

## Module 4.6: Backpropagation: Computing Gradients w.r.t. Hidden Units

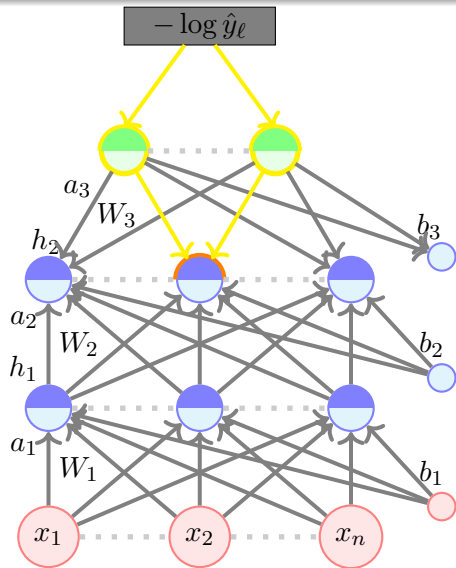
## Quantities of interest (roadmap for the remaining part):

- Gradient w.r.t. output units
- Gradient w.r.t. hidden units
- Gradient w.r.t. weights and biases

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

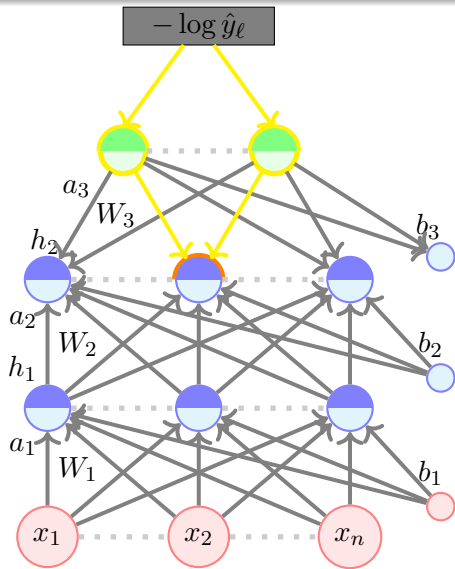
- Our focus is on *Cross entropy loss* and *Softmax* output.

**Chain rule along multiple paths:** If a function  $p(z)$  can be written as a function of intermediate results  $q_i(z)$  then we have :



**Chain rule along multiple paths:** If a function  $p(z)$  can be written as a function of intermediate results  $q_i(z)$  then we have :

$$\frac{\partial p(z)}{\partial z} = \sum_m \frac{\partial p(z)}{\partial q_m(z)} \frac{\partial q_m(z)}{\partial z}$$

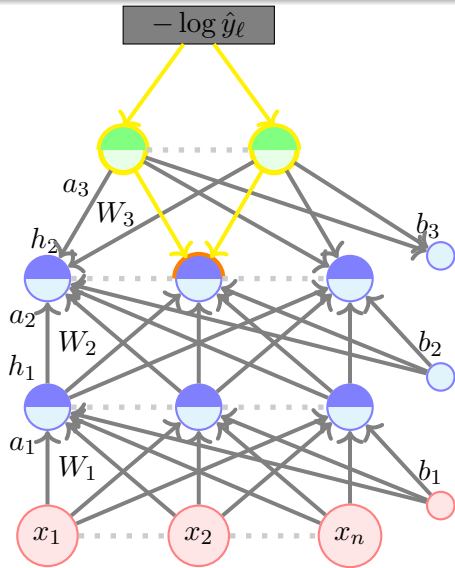


**Chain rule along multiple paths:** If a function  $p(z)$  can be written as a function of intermediate results  $q_i(z)$  then we have :

$$\frac{\partial p(z)}{\partial z} = \sum_m \frac{\partial p(z)}{\partial q_m(z)} \frac{\partial q_m(z)}{\partial z}$$

In our case:

- $p(z)$  is the loss function  $\mathcal{L}(\theta)$

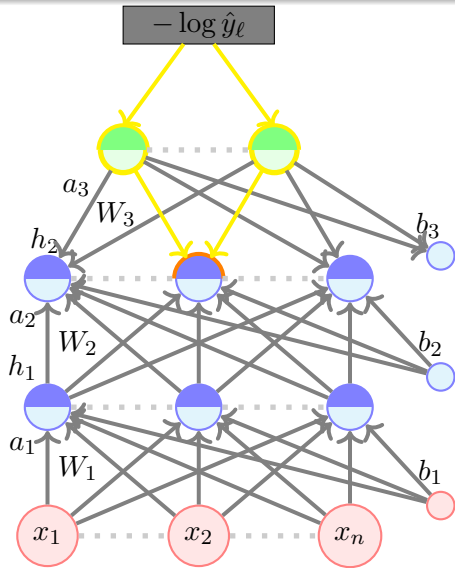


**Chain rule along multiple paths:** If a function  $p(z)$  can be written as a function of intermediate results  $q_i(z)$  then we have :

$$\frac{\partial p(z)}{\partial z} = \sum_m \frac{\partial p(z)}{\partial q_m(z)} \frac{\partial q_m(z)}{\partial z}$$

In our case:

- $p(z)$  is the loss function  $\mathcal{L}(\theta)$
- $z = h_{ij}$

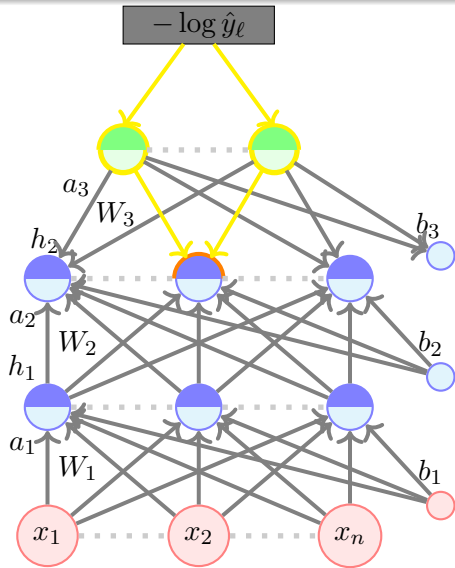


**Chain rule along multiple paths:** If a function  $p(z)$  can be written as a function of intermediate results  $q_i(z)$  then we have :

$$\frac{\partial p(z)}{\partial z} = \sum_m \frac{\partial p(z)}{\partial q_m(z)} \frac{\partial q_m(z)}{\partial z}$$

In our case:

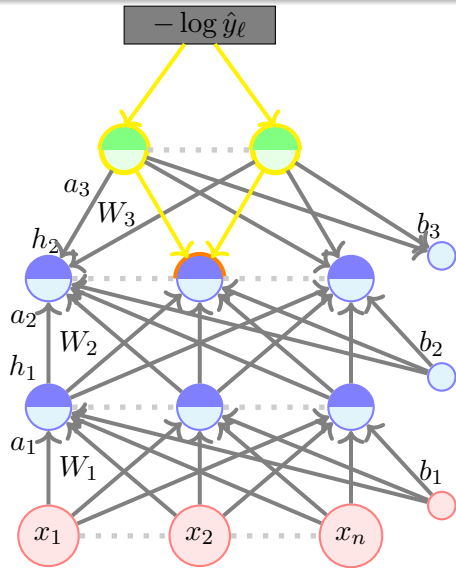
- $p(z)$  is the loss function  $\mathcal{L}(\theta)$
- $z = h_{ij}$
- $q_m(z) = a_{Lm}$



Intentionally left blank

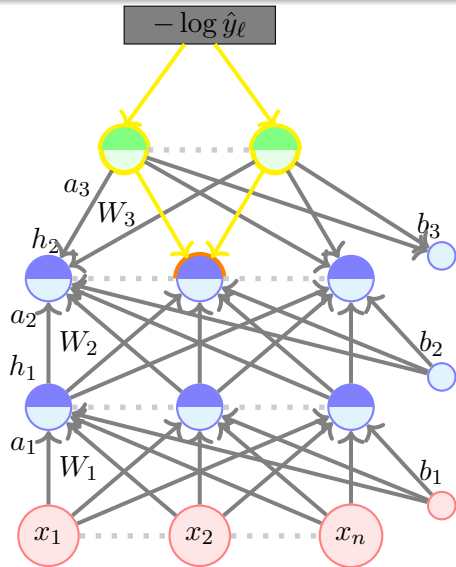


$$\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}}$$



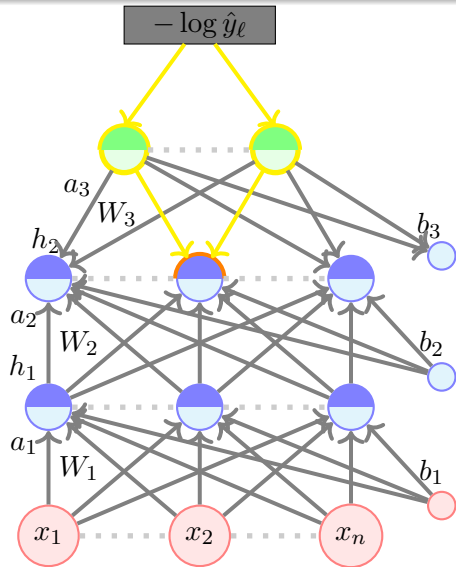
$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}}$$



$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

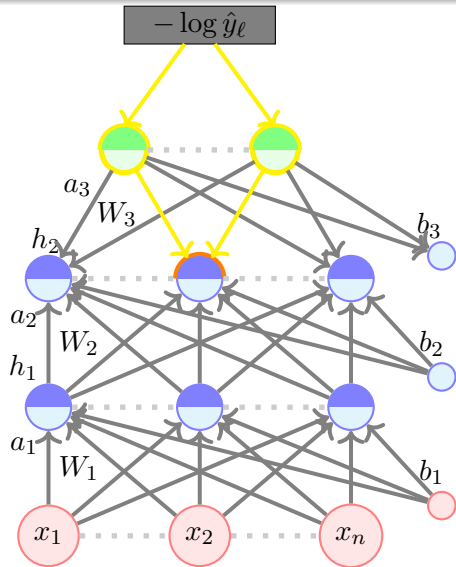
$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$



$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

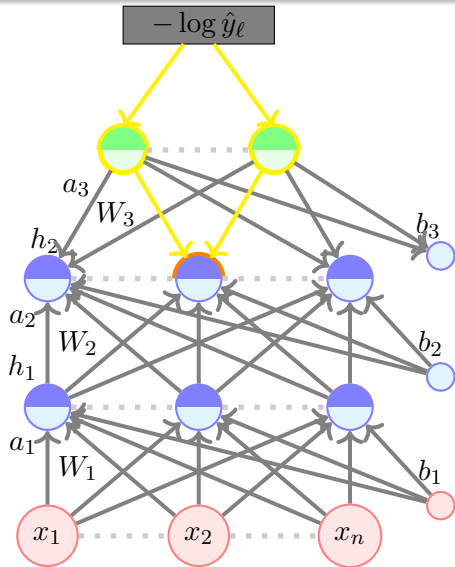


$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \phantom{0} \\ \phantom{0} \\ \phantom{0} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} \phantom{0} \\ \phantom{0} \\ \phantom{0} \end{bmatrix}$$

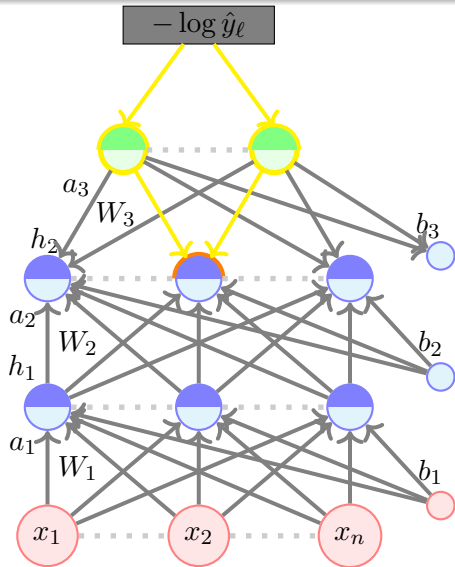


$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

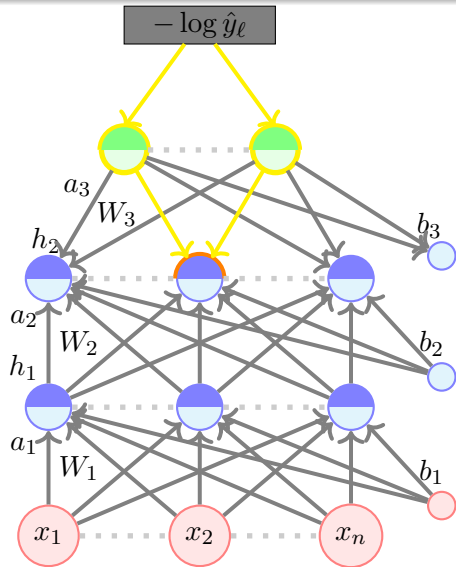


$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

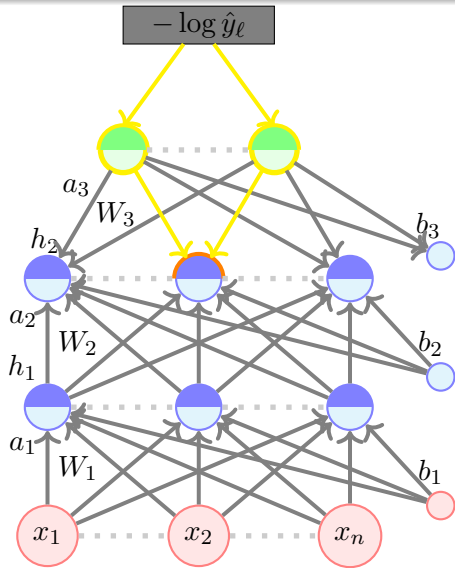


$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \end{bmatrix}$$



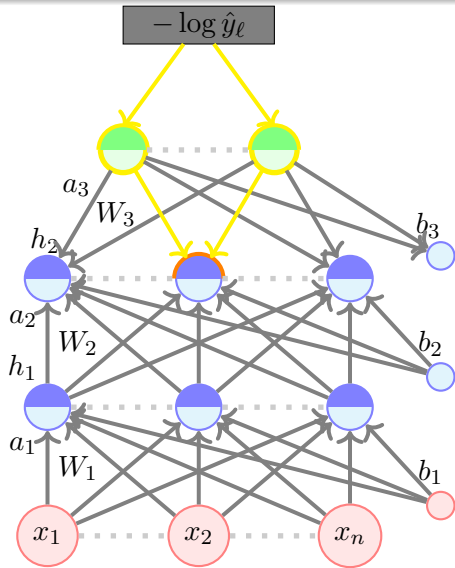
$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$



$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \end{bmatrix}$$

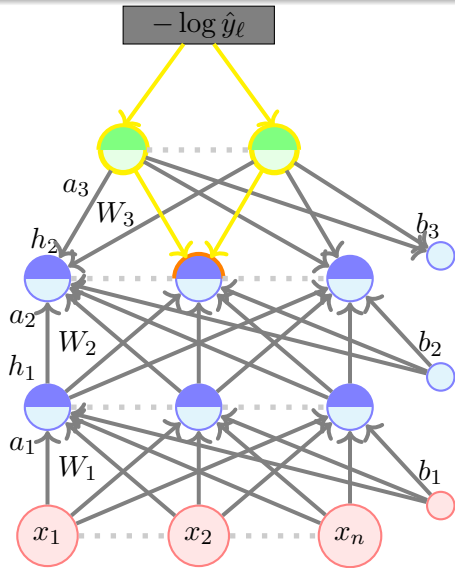


$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

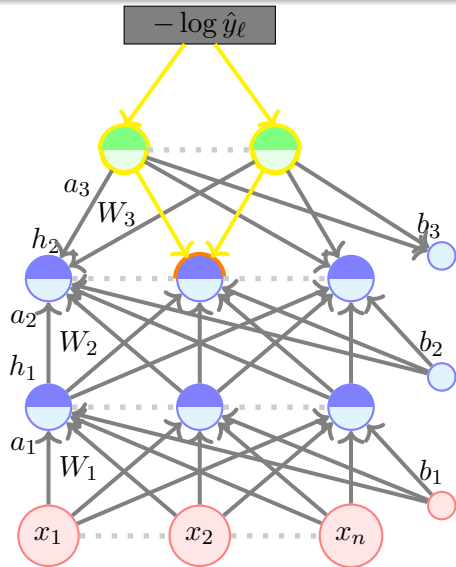


$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$



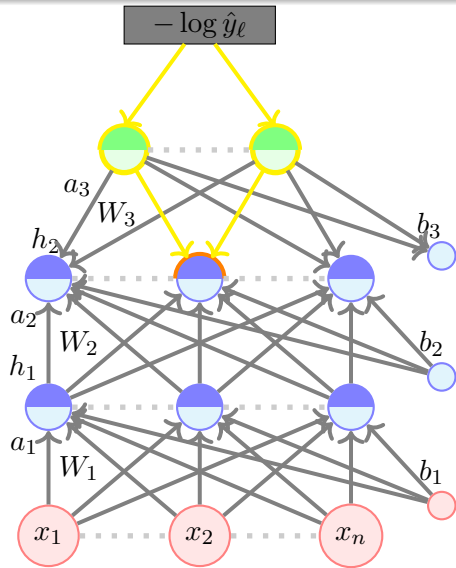
$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

$W_{i+1, \cdot, j}$  is the  $j$ -th column of  $W_{i+1}$ ;



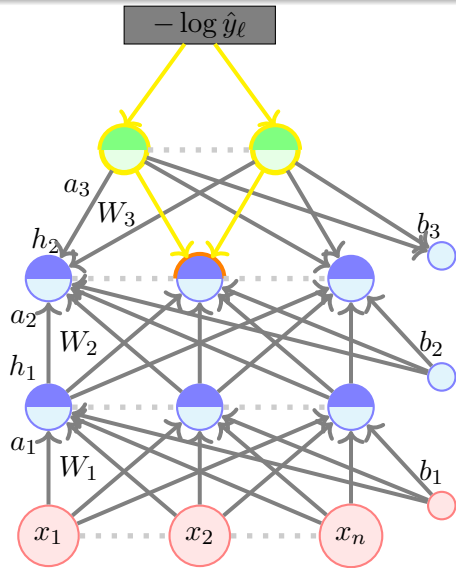
$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

$W_{i+1, \cdot, j}$  is the  $j$ -th column of  $W_{i+1}$ ; see that,



$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

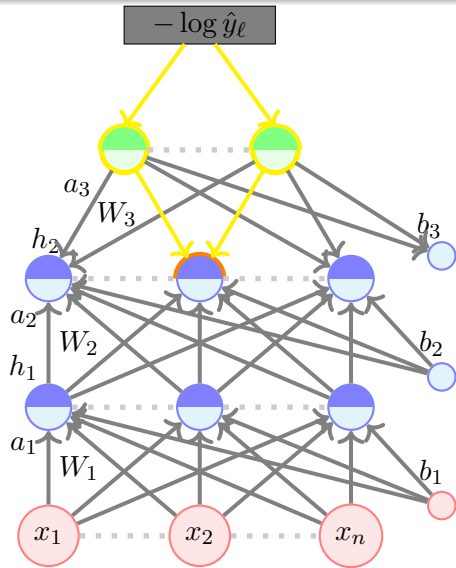
$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

$W_{i+1, \cdot, j}$  is the  $j$ -th column of  $W_{i+1}$ ; see that,

$$(W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) =$$



$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

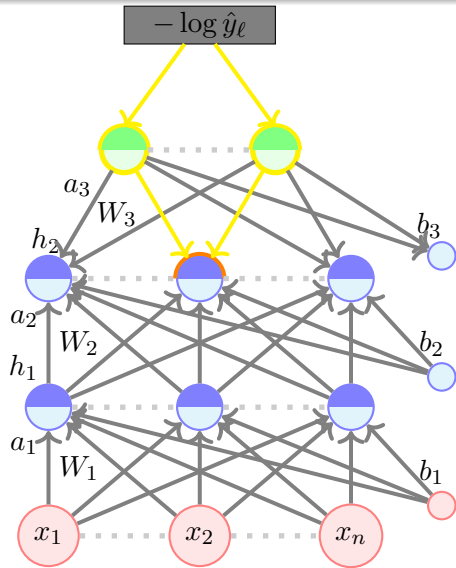
$$\begin{aligned}\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} \frac{\partial a_{i+1,m}}{\partial h_{ij}} \\ &= \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}\end{aligned}$$

Now consider these two vectors,

$$\nabla_{a_{i+1}} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,k}} \end{bmatrix}; W_{i+1, \cdot, j} = \begin{bmatrix} W_{i+1,1,j} \\ \vdots \\ W_{i+1,k,j} \end{bmatrix}$$

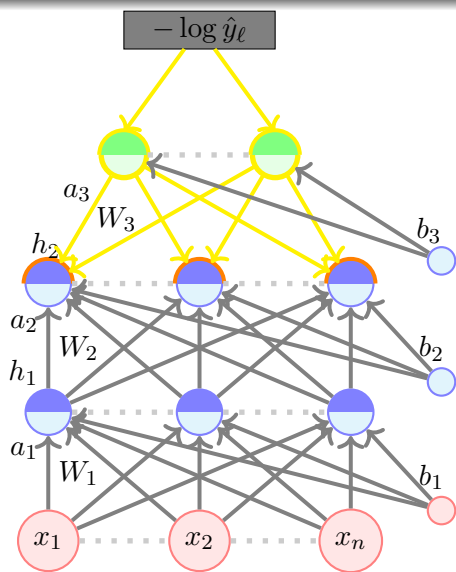
$W_{i+1, \cdot, j}$  is the  $j$ -th column of  $W_{i+1}$ ; see that,

$$(W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) = \sum_{m=1}^k \frac{\partial \mathcal{L}(\theta)}{\partial a_{i+1,m}} W_{i+1,m,j}$$



$$a_{i+1} = W_{i+1} h_{ij} + b_{i+1}$$

We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

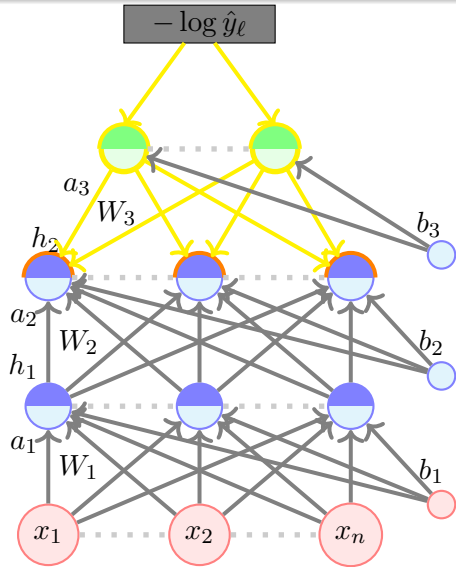




We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1,.,j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

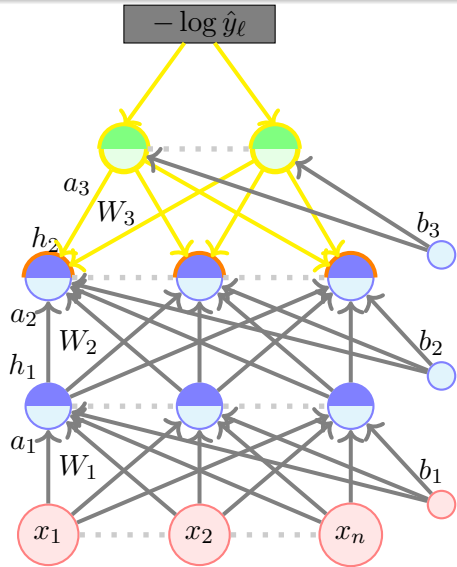
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta)$$



$$\text{We have, } \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$$

We can now write the gradient w.r.t.  $h_i$

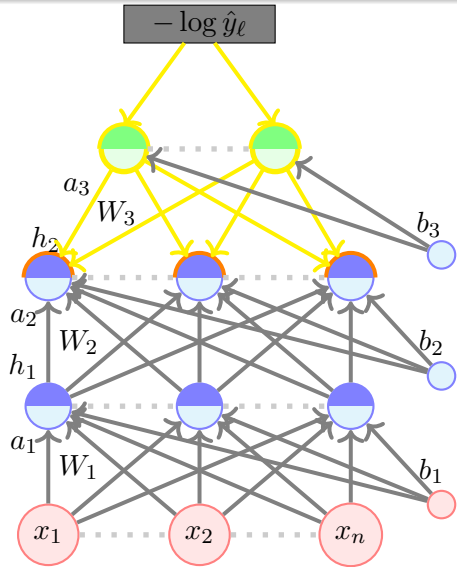
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \\ \\ \end{bmatrix} = \begin{bmatrix} \\ \\ \end{bmatrix}$$



We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1,.,j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

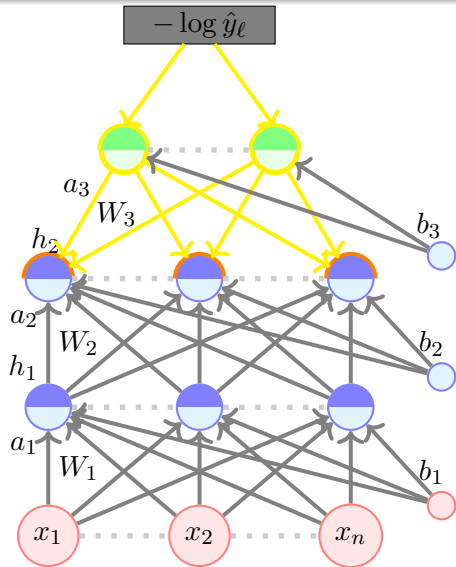
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$



We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

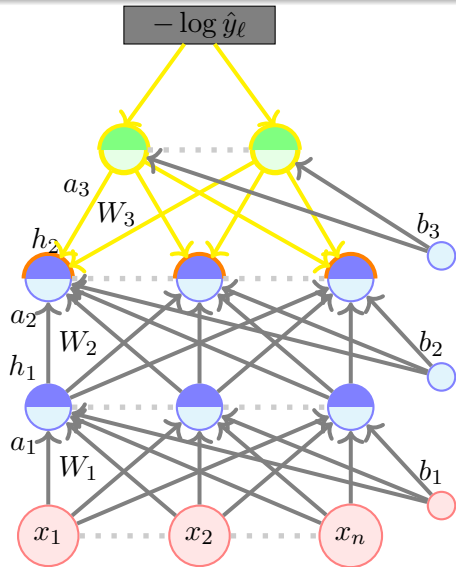
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \end{bmatrix}$$



We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

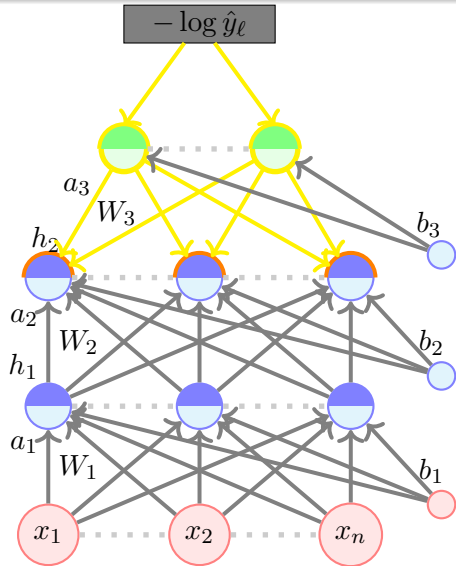
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \end{bmatrix}$$



We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

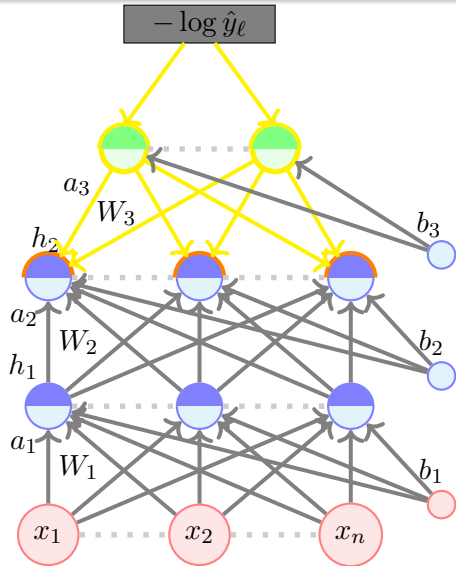
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \end{bmatrix}$$



We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

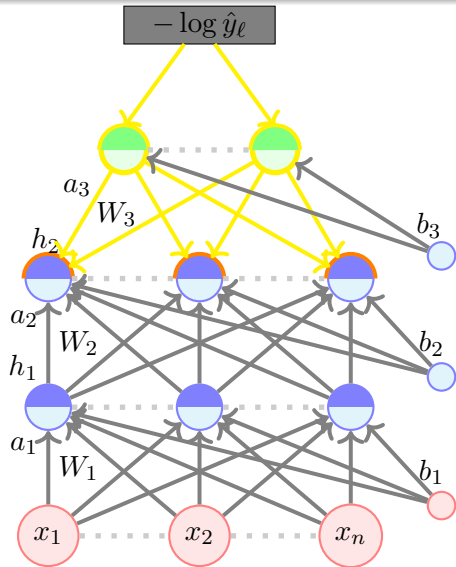
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \\ \vdots \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ \vdots \end{bmatrix}$$



We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ \vdots \end{bmatrix}$$

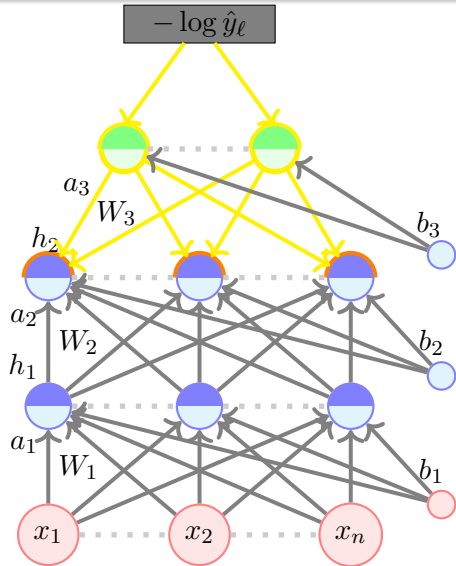




We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

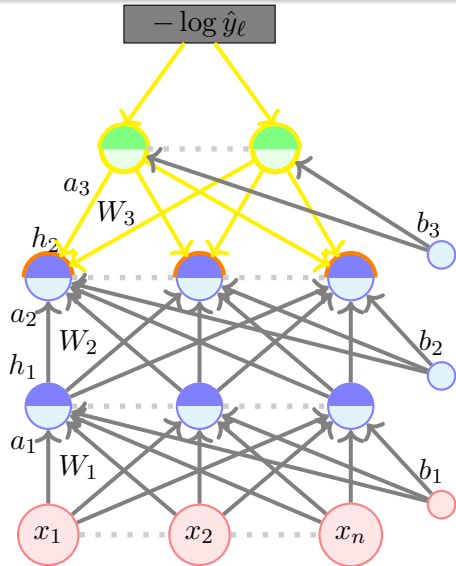
$$\nabla_{\mathbf{h}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ \vdots \\ (W_{i+1, \cdot, n})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \end{bmatrix}$$



We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

$$\begin{aligned} \nabla_{\mathbf{h}_i} \mathcal{L}(\theta) &= \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ \vdots \\ (W_{i+1, \cdot, n})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \end{bmatrix} \\ &= (W_{i+1})^T (\nabla_{a_{i+1}} \mathcal{L}(\theta)) \end{aligned}$$

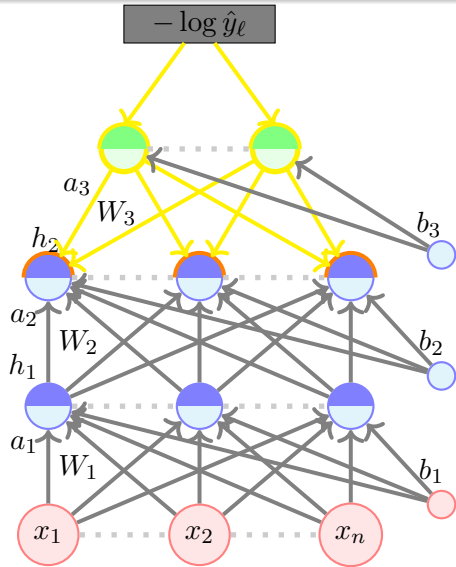


We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

We can now write the gradient w.r.t.  $h_i$

$$\begin{aligned} \nabla_{\mathbf{h}_i} \mathcal{L}(\theta) &= \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ \vdots \\ (W_{i+1, \cdot, n})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \end{bmatrix} \\ &= (W_{i+1})^T (\nabla_{a_{i+1}} \mathcal{L}(\theta)) \end{aligned}$$

- We are almost done except that we do not know how to calculate  $\nabla_{a_{i+1}} \mathcal{L}(\theta)$  for  $i < L-1$

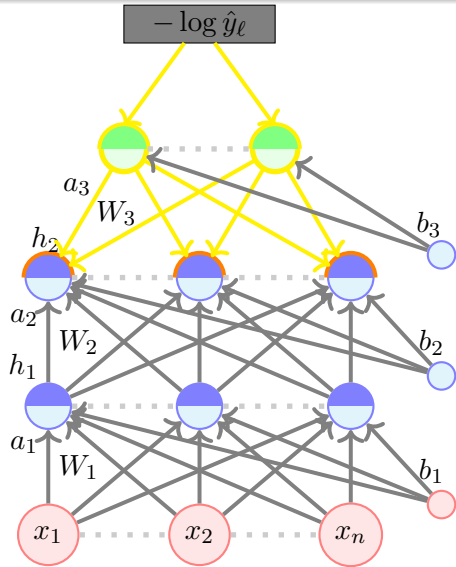


We have,  $\frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} = (W_{i+1, \cdot, j})^T \nabla_{a_{i+1}} \mathcal{L}(\theta)$

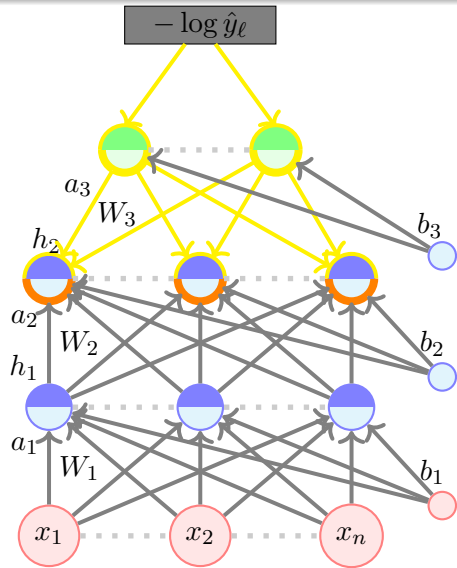
We can now write the gradient w.r.t.  $h_i$

$$\begin{aligned} \nabla_{\mathbf{h}_i} \mathcal{L}(\theta) &= \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} \end{bmatrix} = \begin{bmatrix} (W_{i+1, \cdot, 1})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ (W_{i+1, \cdot, 2})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \\ \vdots \\ (W_{i+1, \cdot, n})^T \nabla_{a_{i+1}} \mathcal{L}(\theta) \end{bmatrix} \\ &= (W_{i+1})^T (\nabla_{a_{i+1}} \mathcal{L}(\theta)) \end{aligned}$$

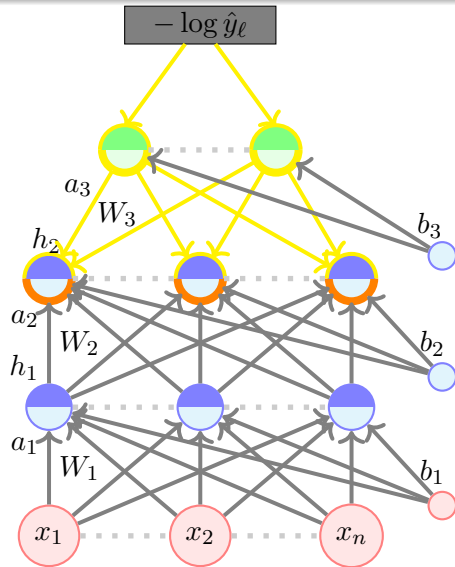
- We are almost done except that we do not know how to calculate  $\nabla_{a_{i+1}} \mathcal{L}(\theta)$  for  $i < L-1$
- We will see how to compute that



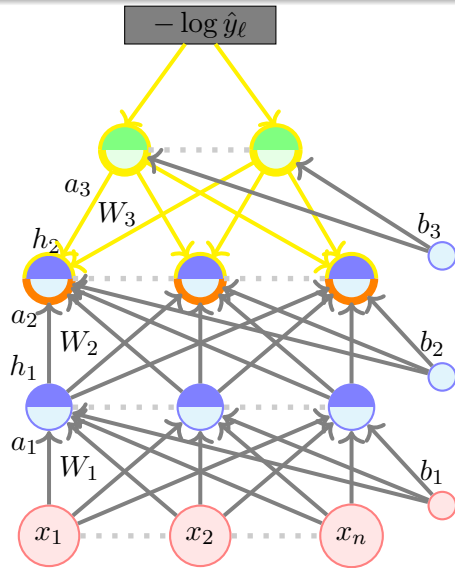
$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta)$$



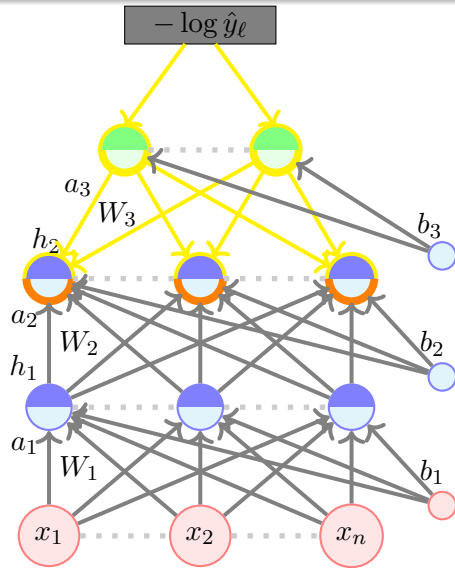
$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \phantom{0} \\ \phantom{0} \\ \phantom{0} \end{bmatrix}$$



$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \end{bmatrix}$$

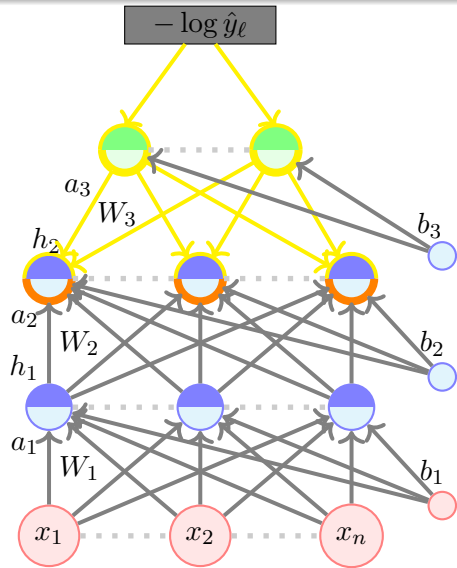


$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \end{bmatrix}$$



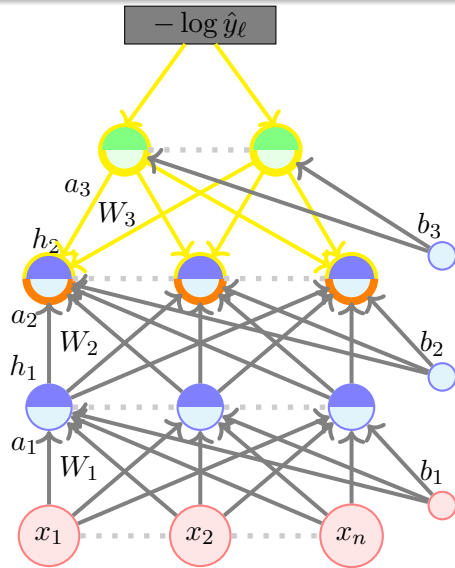


$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$



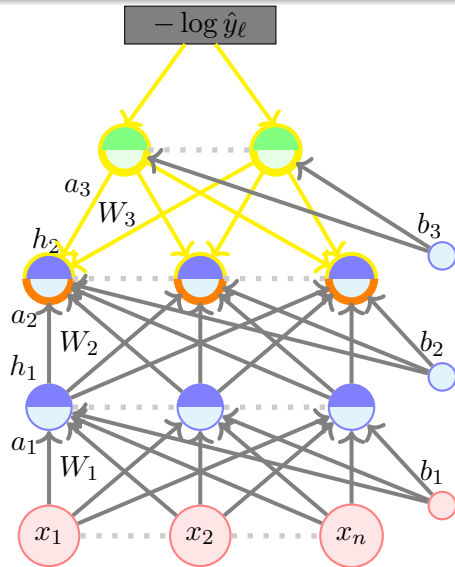
$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}}$$



$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

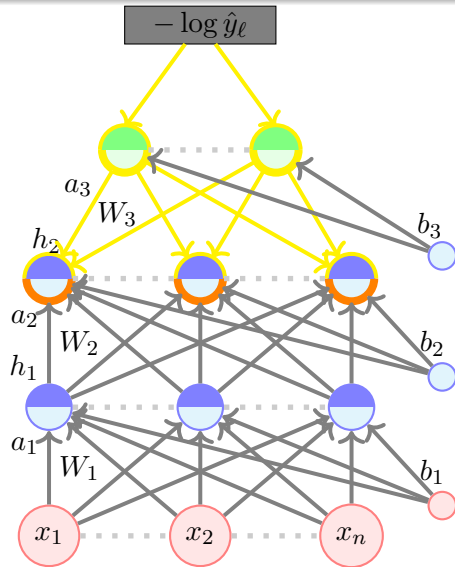
$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$



$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$

$$= \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})]$$

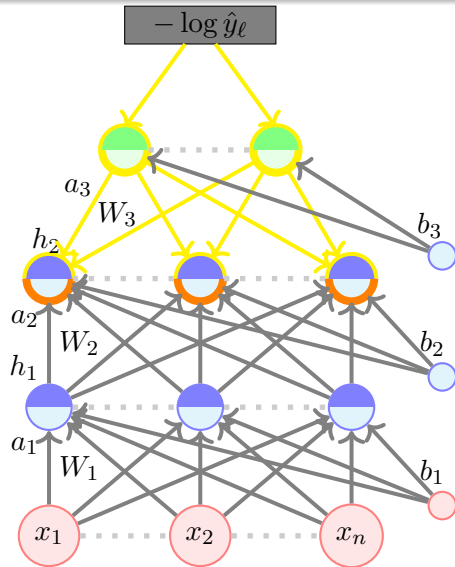


$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$

$$= \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})]$$

$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta)$$

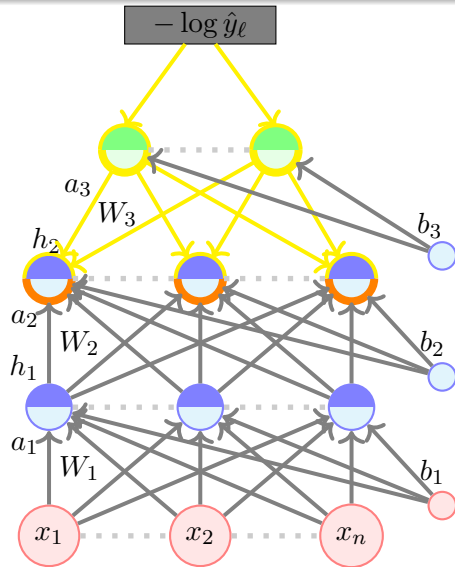


$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$

$$= \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})]$$

$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \phantom{\frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}}} \\ \phantom{\frac{\partial \mathcal{L}(\theta)}{\partial h_{i2}}} \\ \phantom{\frac{\partial \mathcal{L}(\theta)}{\partial h_{i3}}} \end{bmatrix}$$

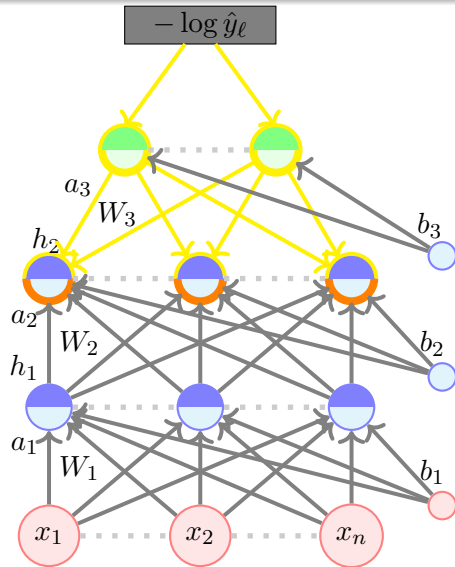


$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$

$$= \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})]$$

$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} g'(a_{i1}) \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} g'(a_{in}) \end{bmatrix}$$

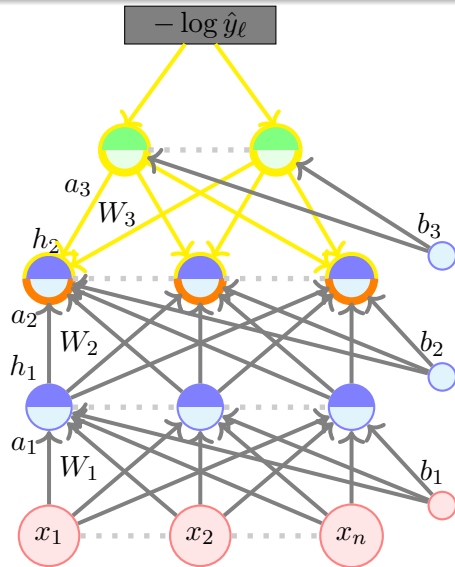


$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$

$$= \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})]$$

$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} g'(a_{i1}) \\ \vdots \end{bmatrix}$$



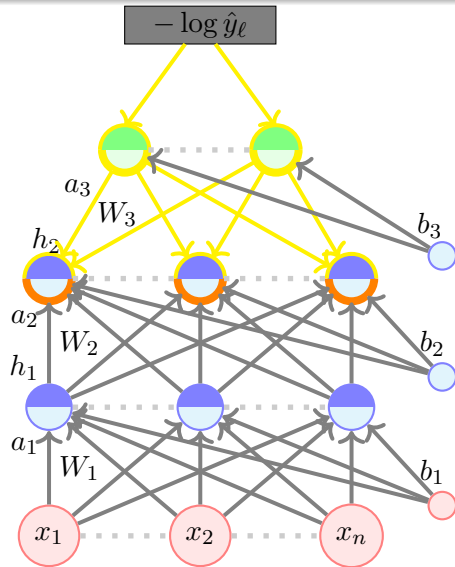


$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$

$$= \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})]$$

$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} g'(a_{i1}) \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} g'(a_{in}) \end{bmatrix}$$



$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial a_{i1}} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial a_{in}} \end{bmatrix}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial a_{ij}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial a_{ij}}$$

$$= \frac{\partial \mathcal{L}(\theta)}{\partial h_{ij}} g'(a_{ij}) \quad [\because h_{ij} = g(a_{ij})]$$

$$\nabla_{\mathbf{a}_i} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial h_{i1}} g'(a_{i1}) \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial h_{in}} g'(a_{in}) \end{bmatrix}$$

$$= \nabla_{h_i} \mathcal{L}(\theta) \odot [\dots, g'(a_{ik}), \dots]$$

