

## 11 Neural Network

### 11.1 Feedforward neural networks

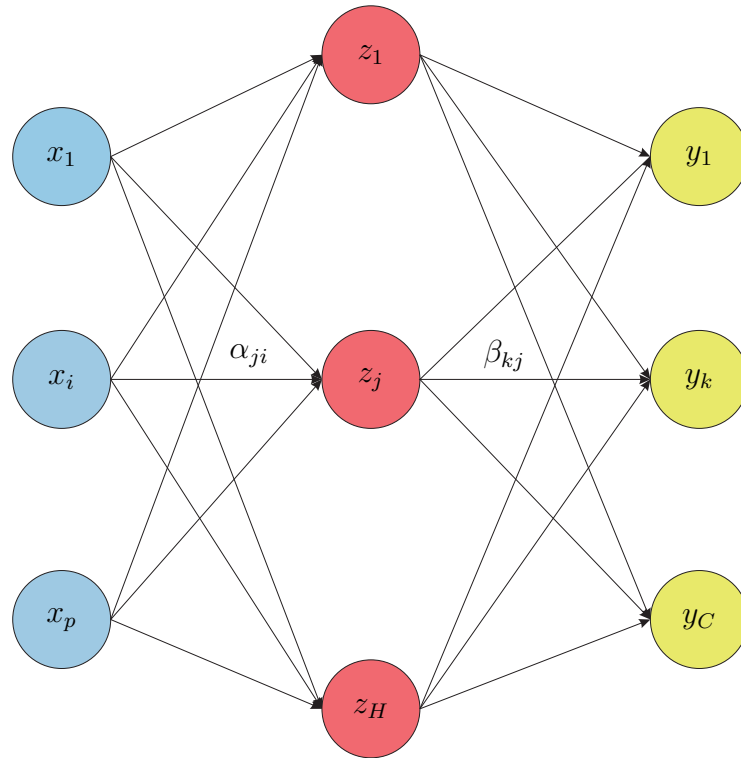


Figure 1: Neural Network with single hidden layer

We denote the complete set of weights by  $\boldsymbol{\theta}$ , which consists of

$$\boldsymbol{\theta} = \begin{cases} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{cases}$$

Each of the  $\alpha_j$  and  $\beta_k$  is a subset vector of  $\alpha$  and  $\beta$ ,

$$\alpha = (\alpha_1 \ \alpha_j \ \dots \ \alpha_H) = \begin{matrix} & j \\ & \alpha_1 \ \alpha_2 \ \dots \ \alpha_H \\ i \left( \begin{matrix} \alpha_{11} & \alpha_{21} & \dots & \alpha_{H1} \\ \alpha_{12} & \alpha_{22} & \dots & \alpha_{H2} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{1D} & \alpha_{2D} & \dots & \alpha_{HD} \end{matrix} \right)$$

$$\alpha^T = \begin{matrix} & i \\ j \left( \begin{matrix} \alpha_{11} & \alpha_{12} & \dots & \alpha_{1D} \\ \alpha_{21} & \alpha_{22} & \dots & \alpha_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{H1} & \alpha_{H2} & \dots & \alpha_{HD} \end{matrix} \right)$$

$$\beta = (\beta_1 \ \beta_k \ \dots \ \beta_C) = \begin{matrix} & k \\ & \beta_1 \ \beta_2 \ \dots \ \beta_C \\ j \left( \begin{matrix} \beta_{11} & \beta_{21} & \dots & \beta_{C1} \\ \beta_{12} & \beta_{22} & \dots & \beta_{C2} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{1H} & \beta_{2H} & \dots & \beta_{CH} \end{matrix} \right)$$

$$\beta^T = \begin{matrix} & j \\ k \left( \begin{matrix} \beta_{11} & \beta_{12} & \dots & \beta_{1H} \\ \beta_{21} & \beta_{22} & \dots & \beta_{2H} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{K1} & \beta_{K2} & \dots & \beta_{CH} \end{matrix} \right)$$

We can write the overall model as follows:

$$x_n \xrightarrow{\alpha} \mathbf{a}_n \xrightarrow{g} \mathbf{z}_n \xrightarrow{\beta} \mathbf{b}_n \xrightarrow{h} \hat{\mathbf{y}}_n$$

$x_n$  is a vector of length  $D$  containing  $D$  variables for the  $n$ th observation. That is,

$$x_n = \begin{pmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{pmatrix} \quad (11.1)$$

In general, we have the following models extract from the overall model,

$$\begin{aligned} a_{nj} &= \boldsymbol{\alpha}_j^T x_n, \quad j = 1, \dots, H \\ z_{nj} &= g(a_{nj}), \quad j = 1, \dots, H \\ b_{nk} &= \boldsymbol{\beta}_k^T \mathbf{z}_n, \quad k = 1, \dots, C \\ \hat{y}_{nk} &= h(b_{nk}), \quad k = 1, \dots, C \end{aligned}$$

From the above equations, it's not hard to see that  $\hat{y}_{nk}$  consisted from a bunch of other functions, this can be interpreted as,

$$\hat{\mathbf{y}}_n = h(\mathbf{b}_n) = h(\boldsymbol{\beta}^T \mathbf{z}_n) = h(\boldsymbol{\beta}^T g(\mathbf{a}_n)) = h(\boldsymbol{\beta}^T g(\boldsymbol{\alpha}^T x_n)) \quad (11.2)$$

Summarize the points above,  $x_n$  is the  $n$ th observation,  $\mathbf{a}_n = \boldsymbol{\alpha}^T x_n$  be the pre-synaptic hidden layer, and  $\mathbf{z}_n = g(\mathbf{a}_n)$  be the post-synaptic hidden layer, where  $g$  is some transfer function. We typically use  $g(a) = \text{sigm}(a)$ , but we may also use  $g(a) = \text{tanh}(a)$ . We now convert this hidden layer to the output layer. Let  $\mathbf{b}_n = \boldsymbol{\beta}^T \mathbf{z}_n$  be the pre-synaptic output layer, and  $\hat{\mathbf{y}}_n = h(\mathbf{b}_n)$  be the post-synaptic output layer, where  $h$  is another non-linearity. For a regression model, we use  $h(\mathbf{b}) = \mathbf{b}$ ; for binary classification, we use  $h(\mathbf{b}) = [\text{sigm}(b_1), \dots, \text{sigm}(b_c)]$ ; for multi-class classification, we use  $h(\mathbf{b}) = \mathcal{S}(\mathbf{b})$ .

## 11.2 Backpropagation neural networks

We now discuss how to compute the gradient vector of the NLL(negative log likelihood) by applying the chain rule of calculus. The resulting algorithm is known as backpropagation.

In the regression case, with  $C$  outputs, the NLL is given by the squared error:

$$J(\boldsymbol{\theta}) = - \sum_n \sum_k (\hat{y}_{nk}(\boldsymbol{\theta}) - y_{nk})^2 \quad (11.3)$$

In the classification case, with  $C$  classes, the NLL is given by the cross entropy,

$$J(\boldsymbol{\theta}) = - \sum_n \sum_k y_{nk} \log \hat{y}_{nk}(\boldsymbol{\theta}) \quad (11.4)$$

Our task is to compute  $\frac{d}{d\boldsymbol{\theta}} J$ . The overall gradient is obtained by summing over  $n$ , we will derive this for each  $n$  separately or we will perform mini-batch and summing the batch size. Let us start by considering the output layer weights. We have

$$\nabla_{\beta_k} J_n = \frac{\partial J_n}{\partial b_{nk}} \nabla_{\beta_k} b_{nk} = \frac{\partial J_n}{\partial b_{nk}} \mathbf{z}_n \quad (11.5)$$

Assuming  $h$  is canonical link function for the output GLM, then

$$\frac{\partial J_n}{\partial b_{nk}} = \delta_{nk}^\beta = (\hat{y}_{nk} - y_{nk}) \quad (11.6)$$

So the overall gradient is

$$\nabla_{\beta_k} J_n = \delta_{nk}^\beta \mathbf{z}_n \quad (11.7)$$

which is the pre-synaptic input to the output layer, namely  $\mathbf{z}_n$ , times the error signal, namely  $\delta_{nk}^\beta$ .

For the input layer weights, we have

$$\nabla_{\alpha_j} J_n = \frac{\partial J_n}{\partial a_{nj}} \nabla_{\alpha_j} a_{nj} = \delta_{nj}^\alpha x_n \quad (11.8)$$

The first level error signal  $\delta_{nj}^\alpha$ ,

$$\delta_{nj}^\alpha = \frac{\partial J_n}{\partial a_{nj}} = \sum_k \frac{\partial J_n}{\partial b_{nk}} \frac{\partial b_{nk}}{\partial a_{nj}} \quad (11.9)$$

where,

$$\frac{\partial b_{nk}}{\partial a_{nj}} = \frac{\partial b_{nk}}{\partial z_{nj}} \frac{\partial z_{nj}}{\partial a_{nj}} = \beta_{kj} g'(a_{nj}) \quad (11.10)$$

So, the first level error signal  $\delta_{nj}^\alpha$  is,

$$\delta_{nj}^\alpha = \frac{\partial J_n}{\partial a_{nj}} = \sum_k \delta_{nk}^\beta \beta_{kj} g'(a_{nj}) \quad (11.11)$$

In the end, the input layer weights looks like

$$\nabla_{\alpha_j} J_n = \delta_{nj}^\alpha x_n = \sum_k \delta_{nk}^\beta \beta_{kj} g'(a_{nj}) x_n \quad (11.12)$$

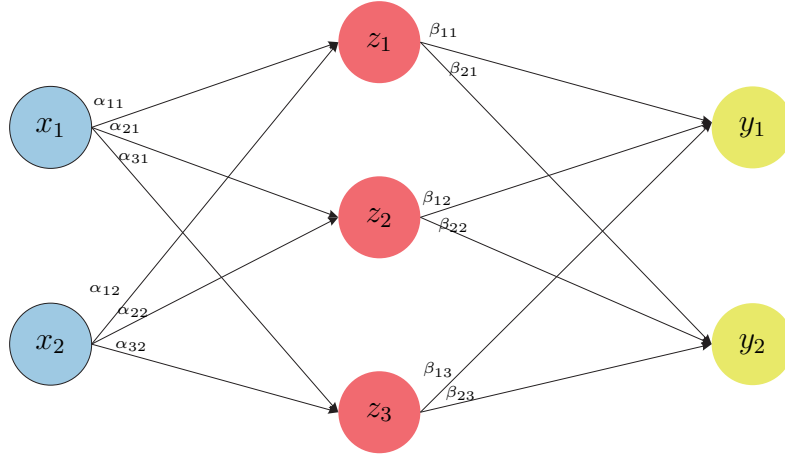
Given the derivatives, a gradient descent update at the  $(r + 1)$ st iteration has the form

$$\begin{aligned} \beta_k^{(r+1)} &= \beta_k^{(r)} - \eta \sum_n \nabla_{\beta_k} J_n \\ \alpha_j^{(r+1)} &= \alpha_j^{(r)} - \eta \sum_n \nabla_{\alpha_j} J_n \end{aligned}$$

Learning can be carried out online - processing each observation one at a time, updating the gradient after each training case, and cycling through the training cases many times. In this case, the sums of  $n$  in the update equations are replaced with single summation  $k$ . Online training allows the network to handle large training sets, and also update the weights as new observations come in.

The learning rate  $\eta$  for batch learning is usually taken to be a constant, and can be optimized by a **line search**.

### 11.2.1 Example 1



$$\boldsymbol{\alpha} = \begin{pmatrix} \alpha_1 & \alpha_2 & \alpha_3 \end{pmatrix} \quad (11.13)$$

$$\boldsymbol{\alpha}^T = \begin{pmatrix} \alpha_{11} & \alpha_{12} \\ \alpha_{21} & \alpha_{22} \\ \alpha_{31} & \alpha_{32} \end{pmatrix} \quad (11.14)$$

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_1 & \beta_2 \end{pmatrix} \quad (11.15)$$

$$\boldsymbol{\beta}^T = \begin{pmatrix} \beta_{11} & \beta_{12} & \beta_{13} \\ \beta_{21} & \beta_{22} & \beta_{23} \end{pmatrix} \quad (11.16)$$

The first output layer,

$$\begin{aligned} a_1 &= \boldsymbol{\alpha}_1^T x = \sum_{i=1}^2 \alpha_{1i} x_i \\ a_2 &= \boldsymbol{\alpha}_2^T x = \sum_{i=1}^2 \alpha_{2i} x_i \\ a_3 &= \boldsymbol{\alpha}_3^T x = \sum_{i=1}^2 \alpha_{3i} x_i \end{aligned}$$

So,

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \quad (11.17)$$

The post-synaptic hidden layer

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix} = \begin{pmatrix} g(a_1) \\ g(a_2) \\ g(a_3) \end{pmatrix} \quad (11.18)$$

The second output layer,

$$\begin{aligned} b_1 &= \boldsymbol{\beta}_1^T \mathbf{z} = \sum_{j=1}^3 \beta_{1j} z_j \\ b_2 &= \boldsymbol{\beta}_2^T \mathbf{z} = \sum_{j=1}^3 \beta_{2j} z_j \end{aligned}$$

The estimated response,

$$\hat{\mathbf{y}} = \begin{pmatrix} h(b_1) \\ h(b_2) \end{pmatrix} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \end{pmatrix}$$

a gradient descent update at the  $(r + 1)$ st iteration has the form

$$\boldsymbol{\beta}_1^{(r+1)} = \boldsymbol{\beta}_1^{(r)} - \eta \sum_n \nabla_{\boldsymbol{\beta}_1} J_n = \begin{pmatrix} \beta_{11} \\ \beta_{12} \\ \beta_{13} \end{pmatrix} - \eta \left( \sum_n (\hat{y}_{n1} - y_{n1}) \begin{pmatrix} z_{n1} \\ z_{n2} \\ z_{n3} \end{pmatrix} \right)$$

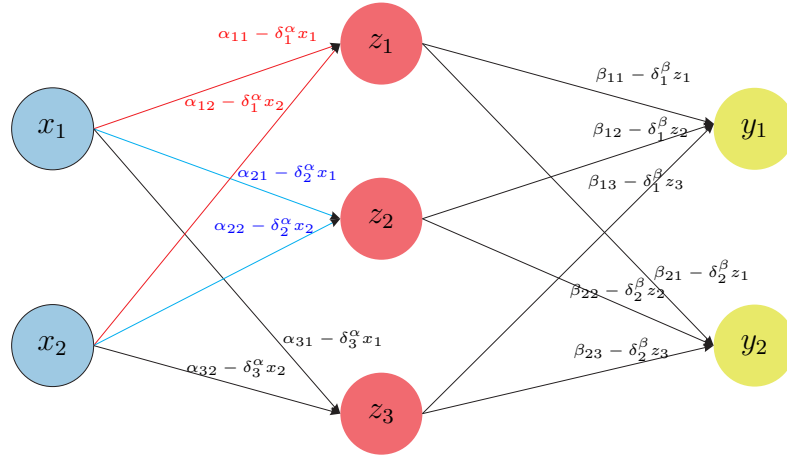
$$\beta_2^{(r+1)} = \beta_2^{(r)} - \eta \sum_n \nabla_{\beta_2} J_n = \begin{pmatrix} \beta_{21} \\ \beta_{22} \\ \beta_{23} \end{pmatrix} - \eta \left( \sum_n (\hat{y}_{n2} - y_{n2}) \begin{pmatrix} z_{n1} \\ z_{n2} \\ z_{n3} \end{pmatrix} \right)$$

$$\alpha_1^{(r+1)} = \alpha_1^{(r)} - \eta \sum_n \nabla_{\alpha_1} J_n = \begin{pmatrix} \alpha_{11} \\ \alpha_{12} \end{pmatrix} - \eta \left( \sum_n \sum_k (\delta_{nk}^\beta \beta_{k1}) g'(a_1) \begin{pmatrix} x_{n1} \\ x_{n2} \end{pmatrix} \right)$$

$$\alpha_2^{(r+1)} = \alpha_2^{(r)} - \eta \sum_n \nabla_{\alpha_2} J_n = \begin{pmatrix} \alpha_{21} \\ \alpha_{22} \end{pmatrix} - \eta \left( \sum_n \sum_k (\delta_{nk}^\beta \beta_{k2}) g'(a_2) \begin{pmatrix} x_{n1} \\ x_{n2} \end{pmatrix} \right)$$

$$\alpha_3^{(r+1)} = \alpha_3^{(r)} - \eta \sum_n \nabla_{\alpha_3} J_n = \begin{pmatrix} \alpha_{31} \\ \alpha_{32} \end{pmatrix} - \eta \left( \sum_n \sum_k (\delta_{nk}^\beta \beta_{k3}) g'(a_3) \begin{pmatrix} x_{n1} \\ x_{n2} \end{pmatrix} \right)$$

The graph of back propagation



I didn't add learning rate on the update, since the tag is too small to add any words on it.



## 11.3 General form of Feed-forward and Back-propagation

### 11.3.1 Notation

- $\vec{X}_\ell$  : denotes the input data matrix in layer  $\ell$  with  $N$  input size and  $D$  input dimensional size.
- $\vec{W}_\ell$  : denotes the weight matrix in layer  $\ell$  with  $D$  input dimensional size and  $H$  output layer size.
- $\vec{b}_\ell$  : denotes the bias vector in layer  $\ell$  with  $H$  output layer size.
- $\vec{a}_\ell$  : denotes the output matrix before activation from layer  $\ell$ .
- $f_\ell(\cdot)$  : denotes the activation function in layer  $\ell$ .

$$\vec{X}_\ell = \begin{matrix} & D \\ N & \left( \begin{matrix} n(x_\ell)_d \end{matrix} \right) \end{matrix}$$

$$\vec{W}_\ell = \begin{matrix} & H \\ D & \left( \begin{matrix} (w_\ell)_h^d \end{matrix} \right) \end{matrix}$$

$$\vec{b}_\ell = \begin{matrix} & H \\ & \left( \begin{matrix} (b_\ell)_h \end{matrix} \right) \end{matrix}$$

$$\vec{a}_\ell = N \begin{pmatrix} H \\ {}^n(a_\ell)_h = {}^n(x_\ell)_d \cdot (w_{\ell-1})_h^d + (b_{\ell-1})_h \end{pmatrix}$$

$$\vec{X}_{\ell+1} = f_\ell(\vec{a}_\ell) \quad (11.19)$$

### 11.3.1.1 Feed-Forward Neural Network

For a two layer fully-connected neural network, the network has the following architecture:

$$\vec{X}_{\ell-1} \mapsto \vec{a}_{\ell-1} = \vec{X}_{\ell-1} \vec{W}_{\ell-2} \rightarrow \vec{X}_\ell = f_{\ell-1}(\vec{a}_{\ell-1}) \mapsto \vec{a}_\ell = \vec{X}_\ell \vec{W}_{\ell-1} \rightarrow \hat{y} = \text{softmax}(\vec{a}_\ell) \quad (11.20)$$

We denote the loss function as,

$$L = \text{loss}(y, \hat{y}) = \sum_n \sum_h \text{loss}({}^n y_h, {}^n \hat{y}_h) \quad (11.21)$$

### 11.3.1.2 Back-propagation Neural Network

Let us start by considering the last layer weights  $(w_{\ell-1})_h^d$  and perform the derivative on the loss function

$$\frac{\partial}{\partial (w_{\ell-1})_h^d} L = \frac{\partial L}{\partial {}^n(a_\ell)_h} \frac{\partial {}^n(a_\ell)_h}{\partial (w_{\ell-1})_h^d} = \frac{\partial L}{\partial {}^n(a_\ell)_h} {}^n(x_\ell)_d = {}^n(\delta_\ell)_h \cdot {}^n(x_\ell)_d = {}_n(x_\ell^T)^d \cdot {}^n(\delta_\ell)_h \quad (11.22)$$

For the ease of notation, we denote  ${}^n(\delta_\ell)_h$  as the error signal in layer  $\ell$ . Now, we derivative of  $(w_{\ell-2})_h^d$  on the loss function,

$$\frac{\partial}{\partial (w_{\ell-2})_h^d} L = \frac{\partial L}{\partial {}^n(a_{\ell-1})_h} \frac{\partial {}^n(a_{\ell-1})_h}{\partial (w_{\ell-2})_h^d} \quad (11.23)$$

$$= \frac{\partial L}{\partial {}^n(a_{\ell-1})_h} {}^n(x_{\ell-1})_d \quad (11.24)$$

$$= {}_n(x_{\ell-1}^T)^d \cdot {}^n(\delta_{\ell-1})_h \quad (11.25)$$

For a two layer( $\ell = 2$ ) fully connected neural network, the error signal for the last layer has the below form,

$$\begin{aligned}
{}^n(\delta_\ell)_h &= \frac{\partial L}{\partial {}^n(a_\ell)_h} \\
&= \text{loss}'(y, \hat{y}) \hat{y}' \\
&= \text{loss}'(y, \hat{y}) \odot f'_\ell({}^n(a_\ell)_h)
\end{aligned} \tag{11.26}$$

For the error signal in the first layer,

$$\begin{aligned}
{}^n(\delta_{\ell-1})_h &= \frac{\partial L}{\partial {}^n(a_{\ell-1})_h} \\
&= \sum_d \frac{\partial L}{\partial {}^n(a_\ell)_d} \cdot \frac{\partial {}^n(a_\ell)_d}{\partial {}^n(a_{\ell-1})_h} \\
&= \sum_d {}^n(\delta_\ell)_d \cdot \frac{\partial {}^n(a_\ell)_d}{\partial {}^n(a_{\ell-1})_h}
\end{aligned}$$

We'll show how to proof  $\frac{\partial {}^n(a_\ell)_d}{\partial {}^n(a_{\ell-1})_h}$ . For  ${}^n(a_\ell)_d$ , we know that

$$\begin{aligned}
{}^n(a_\ell)_d &= {}^n(x_\ell)_h \cdot (w_{\ell-1})_d^h + (b_{\ell-1})_d \\
&= f_{\ell-1}({}^n(a_{\ell-1})_h) (w_{\ell-1})_d^h + (b_{\ell-1})_d
\end{aligned}$$

So,

$$\frac{\partial {}^n(a_\ell)_d}{\partial {}^n(a_{\ell-1})_h} = f'_{\ell-1}({}^n(a_{\ell-1})_h) (w_{\ell-1})_d^h \tag{11.27}$$

Finally, the complete form of the error signal in the first layer,

$$\begin{aligned}
{}^n(\delta_{\ell-1})_h &= f'_{\ell-1}({}^n(a_{\ell-1})_h) \sum_d {}^n(\delta_\ell)_d \cdot (w_\ell)_d^h \\
&= f'_{\ell-1}({}^n(a_{\ell-1})_h) \odot {}^n(\delta_\ell)_d \cdot (w_\ell^T)_h^d
\end{aligned}$$

The gradient of bias is similar with the above proof,

$$\begin{aligned}
\frac{\partial}{\partial (b_\ell)_h} L &= \frac{\partial L}{\partial {}^n(a_\ell)_h} \frac{\partial {}^n(a_\ell)_h}{\partial (b_\ell)_h} = \sum_n {}^n(\delta_\ell)_h \\
\frac{\partial}{\partial (b_{\ell-1})_h} L &= \frac{\partial L}{\partial {}^n(a_{\ell-1})_h} \frac{\partial {}^n(a_{\ell-1})_h}{\partial (b_{\ell-1})_h} = \sum_n {}^n(\delta_{\ell-1})_h
\end{aligned}$$

The loss function, full gradients and L2 regularization,

$$\begin{aligned}
L &= \text{loss}(y, \hat{y}) + \frac{\lambda}{2}(\vec{W}_\ell^2 + \vec{W}_{\ell-1}^2) \\
\frac{\partial}{\partial(w_\ell)_h} L &= {}_n(x_\ell^T)^d \cdot {}^n(\delta_\ell)_h + \lambda(w_\ell)_h^d \\
&= {}_n(x_\ell^T)^d (\text{loss}'(y, \hat{y}) \odot f'_\ell({}^n(a_\ell)_h)) + \lambda(w_\ell)_h^d \\
\frac{\partial}{\partial(w_{\ell-1})_h} L &= {}_n(x_{\ell-1}^T)^d {}^n(\delta_\ell)_h + \lambda(w_{\ell-1})_h^d \\
&= {}_n(x_{\ell-1}^T)^d \cdot (f'_{\ell-1}({}^n(a_{\ell-1})_h) \odot ({}^n(\delta_\ell)_d \cdot (w_\ell^T)_h^d)) + \lambda(w_{\ell-1})_h^d \\
\frac{\partial}{\partial(b_\ell)_h} L &= \frac{\partial L}{\partial {}^n(a_\ell)_h} \frac{\partial {}^n(a_\ell)_h}{\partial (b_\ell)_h} = \sum_n {}^n(\delta_\ell)_h \\
\frac{\partial}{\partial(b_{\ell-1})_h} L &= \frac{\partial L}{\partial {}^n(a_{\ell-1})_h} \frac{\partial {}^n(a_{\ell-1})_h}{\partial (b_{\ell-1})_h} = \sum_n {}^n(\delta_{\ell-1})_h
\end{aligned}$$