**1**

**(1.1)** $\nabla_{\hat{g}} L(\hat{y}, y) = \nabla_{\hat{g}}\left(-y \ln(\hat{g}) - (1-y)\ln(1-\hat{g})\right)$

$$= \frac{-y}{\hat{g}} - \frac{(-1)(1-y)}{(1-\hat{g})} = \frac{-y}{\hat{g}} + \frac{(1-y)}{(1-\hat{g})} = \frac{-y(1-\hat{g}) + \hat{g}(1-y)}{\hat{g}(1-\hat{g})}$$

$$= \frac{-y + \hat{g}y + \hat{g} - \hat{g}y}{\hat{g}(1-\hat{g})} = \frac{-y + \hat{g}}{\hat{g}(1-\hat{g})} \qquad \text{and since } h^{(2)} = \hat{g},$$

$$\boxed{\nabla_{\hat{g}} L(\hat{y}, y) = \frac{-y + \hat{g}}{\hat{g}(1-\hat{g})} = \frac{-y + h^{(2)}}{h^{(2)}(1-h^{(2)})}}$$

**(1.2)** $\nabla_{a^{(2)}} J = g \odot f'(a^{(2)}) \qquad f'(z) = \frac{e^{-z}}{(1+e^{-z})^2} = \frac{(1+e^{-z}-1)}{(1+e^{-z})} \cdot \frac{1}{(1+e^{-z})} = (1-f(z))f(z)$

so $\quad g \odot f'(a^{(2)}) = g \odot f(a^{(2)})(1-f(a^{(2)}))$

$$= g \odot h^{(2)}(1-h^{(2)}) \quad \text{since } h^{(2)} = f(a^{(2)})$$

$$= \underbrace{\frac{-y + h^{(2)}}{h^{(2)}(1-h^{(2)})}}_{g \text{ from } 1.1} \odot h^{(2)}(1-h^{(2)}) \qquad \text{so } \boxed{\nabla_{a^{(2)}} J = h^{(2)} - y}$$

**(1.3)** $\nabla_{b^{(2)}} J = g$

from 1.2, $\quad g = \nabla_{a^{(2)}} J \quad$ so $\quad \boxed{\nabla_{b^{(2)}} J = h^{(2)} - y}$

**(1.4)** $\nabla_{W^{(2)}} J = g\, h^{(2-1)T} = g\, h^{(1)T} \quad$ so $\quad \boxed{\nabla_{W^{(2)}} J = (h^{(2)} - y) \cdot h^{(1)T}}$

**(1.5)** $\nabla_{h^{(1)}} J = W^{(2)T} g$

$g = h^{(2)} - y \quad$ so $\quad \boxed{\nabla_{h^{(1)}} J = W^{(2)T} \cdot (h^{(2)} - y)}$

(1.6) $\quad g \leftarrow \nabla_{a^{(1)}} J = g \odot f'(a^{(1)})$

$$= g \odot f(a^{(1)})(1 - f(a^{(1)})) = g \odot h^{(1)}(1 - h^{(1)})$$

from 1.4, we had: $g \leftarrow \nabla_{h^{(1)}} J = W^{(2)T}(h^{(2)} - y)$

so $\quad g \odot h^{(1)}(1 - h^{(1)}) = W^{(2)T}(h^{(2)} - y) \odot h^{(1)}(1 - h^{(1)}) = \nabla_{a^{(1)}} J = g$

$\nabla_{b^{(1)}} J = g \qquad$ so $\quad \boxed{\nabla_{b^{(1)}} J = W^{(2)T}(h^{(2)} - y) \odot h^{(1)}(1 - h^{(1)})}$

$\nabla_{W^{(1)}} J = g\, h^{(1-1)T} \quad$ so $\quad \boxed{\nabla_{W^{(1)}} J = \left[ W^{(2)T}(h^{(2)} - y) \odot h^{(1)}(1 - h^{(1)}) \right] \cdot h^{(0)T}}$

where $\underline{h^{(0)} = x}$, $\underline{h^{(1)} = \sigma(b^{(1)} + W^{(1)} x)}$, $\underline{h^{(2)} = \sigma(b^{(2)} + W^{(2)} \sigma(b^{(1)} + W^{(1)} x))}$

---

2.2.2 $\quad f = \dfrac{1}{1 + e^{-\omega^T x}} = (1 + e^{-\omega^T x})^{-1}$

$\dfrac{df}{dx} = (-\omega^T e^{-\omega^T x}) \cdot -1 (1 + e^{-\omega^T x})^{-2} \quad$ so $\quad \boxed{\dfrac{df}{dx} = \dfrac{\omega^T e^{-\omega^T x}}{(1 + e^{-\omega^T x})^2}}$

$\dfrac{df}{d\omega} = (-x e^{-\omega^T x}) \cdot -1 (1 + e^{-\omega^T x})^{-2} \quad$ so $\quad \boxed{\dfrac{df}{d\omega} = \dfrac{x e^{-\omega^T x}}{(1 + e^{-\omega^T x})^2}}$

$\omega^T x = \begin{bmatrix} 2 & -3 & -3 \end{bmatrix} \begin{bmatrix} -1 \\ -2 \\ 1 \end{bmatrix} = -2 + 6 - 3 = 1$

so $\quad \dfrac{\omega^T e^{-\omega^T x}}{(1 + e^{-\omega^T x})^2} = \dfrac{\omega^T e^{-1}}{(1 + e^{-1})^2} = \dfrac{.367879}{1.87109} \begin{bmatrix} 2 \\ -3 \\ -3 \end{bmatrix} = \begin{bmatrix} 0.3932 \\ -0.5898 \\ -0.5898 \end{bmatrix} \checkmark$

and $\quad \dfrac{x e^{-\omega^T x}}{(1 + e^{-\omega^T x})^2} = \dfrac{x e^{-1}}{(1 + e^{-1})^2} = \dfrac{0.367879}{1.87109} \begin{bmatrix} -1 \\ -2 \\ 1 \end{bmatrix} = \begin{bmatrix} -0.1966 \\ -0.3932 \\ 0.1966 \end{bmatrix} \checkmark$