

Q2)

$$L(\hat{y}, y) = -y \ln(\hat{y}) - (1-y) \ln(1-\hat{y})$$

$$\begin{aligned} 2.1) \hat{y} \nabla L(\hat{y}, y) &= -\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})} \\ &= \frac{-y(1-\hat{y}) + \hat{y}(1-y)}{\hat{y}(1-\hat{y})} \\ &= \frac{-y + y\hat{y} + \hat{y} - y\hat{y}}{\hat{y}(1-\hat{y})} = \frac{\hat{y} - y}{\hat{y}(1-\hat{y})} = \boxed{\frac{h_2 - y}{h_2(1-h_2)}} = g \end{aligned}$$

$$\begin{aligned} 2.2) \nabla_{a^2} J &= g \odot f'(a_2) \\ &= \frac{h_2 - y}{h_2(1-h_2)} \cdot \cancel{h_2 \cdot (1-h_2)} \quad \begin{array}{l} f'(a_2) = f(a_2)(1-f(a_2)) = h_2 \cdot (1-h_2) \\ \hookrightarrow f(x) \text{ is sigmoid function} \\ f(a_2) = h_2 \end{array} \\ &= \boxed{h_2 - y} = g \end{aligned}$$

$$\begin{aligned} 2.3) \nabla_b^2 J &= g + \lambda \nabla_{w^2} J^2(\theta) \rightarrow \lambda \text{ is zero} \\ \nabla_b^2 J &= h_2 - y = g \end{aligned}$$

$$\begin{aligned} 2.4) \nabla_{w^2} J &= g \cdot h^{(k-1)T} + 0 \rightarrow k=2 \\ &= (h_2 - y) h^{(1)T} \end{aligned}$$

$$2.5) \nabla_{h^1} J = w^{(k)T} \cdot g = w^{2T} (h_2 - y) \rightarrow \text{new } g$$

$$\begin{aligned} 2.6) \nabla_{b^1} J &= g \cdot f(a_1)(1-f(a_1)) \\ &= w^{2T} \cdot (h_2 - y) \odot h_1(1-h_1) \Rightarrow \text{new } g \\ &\quad \rightarrow k=1 \end{aligned}$$

$$\nabla_{w_1} J = g \cdot h^{(k-1)T} + \cancel{\lambda \nabla_{w^2} J^2(\theta)} \rightarrow w^{2T} (h_2 - y) h_1(1-h_1) \cdot h_0^T$$

2.7) Backpropagation would be much slower w/o forward pass. As beginning weights have the most function implemented like $g(h(f(kx)))$ So getting weights requires a lot of computations in chain rule. Forward pass allows you to compute gradient from the end, so gradient calculation considers only previous calculated gradient, which is much faster.