alpha u_k}$$

4. Sentiment Analysis

(b). Reason: avoid parameters being too large; also this can avoid overfitting.

(d). Reason:

The glove model is trained by large contexts, so the pretrained vectors can better describe a word in vector.

The pretrained vectors have more dimensions(50-dimension) than what we trained(10-dimension).

The glove model is trained by statistical methods, which is different from skip-gram methods. The glove model has a better model for word vector training.

(e). As regularization value grows larger, training accuracy is dropping while dev accuracy is first increasing then decreasing.

We can see the regularization here plays a role for avoid overfitting.

But when the regularization growing too large, the cost function seems will be smaller when the parameter is small, so in this case, too large regularization value will cause parameter becoming very small.

(f). The model tends to predict sentences as '+' or '-'. This may because the sentiment analysis criteria is obscure and the labelled '+' or '-' training sets can cross over the whole vector space. And the vector summing up for average seems not a good idea, because when "not good" with "bad" vs. "not bad" with "good" will sometimes generate the same average vector, but these are two different situations.

(g). Sentence1(line71): it 's refreshing to see a girl-power movie that does n't feel it has to prove anything .

Analyze: the most important attitude word here is "refreshing", however, the whole words after "see" have words like "doesn't feel", "anything" can lead the sentence to a negative space. The reason is that all vectors weights the same in sentence average vector.