1. (a) (i) $\delta_x^{(1)} = \nabla_{\theta_x} l_x(x; \omega^{(2)}, \omega_x^{(1)})$

$$= \nabla_{\theta_x}\left(\tfrac{1}{2}\|x\omega^{(2)}\omega_x^{(1)} - x\|_2^2\right)$$

$\Rightarrow$ In terms of $\theta_x \Rightarrow \nabla_{\theta_x}\left(\tfrac{1}{2}\|\theta_x - x\|_2^2\right)$, $\theta_x = h\omega_x^{(1)}$, $h = x\omega^{(2)}$

$\Rightarrow \tfrac{1}{2}\nabla_{\theta_x}\left[(\theta_x - x)^T(\theta_x - x)\right] = \tfrac{1}{2}\nabla_{\theta_x}\left[(\theta_x^T - x^T)(\theta_x - x)\right]$

$= \tfrac{1}{2}\nabla_{\theta_x}\left[\theta_x^T\theta_x - \theta_x^T x - x^T\theta_x + \cancel{x^T x}^{\,0}\right]$

$= \tfrac{1}{2}\left[2\theta_x - \theta_x - x^T\right] = \boxed{\tfrac{1}{2}(\theta_x - x^T) \in \mathbb{R}^{d\times 1}}$

(ii) $\nabla_{\omega_y^{(1)}} l(x, y; \omega^{(2)}, \omega_x^{(1)}, \omega_y^{(1)}) = \nabla_{\omega_y^{(1)}}(l_y(\text{softmax}(x\omega^{(2)}\omega_y^{(1)}), y)) + 0$

$\hookrightarrow$ since $l_x$ does not have $\omega_y^{(1)}$ $\nabla l_x = 0$

$\nabla_y^{(1)} l_y = \nabla_{\omega_y^{(1)}}\theta_y \nabla_{\theta_y} l_y = h\delta_y^{(1)} = \boxed{x\omega^{(2)}(\text{softmax}(x\omega^{(2)}\omega_y^{(1)}) - y) \in \mathbb{R}^{p\times m}}$


(iii) $\nabla_h\, l_y(\text{softmax}(h\omega_y^{(1)}), y) + \beta\, l_x(h\omega_x^{(1)}, x)$

$= (\nabla_h\theta_y)(\text{softmax}(h\omega_y^{(1)}) - y) + \beta(\nabla_h\theta_x)\left(\tfrac{1}{2}(\theta_x - x^T)\right)$

$= \boxed{\omega_y^{(1)}(\text{softmax}(h\omega_y^{(1)}) - y) + \tfrac{\beta}{2}\omega_x^{(1)}(h\omega_x^{(1)} - x^T)) \in \mathbb{R}^{p\times 1}}$

(iv)

$\nabla_{\omega^{(2)}} l = \nabla_{\omega^{(2)}} h \cdot \nabla_h(l)$

$= \boxed{x(\omega_y^{(1)}(\text{softmax}(h\omega_y^{(1)}) - y) + \tfrac{\beta}{2}\omega_x^{(1)}(h\omega_x^{(1)} - x^T)) \in \mathbb{R}^{p\times m}}$

b) When $\beta = 0$, $\hat{y} = \text{softmax}(x\omega^{(2)}\omega_y^{(1)})$ is equivalent to $\hat{y} = \sigma(x\omega)$

which ends up being a linear classifier.

c)


d)