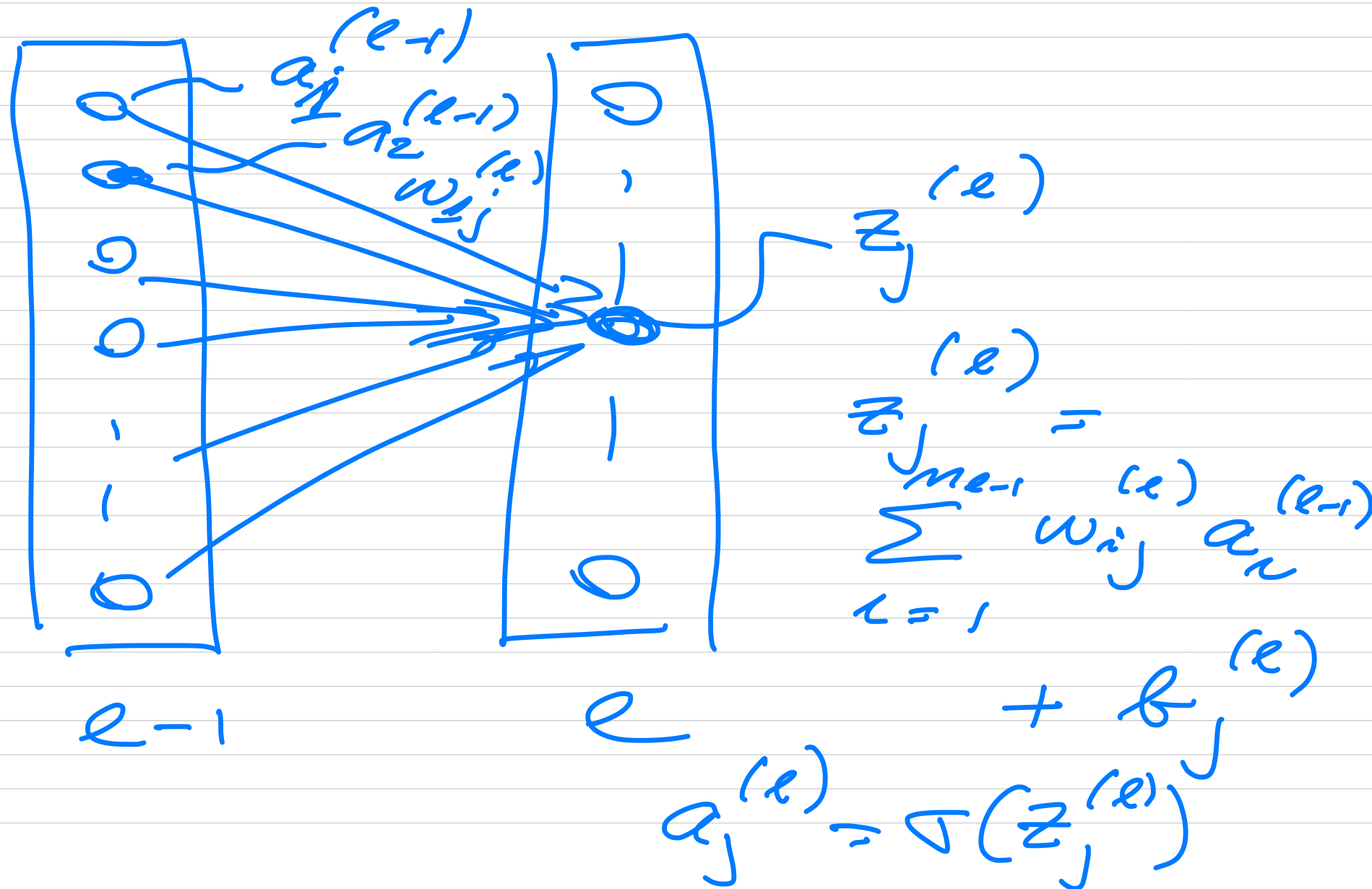


FYS-STK3155/4155, lecture  
October 13, 2025

# FYS-STK3155/4155 October 13



$$l = L$$

$$a_j^{(l)} \rightarrow C(a_j^{(l)}, y_j, \theta)$$

$$\frac{MSE}{C(\theta)} = \frac{1}{2} \sum_{i=1}^n (a_j^{(l)} - y_j)^2$$

$$\frac{\partial C}{\partial \theta} = 0$$

$$\frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} = a_i^{(l-1)}$$

$$\frac{\partial z_j^{(l)}}{\partial a_k^{(l-1)}} = w_{kj}^{(l)}$$

$$\frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \stackrel{\uparrow}{=} a_j^{(l)} (1 - a_j^{(l)})$$

Sigmoid  
for  $\nabla$

start with  $l = L$  (final)

$$C(\theta) = \frac{1}{2} \sum_i (a_i^{(L)} - y_i)^2$$

$$\frac{\partial C}{\partial w_{jk}^{(L)}} = \underbrace{(a_j^{(L)} - y_j)}_{\text{From } C(\theta)} \leftarrow \frac{\partial}{\partial a_j^{(L)}}$$

$$\times \underbrace{\left[ a_j^{(L)} (1 - a_j^{(L)}) \right]}_{\sigma(z_j^{(L)})} a_k^{(L-1)}$$

$$\delta_j^{(L)} = \frac{\partial C}{\partial a_j^{(L)}} \Delta_j^{(L)} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}}$$

$$\frac{\partial C}{\partial w_{jk}}^{(L)} = \delta_j^{(L)} a_k^{(L-1)}$$

$$\frac{\partial C}{\partial b_j^{(L)}} = \delta_j^{(L)}$$

$$w_{jk}^{(L)} \leftarrow w_{jk}^{(L)} - \eta \delta_j^{(L)} a_k^{(L-1)}$$

$$b_j^{(L)} \leftarrow b_j^{(L)} - \eta \delta_j^{(L)}$$

$$\begin{aligned} L &\rightarrow L \\ \delta_j^{(L)} &= \frac{\partial C}{\partial z_j^{(L)}} \left[ \delta_k^{(L+1)} \right] \\ &= \sum_k \left( \frac{\partial C}{\partial z_k^{(L+1)}} \right) \left( \frac{\partial z_k^{(L+1)}}{\partial z_j^{(L)}} \right) \end{aligned}$$

$$\bar{z}_j^{(l+1)} = \sum_{i=1}^{N_l} w_{ij}^{(l+1)} a_i^{(l)} + b_j^{(l+1)}$$

$$a_i^{(l)} = \Delta(\bar{z}_i^{(l)})$$

$$\frac{\partial \bar{z}_k^{(l+1)}}{\partial \bar{z}_j^{(l)}} = w_{kj}^{(l+1)} \Delta'(\bar{z}_j^{(l)})$$

$$\delta_j^{(l)} = \sum_k \delta_k^{(l+1)} w_{kj}^{(l+1)} \Delta'(\bar{z}_j^{(l)})$$



updates of gradients

$$w_{jk}^{(l)} \leftarrow w_{jk}^{(l)} - \eta \delta_j^{(l)} a_k^{(l-1)}$$

$$f_j^{(l)} \leftarrow f_j^{(l)} - \eta \underbrace{\frac{\partial C}{\partial b_j^{(l)}}}_{\delta_j^{(l)}}$$

- algorithm

- Define architecture  
(model)

- # nodes

- # hidden layers

- activation functions

- cost/loss function

- initialize  $\Theta = \{W, b\}$

- set up design matrix

X

- learning rate, hyper-parameter
- gradient methods

(i) Perform first FF step and continue from hidden layer  $l=1$  to  $l=L$   
 $q^{(L)}$

- compute  $\delta_j^{(c)}$

- Then backpropagate

$l = L-1, l = L-2, \dots$

$l = 1$

compute  $\delta_j^{(c)}$

Train gradients, update

$$w_{jk}^{(c)} \leftarrow w_{jk}^{(c)} - \eta \delta_j^{(c)} a_k^{(c)}$$

$$b_j^{(c)} \leftarrow b_j^{(c)} - \eta \delta_j^{(c)}$$

- continue till convergence