In the lectures notes, one has $\theta^{[l]}_{\text{here}} = \theta^{[l-1]}_{\text{there}}$, where there refers to the machine learning learning course by Andrew. We do this change to make the connexion with the deep learning specialization lectures.

# 1 Structure of dense layers

$$a^{[0]}_\mu \equiv \delta_{0\mu} + x_i \delta_{i\mu} \tag{1}$$

$$a^{[l]}_\mu \equiv \delta_{0\mu} + g^{[l]} \left( a^{[l-1]}_\nu \theta^{[l]}_{\nu i} \right) \, \delta^{[l]}_{i\mu} \tag{2}$$

$$z^{[l]}_i \equiv a^{[l-1]}_\nu \theta^{[l]}_{\nu i} \tag{3}$$

It should be understood that in $a^{[l]}_i$ one has $i \in \{1, ..., s_l\}$, and that in $a^{[l]}_\mu$ one has $\mu \in \{0, ..., s_l\}$. The binary cross entropy cross function is

$$cost = \frac{1}{m} \sum_n^m J \Big|_{x=x_n, y=y_n} \tag{4}$$

where

$$J \equiv - \sum_k^{s_L} \left( y_k \log \left( a^{[L]}_k \right) + (1 - y_k) \log \left( 1 - a^{[L]}_k \right) \right) \tag{5}$$

# 2 Backward propagation. Recursion relation

The following recursion relations do the job

$$\underbrace{\frac{\partial J}{\partial \theta^{[l]}_{j_{l-1} j_l}}}_{\mathrm{d}\theta^{[l]}_{j_{l-1} j_l}} = \underbrace{\frac{\partial z^{[l]}_{j_l}}{\partial \theta^{[l]}_{j_{l-1} j_l}}}_{a^{[l-1]}_{j_{l-1}}} \underbrace{\frac{\partial J}{\partial z^{[l]}_{j_l}}}_{\mathrm{d}z^{[l]}_{j_l}} \tag{6}$$

$$\underbrace{\frac{\partial J}{\partial z^{[l]}_{j_l}}}_{\mathrm{d}z^{[l]}_{j_l}} = \underbrace{\frac{\partial J}{\partial z^{[l+1]}_{j_{l+1}}} \frac{\partial z^{[l+1]}_{j_{l+1}}}{\partial a^{[l]}_{j_l}} \frac{\partial a^{[l]}_{j_l}}{\partial z^{[l]}_{j_l}}}_{\mathrm{d}a^{[l]}_{j_l}} = \underbrace{\frac{\partial J}{\partial z^{[l+1]}_{j_{l+1}}} \theta^{[l+1]}_{j_l j_{l+1}}}_{\mathrm{d}a^{[l]}_{j_l} = \theta^{[l+1]}_{j_l j_{l+1}} \mathrm{d}z^{[l+1]}_{j_{l+1}}} \, g^{[l]\prime} \left( z^{[l]}_{j_l} \right) \tag{7}$$

Hence, in an abbreviated notation one has:

$$\mathrm{d}\theta^{[l]}_{j_{l-1} j_l} = \mathrm{d}z^{[l]}_{j_l} a^{[l-1]}_{j_{l-1}} \tag{8}$$

$$\mathrm{d}a^{[l]}_{j_l} = \theta^{[l+1]}_{j_l j_{l+1}} \mathrm{d}z^{[l+1]}_{j_{l+1}} \tag{9}$$

$$\mathrm{d}z^{[l]}_{j_l} = \mathrm{d}a^{[l]}_{j_l} \, g^{[l]\prime} \left( z^{[l]}_{j_l} \right) \tag{10}$$

The initial condition is

$$\mathrm{d}\theta_{\mu j_L}^{[L]} = \frac{\partial J}{\partial \theta_{\mu j_L}^{[L]}} = \underbrace{\frac{(a_{j_L}^{[L]} - y_{j_L})}{a_{j_L}^{[L]}(1 - a_{j_L}^{[L-1]})}}_{\mathrm{d}a_{j_L}^{[L]}} \ g^{[L]'}\left(z_{j_L}^{[L]}\right) a_\mu^{[L-1]} \tag{11}$$

$$\mathrm{d}z_{j_L}^{[L]} = \mathrm{d}a_{j_L}^{[L]} \, g^{[L]'}\left(z_{j_L}^{[L]}\right) \tag{12}$$

$$\mathrm{d}\theta_{j_L\mu}^{[L]} = \mathrm{d}z_{j_L}^{[L]} \, a_\mu^{[L-1]} \tag{13}$$

This layer actually simplifies because one chooses $g^{[L]} = \sigma_{\mathrm{sigmoid}}$, which means that $\sigma'_{\mathrm{sigmoid}} = \sigma_{\mathrm{sigmoid}}(1 - \sigma_{\mathrm{sigmoid}})$.

## 3 Propagation with softmax

$$a_\mu^{[0]} \equiv \delta_{0\mu} + x_i \delta_{i\mu} \tag{14}$$

$$a_\mu^{[l]} \equiv \delta_{\mu 0} + \delta_{\mu i}^{[l]} \, g^{[l]}\left(a_\nu^{[l-1]}\theta_{\nu i}^{[l]}\right), \quad l < L \tag{15}$$

$$z_i^{[l]} \equiv a_\nu^{[l-1]}\theta_{\nu i}^{[l]}, \quad l \leq L \tag{16}$$

$$\hat{y}_i \equiv \frac{e^{z_i^{[L]}}}{\sum_c e^{z_c^{[L]}}} \tag{17}$$

It should be understood that in $a_i^{[l]}$ one has $i \in \{1, ..., s_l\}$, and that in $a_\mu^{[l]}$ one has $\mu \in \{0, ..., s_l\}$. The loss function is

$$J \equiv -\sum_i y_i \log\left(\hat{y}_i\right) \tag{18}$$

## 4 Back propagation

The following recursion relations do the job

$$\underbrace{\frac{\partial J}{\partial \theta_{j_{l-1}j_l}^{[l]}}}_{\mathrm{d}\theta_{j_{l-1}j_l}^{[l]}} = \underbrace{\frac{\partial z_{j_l}^{[l]}}{\partial \theta_{j_{l-1}j_l}^{[l]}}}_{a_{j_{l-1}}^{[l-1]}} \underbrace{\frac{\partial J}{\partial z_{j_l}^{[l]}}}_{\mathrm{d}z_{j_l}^{[l]}} \tag{19}$$

$$\underbrace{\frac{\partial J}{\partial z_{j_l}^{[l]}}}_{\mathrm{d}z_{j_l}^{[l]}} = \underbrace{\frac{\partial J}{\partial z_{j_{l+1}}^{[l+1]}} \frac{\partial z_{j_{l+1}}^{[l+1]}}{\partial a_{j_l}^{[l]}} \frac{\partial a_{j_l}^{[l]}}{\partial z_{j_l}^{[l]}}}_{\mathrm{d}a_{j_l}^{[l]}} = \underbrace{\frac{\partial J}{\partial z_{j_{l+1}}^{[l+1]}} \theta_{j_l j_{l+1}}^{[l+1]}}_{\mathrm{d}a_{j_l}^{[l]} = \theta_{j_l j_{l+1}}^{[l+1]} \mathrm{d}z_{j_{l+1}}^{[l+1]}} \; g^{[l]\prime}\left(z_{j_l}^{[l]}\right) \tag{20}$$

Hence, in an abbreviated notation one has:

$$\mathrm{d}\theta_{j_{l-1} j_l}^{[l]} = \mathrm{d}z_{j_l}^{[l]} a_{j_{l-1}}^{[l-1]} \tag{21}$$

$$\mathrm{d}a_{j_l}^{[l]} = \theta_{j_l j_{l+1}}^{[l+1]} \mathrm{d}z_{j_{l+1}}^{[l+1]} \tag{22}$$

$$\mathrm{d}z_{j_l}^{[l]} = \mathrm{d}a_{j_l}^{[l]} \, g^{[l]\prime}\left(z_{j_l}^{[l]}\right) \tag{23}$$

The initial condition is $(\hat{y}_i = g(z_i))$

$$\mathrm{d}\theta_{j_{L-1} j_L}^{[L]} = \frac{\partial J}{\partial \theta_{j_{L-1} j_L}^{[L]}} = \left(-\frac{y_i}{\hat{y}_i}\right)\left(\frac{\partial \hat{y}_i}{\partial z_{j_L}^{[L]}}\right)\left(\frac{\partial z_{j_L}^{[L]}}{\partial \theta_{j_{L-1} j_L}^{[L]}}\right) \tag{24}$$

$$\frac{\partial \hat{y}_i}{\partial z_{j_L}^{[L]}} = \hat{y}_i(\delta_{i j_L} - \hat{y}_{j_L}) \tag{25}$$

$$\frac{\partial z_{j_L}^{[L]}}{\partial \theta_{j_{L-1} j_L}^{[L]}} = a_{j_{L-1}}^{L-1} \tag{26}$$

Note that after summing over $i$, the first two terms in (24) simplify as:

$$\left(-\frac{y_i}{\hat{y}_i}\right)\left(\frac{\partial \hat{y}_i}{\partial z_{j_L}^{[L]}}\right) = -y_i + \hat{y}_i \tag{27}$$

****************************************************