# Assignment 1

Chen Luo

March 21, 2018

#### 1 Softmax

#### 1.1 $\mathbf{a}$

$$softmax(\mathbf{x}+c) = \frac{e^{\mathbf{x}+\mathbf{c}}}{\sum_{j=1}^{d} e^{x_{j}+c}}$$

$$= \frac{e^{\mathbf{x}}e^{c}}{e^{c}\sum_{j=1}^{d} e^{x_{j}}} = softmax(\mathbf{x})$$
(2)

$$= \frac{e^{\boldsymbol{x}}e^{c}}{e^{c}\sum_{j=1}^{d}e^{x_{j}}} = softmax(\boldsymbol{x})$$
 (2)

#### NN Basic $\mathbf{2}$

#### 2.1a

$$\frac{\partial \sigma(x)}{\partial x} = -1(1 + e^{-x})^{-2} \frac{\partial (1 + e^{-x})}{\partial x}$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2}$$
(4)

$$=\frac{e^{-x}}{(1+e^{-x})^2}\tag{4}$$

$$=\frac{e^{-x}+1-1}{(1+e^{-x})^2}\tag{5}$$

$$= \sigma(x) - \sigma(x)^2 \tag{6}$$

2.2b

$$CE(\boldsymbol{y}, \hat{\boldsymbol{y}}) = -\sum_{i} y_{i} \log(\hat{y}_{i})$$
(7)

$$\hat{\mathbf{y}} = softmax(\boldsymbol{\theta}) = \frac{e^{\boldsymbol{\theta}}}{\sum_{j} e^{\theta_{j}}}$$
 (8)

$$\frac{\partial CE(\boldsymbol{y}, \hat{\boldsymbol{y}})}{\partial \boldsymbol{\theta}} = -\frac{\partial \sum_{i} 1_{[y_{i}=1]} \left\{ \log e^{\theta_{i}} - \log \sum_{j} e^{\theta_{j}} \right\}}{\partial \boldsymbol{\theta}}$$

$$= -\frac{\partial \boldsymbol{y}^{T} \boldsymbol{\theta} - \log \sum_{j} e^{\theta_{j}}}{\partial \boldsymbol{\theta}}$$
(9)

$$= -\frac{\partial \boldsymbol{y}^T \boldsymbol{\theta} - \log \sum_j e^{\theta_j}}{\partial \boldsymbol{\theta}} \tag{10}$$

$$= -y + \frac{\log \mathbf{1}^T e^{\boldsymbol{\theta}}}{\partial \boldsymbol{\theta}} \tag{11}$$

$$= -\boldsymbol{y} + \frac{e^{\boldsymbol{\theta}}}{\mathbf{1}^{T_{\boldsymbol{\rho}}\boldsymbol{\theta}}} \tag{12}$$

$$= -y + softmax(\boldsymbol{\theta}) \tag{13}$$

$$= \hat{\boldsymbol{y}} - \boldsymbol{y} \tag{14}$$

2.3 $\mathbf{c}$ 

$$J = CE(\boldsymbol{y}, \hat{\boldsymbol{y}}) = CE(\boldsymbol{y}, softmax(\boldsymbol{h}\boldsymbol{W}_2 + \boldsymbol{b}_2))$$
(15)

$$= CE(\boldsymbol{y}, softmax(\sigma(\boldsymbol{x}\boldsymbol{W}_1 + \boldsymbol{b}_1)\boldsymbol{W}_2 + \boldsymbol{b}_2))$$
 (16)

Actually,

$$J = -\mathbf{y}_{1 \times D_n} \log(\mathbf{h}_{1 \times H} \mathbf{W}_{H \times D_n} + \mathbf{b}_{1 \times D_n})^T$$
(17)

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial (h W_2 + b_2)} \frac{\partial (h W_2 + b_2)}{\partial h} \frac{\partial h}{\partial x}$$
(18)

$$= (\hat{\boldsymbol{y}} - \boldsymbol{y}) \boldsymbol{W}_{2}^{T} \sigma(\boldsymbol{x} \boldsymbol{W}_{1} + \boldsymbol{b}_{1}) \{ 1 - \sigma(\boldsymbol{x} \boldsymbol{W}_{1} + \boldsymbol{b}_{1}) \}^{T} \boldsymbol{W}_{1}^{T}$$
(19)

, where  $\frac{\partial (hW_2 + b_2)}{\partial h}$  is Jaccobian matrix.

Or, [TODO]

$$d(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2) = d\mathbf{h}\mathbf{W}_2 + \mathbf{h}d\mathbf{W}_2 + d\mathbf{b}_2 \tag{20}$$

$$vec(d(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2)) = vec(d\mathbf{h}\mathbf{W}_2 + \mathbf{h}d\mathbf{W}_2 + d\mathbf{b}_2)$$
(21)

$$= vec(\mathbf{1}d\mathbf{h}\mathbf{W}_2 + \mathbf{h}d\mathbf{W}_2\mathbf{I} + \mathbf{1}d\mathbf{b}_2\mathbf{I}) \rightarrow preparation for \qquad (22)$$

$$vec(AXB) = (B^T \otimes A)vec(X) \tag{23}$$

$$= \boldsymbol{W}_{2}^{T} \otimes 1 vec(d\boldsymbol{h}) + \boldsymbol{I} \otimes \boldsymbol{h} vec(d\boldsymbol{W}_{2}) + \boldsymbol{I} \otimes 1 vec(d\boldsymbol{b}_{2})$$
 (24)

#### 2.4 d

$$\mathbf{W}_1: D_x \times H \tag{25}$$

$$\boldsymbol{b}_1: 1 \times H \tag{26}$$

$$\mathbf{W}_2: H \times D_y \tag{27}$$

$$\boldsymbol{b}_2: 1 \times D_y \tag{28}$$

$$H(D_x + D_y + 1) + D_y (29)$$

# 2.5 g

$$dJ = (\hat{\boldsymbol{y}} - \boldsymbol{y})^T d(\boldsymbol{h} \boldsymbol{W}_2 + \boldsymbol{b}_2) = (\hat{\boldsymbol{y}} - \boldsymbol{y})^T (d\boldsymbol{h} \boldsymbol{W}_2 + \boldsymbol{h} d\boldsymbol{W}_2 + d\boldsymbol{b}_2)$$
(30)

$$= (\hat{\boldsymbol{y}} - \boldsymbol{y})^T d\boldsymbol{h} \boldsymbol{W}_2 + (\hat{\boldsymbol{y}} - \boldsymbol{y})^T \boldsymbol{h} d\boldsymbol{W}_2 + (\hat{\boldsymbol{y}} - \boldsymbol{y})^T d\boldsymbol{b}_2)$$
(31)

$$= trace((\hat{\boldsymbol{y}} - \boldsymbol{y})^T d\boldsymbol{h} \boldsymbol{W}_2 + (\hat{\boldsymbol{y}} - \boldsymbol{y})^T \boldsymbol{h} d\boldsymbol{W}_2 + (\hat{\boldsymbol{y}} - \boldsymbol{y})^T d\boldsymbol{b}_2))$$
(32)

$$\frac{\partial J}{\partial \mathbf{W}_2} = (\mathbf{W}_2(\hat{\mathbf{y}} - \mathbf{y})^T)^T \tag{33}$$

$$= \boldsymbol{h}^T (\hat{\boldsymbol{y}} - \boldsymbol{y}) \tag{34}$$

$$\frac{\partial J}{\partial \mathbf{h}_2} = \hat{\mathbf{y}} - \mathbf{y} \tag{35}$$

$$d\mathbf{h} = d(\sigma(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)) = \sigma'(\cdot) \odot (\mathbf{x}d\mathbf{W}_1 + d\mathbf{b}_1)$$
(36)

$$dJ = trace((\hat{\boldsymbol{y}} - \boldsymbol{y})^T \sigma'(\cdot) \odot (\boldsymbol{x} d\boldsymbol{W}_1 + d\boldsymbol{b}_1) \boldsymbol{W}_2)$$
(37)

$$= trace(\mathbf{W}_2(\hat{\mathbf{y}} - \mathbf{y})^T \sigma'(\cdot) \odot (\mathbf{x} d\mathbf{W}_1 + d\mathbf{b}_1))$$
(38)

Apply law of trace and element – wise product: 
$$(39)$$

$$trace(A^{T}(B \odot C)) = trace((A \odot B)^{T}C)$$
(40)

$$= trace([(\hat{\boldsymbol{y}} - \boldsymbol{y})\boldsymbol{W}_2^T \odot \sigma'(\cdot)]^T (\boldsymbol{x} d\boldsymbol{W}_1 + d\boldsymbol{b}_1))$$
(41)

$$\frac{\partial J}{\partial \mathbf{W}_1} = \{ [(\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^T \odot \sigma'(\cdot)]^T \mathbf{x} \}^T$$
(42)

$$= \boldsymbol{x}^{T}[(\hat{\boldsymbol{y}} - \boldsymbol{y})\boldsymbol{W}_{2}^{T} \odot \sigma'(\cdot)] \tag{43}$$

$$\frac{\partial J}{\partial \mathbf{b}_1} = [(\hat{\mathbf{y}} - \mathbf{y}) \mathbf{W}_2^T \odot \sigma'(\cdot)] \tag{44}$$

# 3 word2vec

Notice:  $\nu$  and  $\mu$  are column vectors. Previously, all vectors are row vectors.

#### 3.1 a

$$J_{softmax-CE}(o, \nu_c, U) = CE(y, \hat{y})$$
(45)

$$\hat{\boldsymbol{y}} = softmax(\{\boldsymbol{U}_{\{d \times V\}}^T \times \boldsymbol{\nu}_{c\{d \times 1\}}\}^T)$$
(46)

$$dCE(\cdot) = (\hat{\boldsymbol{y}} - \boldsymbol{y})d\boldsymbol{\nu}_c^T \boldsymbol{U} \tag{47}$$

$$= trace[(\hat{\boldsymbol{y}} - \boldsymbol{y})d\boldsymbol{\nu}_c^T \boldsymbol{U}] \tag{48}$$

$$= trace[\boldsymbol{U}(\hat{\boldsymbol{y}} - \boldsymbol{y})d\boldsymbol{\nu}_c^T] \tag{49}$$

$$\frac{\partial CE(\cdot)}{\partial \boldsymbol{\nu}_{c}^{T}} = (\hat{\boldsymbol{y}} - \boldsymbol{y})^{T} \boldsymbol{U}^{T}$$
(50)

$$\frac{\partial CE(\cdot)}{\partial \nu_c} = U(\hat{y} - y) \tag{51}$$

### 3.2 b

$$dCE(\cdot) = (\hat{\mathbf{y}} - \mathbf{y})\boldsymbol{\nu}_c^T d\mathbf{U}$$
 (52)

$$= trace[(\hat{\boldsymbol{y}} - \boldsymbol{y})\boldsymbol{\nu}_c^T d\boldsymbol{U}] \tag{53}$$

$$\frac{\partial CE}{\partial U} = \nu_c (\hat{\boldsymbol{y}} - \boldsymbol{y})^T \tag{54}$$

#### 3.3 c

Assume that for  $\nu_c$ ,  $k_{neg}$  is a row vector of K-hot (ensure  $k_{neg\{o\}} \neq 1$ ) and k is one-hot vector that  $k_o = 1$ . Rewrite  $J_{neg-sample}(\cdot)$  as follow,

$$J_{neg-sample}(\cdot) = -\log(\sigma(\boldsymbol{\nu}_c^T \boldsymbol{U})) \boldsymbol{k}^T - \log(\sigma(-\boldsymbol{\nu}_c^T \boldsymbol{U})) \boldsymbol{k}_{neg}^T$$
(55)

$$= -\log(\sigma(\boldsymbol{\nu}_c^T \boldsymbol{U} \boldsymbol{k}^T)) - \log(\sigma(-\boldsymbol{\nu}_c^T \boldsymbol{U} \boldsymbol{k}_{neg}^T))$$
 (56)

In the last equation, we put all vectors into element-wise operator.

$$dJ_{neg-sample}(\cdot) = \frac{\sigma(\boldsymbol{\nu}_c^T \boldsymbol{U} \boldsymbol{k}^T)(1 - \sigma(\boldsymbol{\nu}_c^T \boldsymbol{U} \boldsymbol{k}^T))}{\sigma(\boldsymbol{\nu}_c^T \boldsymbol{U} \boldsymbol{k}^T)}$$
(57)

$$\frac{\partial J_{neg-sample}(\cdot)}{\partial \nu_c} = \tag{58}$$

TODO

$$\frac{\partial J}{\partial \boldsymbol{\nu}_c} = -\boldsymbol{\mu}_o \frac{\sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c)(1 - \sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c))}{\sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c)} + \sum_{k=1}^K \boldsymbol{\mu}_k (1 - \sigma(-\boldsymbol{\mu}_k^T \boldsymbol{\nu}_c))$$
 (59)

$$= -\boldsymbol{\mu}_o(1 - \sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c)) + \sum_{k=1}^K \boldsymbol{\mu}_k (1 - (1 - \sigma(\boldsymbol{\mu}_k^T \boldsymbol{\nu}_c))$$
 (60)

$$= -\boldsymbol{\mu}_o(1 - \sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c)) + \sum_{k=1}^K \boldsymbol{\mu}_k(\sigma(\boldsymbol{\mu}_k^T \boldsymbol{\nu}_c))$$
(61)

$$dJ(\cdot) = -\frac{d(\sigma(\cdot))}{\sigma(\cdot)} + \sum_{k=1}^{K} (1 - \sigma(-\boldsymbol{\mu}_k^T \boldsymbol{\nu}_c)) d\boldsymbol{\mu}_k^T \boldsymbol{\nu}_c$$
 (62)

$$= -trace((1 - \sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c))d\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c)$$
(63)

$$= trace(\boldsymbol{\nu}_c(\sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c) - 1)d\boldsymbol{\mu}_o^T)$$
(64)

$$\frac{\partial J}{\partial \boldsymbol{\mu}_o} = \boldsymbol{\nu}_c(\sigma(\boldsymbol{\mu}_o^T \boldsymbol{\nu}_c) - 1) \tag{65}$$

(66)

$$\frac{\partial J}{\partial \boldsymbol{\mu}_k} = \boldsymbol{\nu}_c \sigma(\boldsymbol{\mu}_k^T \boldsymbol{\nu}_c), for \ k = [1, 2, ..., K]$$
(67)

Computation efficiency of negtive-sampling loss, speed up ratio is approximately

$$\frac{V}{K+1}$$

#### 3.4 d

For skip-gram:

As the definition of word-vector, there are no items of  $\nu_k$  where  $k \neq c$  in cost function  $J(o, \nu_c, U)$  or  $F(\cdot)$ .

$$\frac{\partial J_{skip-gram}(\omega_{t-m}, ... \omega_{t+m})}{\partial \nu_{t}} = \mathbf{0}, for \ any \ k \neq c$$
 (68)

(69)

For the remaining parts is only applying the gradients derived by previous questions into the sum operator. The inner parts depend on the specified definition of F(o,c).

# For CBOW:

If softmax-cross-entropy cost,

$$p(\omega_{t}, \hat{\boldsymbol{\nu}}) = \frac{\exp(\boldsymbol{\mu}_{\omega_{t}}^{T} \sum_{-m \leq j \leq m, j \neq 0} \boldsymbol{\nu}_{\omega_{j}})}{\sum_{v=1}^{V} \exp(\boldsymbol{\mu}_{v}^{T} \sum_{-m \leq j \leq m, j \neq 0} \boldsymbol{\nu}_{\omega_{j}})}$$

$$J_{CBOW} = F(\omega_{t}, \hat{\boldsymbol{\nu}}) = CE(\boldsymbol{y}, \hat{\boldsymbol{y}} = softmax(\boldsymbol{U}_{\{d \times V\}}^{T} \times \hat{\boldsymbol{\nu}}_{\{d \times 1\}}^{T}))$$
(70)

$$J_{CBOW} = F(\omega_t, \hat{\boldsymbol{\nu}}) = CE(\boldsymbol{y}, \hat{\boldsymbol{y}} = softmax(\boldsymbol{U}_{\{d \times V\}}^T \times \hat{\boldsymbol{\nu}}_{\{d \times 1\}}^T))$$
(71)

$$dJ_{CBOW} = U(\hat{y} - y)d\hat{\nu}$$
(72)

$$= \sum_{-m \le j \le m, j \ne 0} \mathbf{U}(\hat{\mathbf{y}} - \mathbf{y}) d\mathbf{\nu}_{\omega_j}$$
 (73)