

- **Task 1 (25 p):** Show that

$$\delta^{[l]} = g^{[l]'}(Z^{[l]}) \odot \left(W^{[l+1]T} * \delta^{[l+1]} \right), \quad l = 2, \dots, L-1 \quad (i)$$

and

$$\delta^{[L]} = J'(A^{[L]}) \odot g^{[L]'}(Z^{[L]}) \quad (ii)$$

where \odot denotes element-wise multiplication and $*$ denotes ordinary matrix multiplication. Here $g^{[L]'}(Z^{[L]})$ and $J'(A^{[L]})$ denote the derivatives of the univariate functions applied element-wise to the vectors. If any operations with tensors are used, these need to be clearly defined.

$$\begin{aligned}
 (i) \quad \frac{\partial J}{\partial z^{[0]}} &= \frac{\partial J}{\partial z^{[L]}} \frac{\partial z^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L-1]}}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial g^{[L-2]}} \dots \frac{\partial z^{[L-1]}}{\partial g^{[2]}} \frac{\partial g^{[2]}}{\partial z^{[0]}} \\
 &= \frac{\partial J}{\partial g^{[2]}} = \frac{\partial J}{\partial z^{[2+1]}} \frac{\partial z^{[2+1]}}{\partial g^{[2]}} = g^{[2]'}(z^{[2]}) \\
 &= \delta^{[2+1]} \frac{\partial}{\partial g^{[2]}} [W^{[2+1]} g^{[2]}(z^{[2]}) + b] \\
 &= \delta^{[2+1]} W^{[2+1]}
 \end{aligned}$$

$$\Rightarrow \delta^{[2]} = \frac{\partial g^{[2]}}{\partial z^{[2]}} \odot \left(W^{[2+1]T} * \delta^{[2+1]} \right)$$



$$(ii) \quad \text{show } \delta^{[L]} = \frac{\partial J}{\partial z^{[L]}} = \frac{\partial J}{\partial g^{[L]}} \frac{\partial g^{[L]}}{\partial z^{[L]}}$$



simple application of chain rule.

- **Task 2 (10 p):** Derive the corresponding expression for the gradients under the following setting:

- a) – Cost function: Mean Squared Error (MSE)
- b) – Activation function for hidden layers: Rectified Linear Unit (ReLU)
- c) – Activation function for the output layer: Identity function

$$a) \text{ MSE}(\theta) = \frac{1}{m} \sum_{i=1}^m \|Y - \hat{Y}\|_2^2 = \left\{ \begin{array}{l} \text{from task (c) we know } \hat{Y} = Z^{[L]} \\ \Rightarrow \frac{1}{m} \sum_{i=1}^m \|Y - Z^{[L]}\|_2^2 \end{array} \right\}$$

$$\frac{\partial J}{\partial Z^{[L]}} = \frac{\partial J}{\partial Z^{[L]}} \left[\frac{1}{m} \sum_{i=1}^m \sum_{p=1}^D (y_i - z_i^{[L]})^2 \right] = -\frac{2}{m} \sum_{p=1}^D (y_i - z_i^{[L]}) \quad (*)$$

$$\nabla_{\theta} J = \frac{\partial J}{\partial \theta} = \left(\frac{\partial J}{\partial w^{[0]}}, \dots, \frac{\partial J}{\partial w^{[L]}}, \frac{\partial J}{\partial b^{[0]}}, \dots, \frac{\partial J}{\partial b^{[L]}} \right)$$

$$\begin{aligned} (i) \quad \frac{\partial J}{\partial w^{[L]}} &= \frac{\partial J}{\partial Z^{[L]}} \frac{\partial Z^{[L]}}{\partial w^{[L]}} = \left(\frac{\partial g^{[L]}}{\partial Z^{[L]}} \right)^T \left(W^{[L+1]T} \delta^{[L+1]} \right) \cdot \frac{\partial Z^{[L]}}{\partial w^{[L]}} \\ &= \delta^{[L]} \cdot g^{[L-1]}(Z^{[L-1]}) \end{aligned} \quad Z^{[L]} = W^{[L]} A^{[L-1]} + b^{[L]}$$

$$(ii) \quad \frac{\partial J}{\partial b^{[L]}} = \frac{\partial J}{\partial Z^{[L]}} \frac{\partial Z^{[L]}}{\partial b^{[L]}} = \frac{\partial g^{[L]}}{\partial Z^{[L]}} \left(W^{[L+1]T} \delta^{[L+1]} \right) = \delta^{[L]}$$

$$(iii) \quad \frac{\partial J}{\partial w^{[L-1]}} = \frac{\partial J}{\partial Z^{[L]}} \frac{\partial Z^{[L]}}{\partial w^{[L-1]}} = \frac{\partial J}{\partial Z^{[L]}} \cdot g^{[L-1]}(Z^{[L-1]}) = \delta^{[L]} \cdot g^{[L-1]}(Z^{[L-1]})$$

$$(iv) \quad \frac{\partial J}{\partial b^{[L-1]}} = \frac{\partial J}{\partial Z^{[L]}} \frac{\partial Z^{[L]}}{\partial b^{[L-1]}} = \dots = \delta^{[L]}$$



$$b) \frac{\partial g}{\partial z^{[0]}} = \frac{\partial}{\partial z^{[0]}} [\max(0, z^{[0]})] = \max(0, 1) = 1$$

$$c) g(z^{[L]}) = z^{[L]} \text{ identity}$$

plug into
(i)-(iv) to
get final
gradient

- **Task 3 (5 p):** How would you avoid an exponential blowup of computation when computing the gradients?

Use back-propagation.

Refer to page 205, equations (6.49)-(6.52) in the course literature, and also what we did in the previous task. Several parts need only be computed once and then stored in memory.

$$\frac{\partial z}{\partial w} \quad (6.49)$$

$$= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \quad (6.50)$$

$$= f'(y) f'(x) f'(w) \quad (6.51) \quad \leftarrow \text{we compute this}$$

$$= f'(f(f(w))) f'(f(w)) f'(w). \quad (6.52) \quad \leftarrow \text{not this.}$$

exponential
blowup.