1.

$$y = \phi_0 + \phi_1 a\Big[\psi_{01} + \psi_{11}a[\theta_{01} + \theta_{11}x] + \psi_{21}a[\theta_{02} + \theta_{12}x]\Big]$$
$$+\phi_2 a\Big[\psi_{02} + \psi_{12}a[\theta_{01} + \theta_{11}x] + \psi_{22}a[\theta_{02} + \theta_{12}x]\Big], \qquad (7.35)$$

$\dfrac{\partial y}{\partial \phi_0} = 1$

$\dfrac{\partial y}{\partial \phi_1} = a[\psi_{01} + \psi_{11} a[\theta_{01} + \theta_{11}x] + \psi_{21}a[\theta_{02} + \theta_{12}x]]$

similar for $\phi_2$

$\dfrac{\partial y}{\partial \psi_{01}} = \phi_1 \mathbb{1}[\psi_{01} + \psi_{11} a[\theta_{01} + \theta_{11}x] + \psi_{21}a[\theta_{02} + \theta_{12}x] > 0]$

similar for $\psi_{02}$

$\dfrac{\partial y}{\partial \psi_{11}} = \phi_1 \, a[\theta_{01} + \theta_{11}x]\mathbb{1}[\psi_{01} + \psi_{11} a[\theta_{01} + \theta_{11}x] + \psi_{21} a[\theta_{02} + \theta_{12}x] > 0]$

similar for $\psi_{12}, \psi_{21}, \psi_{22}$

$\dfrac{\partial y}{\partial \theta_{01}} = \phi_1 \psi_{11}\mathbb{1}[\theta_{01} + \theta_{11}x > 0]\mathbb{1}[\psi_{01} + \psi_{11} a[\theta_{01} + \theta_{11}x] + \psi_{21}a[\theta_{02} + \theta_{12}x] > 0]$
$\qquad + \phi_2 \psi_{12}\mathbb{1}[\theta_{01} + \theta_{11}x > 0]\mathbb{1}[\psi_{02} + \psi_{12} a[\theta_{01} + \theta_{11}x] + \psi_{22}a[\theta_{02} + \theta_{12}x] > 0]$

similar for $\theta_{02}$

$\dfrac{\partial y}{\partial \theta_{11}} = \phi_1 \psi_{11} x \mathbb{1}[\theta_{01} + \theta_{11} > 0]\mathbb{1}[\psi_{01} + \psi_{11} a[\theta_{01} + \theta_{11}x] + \psi_{21}a[\theta_{02} + \theta_{12}x] > 0]$
$\qquad + \phi_2 \psi_{12} x \mathbb{1}[\theta_{01} + \theta_{11}x]\mathbb{1}[\psi_{02} + \psi_{12} a[\theta_{01} + \theta_{11}x] + \psi_{22}a[\theta_{02} + \theta_{12}x] > 0]$

similar for $\theta_{12}$

3.

$$\frac{\partial \ell_i}{\partial f_0} = \frac{\partial h_1}{\partial f_0}\frac{\partial f_1}{\partial h_1}\left(\frac{\partial h_2}{\partial f_1}\frac{\partial f_2}{\partial h_2}\frac{\partial h_3}{\partial f_2}\frac{\partial f_3}{\partial h_3}\frac{\partial \ell_i}{\partial f_3}\right). \qquad (7.20)$$

$\dfrac{\partial \ell_i}{\partial f_3} : \; D_f \times 1$

$\dfrac{\partial f_3}{\partial h_3} : \; D_3 \times D_f$

$\dfrac{\partial h_3}{\partial f_2} : \; D_3 \times D_3$

$\dfrac{\partial f_2}{\partial h_2} : D_2 \times D_3$

$\dfrac{\partial h_2}{\partial f_1} : \; D_2 \times D_2$

$\dfrac{\partial f_1}{\partial h_1} : \; D_1 \times D_2$

$\dfrac{\partial h_1}{\partial f_0} : \; D_1 \times D_1$

5. $l_i = -(1-y_i) \log(1-\text{sig}[f(x_i \cdot \phi]) - y_i \log[\text{sig}[f(x_i, \phi]]$

$$\frac{\partial \text{sig}(z)}{\partial z} = \frac{e^{-z}}{(1+e^{-z})^2} = \text{sig}(z)(1-\text{sig}(z))$$

$$\frac{\partial}{\partial f}\left(\log[1-\text{sig}f]\right) = \frac{-1}{1-\text{sig}[f]} \quad \frac{\text{sig}(f)(1-\text{sig}(f))}{}$$

$$= -\text{sig}[z]$$

$$\frac{\partial}{\partial f}\left(\log[\text{sig} f]\right) = \frac{1}{\text{sig}[f]} \quad \text{sig}[f](1-\text{sig}[f])$$

$$= 1-\text{sig}[f]$$

$$\frac{\partial l_i}{\partial f} = (1-y_i)\,\text{sig}[f] - y_i(1-\text{sig}[f])$$
$$= \text{sig}[f] - y_i$$

6. $z = \beta + \Omega h$

$h \to (D_i \times 1)$
$z \to (D_0 \times 1)$
$\Omega \to (D_0 \times D_i)$

$$\frac{\partial z}{\partial h} = \begin{pmatrix} \frac{\partial z_1}{\partial h_1} & \frac{\partial z_2}{\partial h_1} & \cdots \frac{\partial z_j}{\partial h_1} \\ \frac{\partial z_1}{\partial h_2} & & \\ & \frac{\partial z_j}{\partial h_i} & \\ \frac{\partial z_1}{\partial h_i} & \cdots & \frac{\partial z_{D_0}}{\partial h_{D_i}} \end{pmatrix}$$

But, $z_j = \beta_j + \sum_i^{D_i} w_{ji} h_i$

so $\quad \frac{\partial z_j}{\partial h_i} = w_{ji}$

which makes up the transpose of $\Omega$.

$$\frac{\partial z}{\partial h} = \begin{pmatrix} w_{11} & w_{21} & \cdots w_{j1} \cdots w_{D_0 1} \\ w_{12} & & \\ w_{ii} & & w_{ji} \\ w_{1 D_i} & & \cdots \quad w_{D_0 D_i} \end{pmatrix} = \Omega^T$$

7. $\text{sig}[f] = \dfrac{1}{1+e^{-f}} = h$

$\dfrac{\partial h}{\partial f} = \text{sig}[f]\,[1-\text{sig}[f]]$

$= \dfrac{1}{1+e^{-f}}\left(1 - \dfrac{1}{1+e^{-f}}\right)$

The gradients vanish in either case



the gradients get smaller and smaller.

9. $\ell[f]$

$\underline{f} = \underline{\beta} + \underline{\Omega}\,\underline{h}$

$f_i = \beta_i + \sum_j \Omega_{ij}\, h_j$

$\dfrac{\partial f_i}{\partial \Omega_{ij}} = h_j \quad\to\quad \dfrac{\partial \underline{f}}{\partial \underline{\Omega}} = \underline{h}^T$

$\dfrac{\partial \ell}{\partial \underline{\Omega}} = \dfrac{\partial \ell}{\partial \underline{f}}\dfrac{\partial \underline{f}}{\partial \underline{\Omega}} = \dfrac{\partial \ell}{\partial \underline{f}}\,\underline{h}^T$

10. For regular ReLu, $\underline{h}' = a[\underline{f}]$, $\underline{f} = \underline{\beta} + \underline{\Omega}\,\underline{h}$

$\dfrac{\partial}{\partial f}\, a[f] = \mathbb{1}[f>0] \in$ indicator function.

Now, instead of the derivative becoming $0$, it is a small constant $\alpha$.

so $\dfrac{\partial \underline{h}'}{\partial \underline{h}} = \mathbb{1}[\underline{f}>0]\,\underline{\Omega}^T + \mathbb{1}[\underline{f}<0]\,\alpha\,\underline{\Omega}^T$

14. $\text{Var}[q] = \mathbb{E}[q^2] - \mathbb{E}[q]$

$\therefore \mathbb{E}[q^2] = \delta^2$.

$\mathbb{E}[b^2] = \mathbb{E}[q^2 \cdot \mathbb{1}(q > 0)]$

Because $q$ is centred around the mean, probability of $q$ being positive is $\frac{1}{2}$. $\mathbb{E}[b] = 0$, because $\mathbb{E}[q] = 0$. Hence $\mathbb{E}[b] = \text{Var}[b^2]$

$\mathbb{E}[b^2] = \frac{1}{2}\mathbb{E}[q^2] = \frac{1}{2}\delta^2$.

15. The network won't train well because.

eg.

$$\frac{\partial \ell_i}{\partial f_0} = \frac{\partial h_1}{\partial f_0}\frac{\partial f_1}{\partial h_1}\left(\frac{\partial h_2}{\partial f_1}\frac{\partial f_2}{\partial h_2}\frac{\partial h_3}{\partial f_2}\frac{\partial f_3}{\partial h_3}\frac{\partial \ell_i}{\partial f_3}\right).$$

Network with 3 hidden layers.

$\frac{\partial f_i}{\partial h_j} = \frac{\partial}{\partial h_j}[\beta_i + \underline{\beta}_i \underline{h}_i]$

$= \underline{\beta}_i^T$

If the weights are all initialised to zero then the gradients of the loss will be zero too.