# QBUS6850

# Lecture 4

## Neural Network and Deep Learning- II

© *Discipline of Business Analytics*

**BUSINESS SCHOOL**

*QBUS6850 Team*

THE UNIVERSITY OF
SYDNEY

## ❑ Topics covered

- Neural Network regression and classification loss functions

- Backward propagation Neural Network

## ❑ References

- Alpaydin (2014), Chapter 11
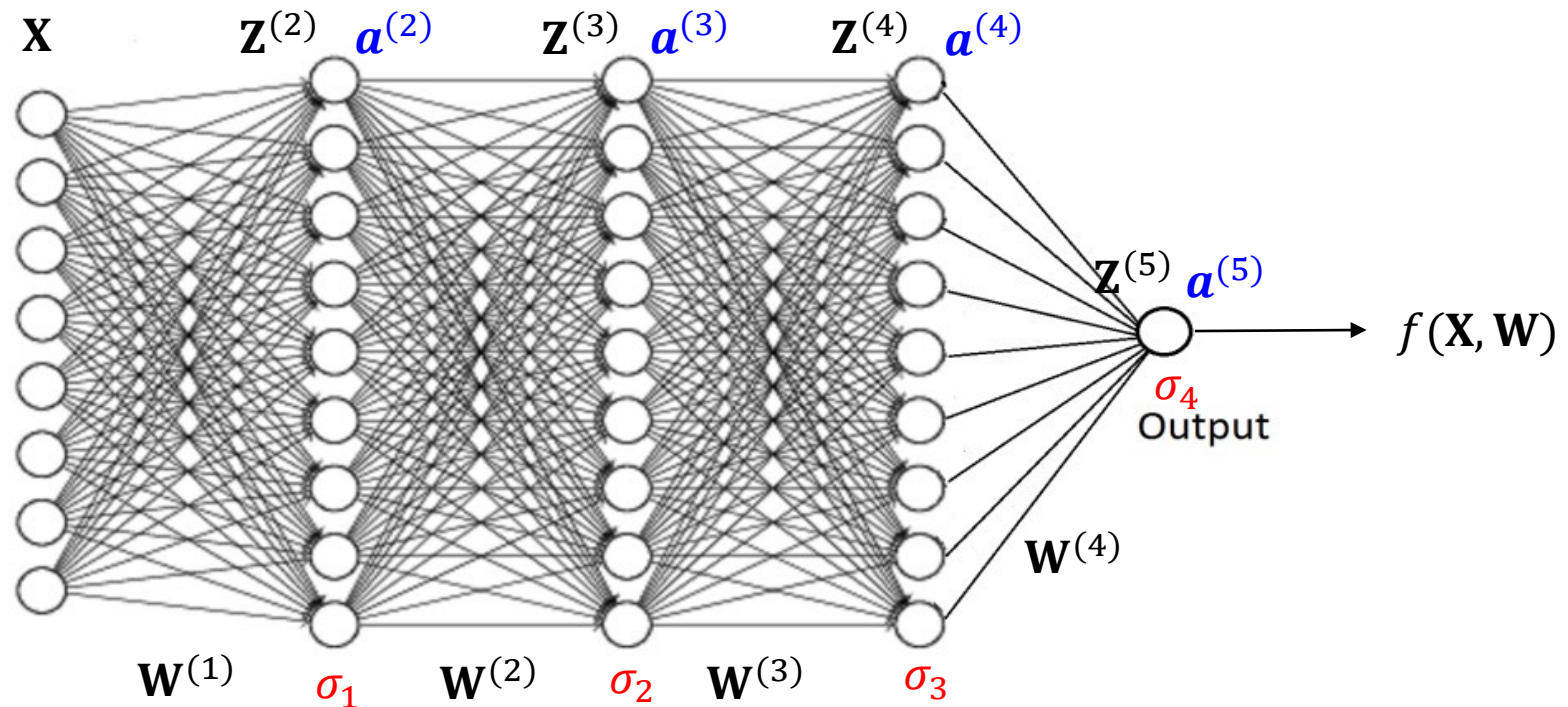
- Bishop (2006), Chapter 5

# Learning Objectives

➢ Understand the Neural Network regression and classification loss functions

➢ Understand why and how to add regularization terms into Neural Network

➢ Understand the intuition of backpropagation

➢ Understand the process of backpropagation

# Recap Week 3

$$\mathbf{Z}^{(2)} = \mathbf{X}\mathbf{W}^{(1)}; \boldsymbol{a}^{(2)} = \sigma_1\big(\mathbf{Z}^{(2)}\big); \mathbf{Z}^{(3)} = \boldsymbol{a}^{(2)}\mathbf{W}^{(2)}; \boldsymbol{a}^{(3)} = \sigma_2\big(\mathbf{Z}^{(3)}\big);$$

$$\mathbf{Z}^{(4)} = \boldsymbol{a}^{(3)}\mathbf{W}^{(3)}; \boldsymbol{a}^{(4)} = \sigma_3\big(\mathbf{Z}^{(4)}\big); \mathbf{Z}^{(5)} = \boldsymbol{a}^{(4)}\mathbf{W}^{(4)}; \boldsymbol{a}^{(5)} = \sigma_4\big(\mathbf{Z}^{(5)}\big); f(\mathbf{X}, \mathbf{W}) = \boldsymbol{a}^{(5)}$$

$$f(\mathbf{X}, \mathbf{W}) = \sigma_4(\sigma_3(\sigma_2(\sigma_1(\mathbf{X}\mathbf{W}^{(1)})\mathbf{W}^{(2)})\mathbf{W}^{(3)})\mathbf{W}^{(4)})$$     A composition function

Another way to look at biases where $\mathbf{b}^{(l)}$ in row shape

$$\mathbf{Z}^{(2)} = \mathbf{X}\mathbf{W}^{(1)} + \mathbf{1}_N \mathbf{b}^{(1)}; \quad \boldsymbol{a}^{(2)} = \sigma_1\big(\mathbf{Z}^{(2)}\big);$$

$$\mathbf{Z}^{(3)} = \boldsymbol{a}^{(2)}\mathbf{W}^{(2)} + \mathbf{1}_N \mathbf{b}^{(2)}; \quad \boldsymbol{a}^{(3)} = \sigma_2\big(\mathbf{Z}^{(3)}\big);$$

$$\mathbf{Z}^{(4)} = \boldsymbol{a}^{(3)}\mathbf{W}^{(3)} + \mathbf{1}_N \mathbf{b}^{(3)}; \boldsymbol{a}^{(4)} = \sigma_3\big(\mathbf{Z}^{(4)}\big);$$

$$\mathbf{Z}^{(5)} = \boldsymbol{a}^{(4)}\mathbf{W}^{(4)} + \mathbf{1}_N \mathbf{b}^{(4)}; \quad \boldsymbol{a}^{(5)} = \sigma_4\big(\mathbf{Z}^{(5)}\big);$$

$$f(\mathbf{X}, \mathbf{W}) = \boldsymbol{a}^{(5)}$$

$$f(\mathbf{X}, \mathbf{W}) = \sigma_4\big(\sigma_3\big(\sigma_2\big(\sigma_1\big(\mathbf{X}\mathbf{W}^{(1)} + \mathbf{1}_N \mathbf{b}^{(1)}\big)\mathbf{W}^{(2)} + \mathbf{1}_N \mathbf{b}^{(2)}\big)\mathbf{W}^{(3)} + \mathbf{1}_N \mathbf{b}^{(3)}\big)\mathbf{W}^{(4)} + \mathbf{1}_N \mathbf{b}^{(4)}\big)$$

# Neural Network Loss Function

$$L(\mathbf{W}) = \frac{1}{2N} \sum_{n=1}^{N} (f(\mathbf{x}_n, \mathbf{W}) - t_n)^2$$

How NN and linear regression loss functions are different?

Normally for regression, the output neuron takes the identity activation

As in the regression we did in week 2&3, we can use gradient descent to try to minimize the loss function as a function of parameters $\mathbf{W}$. $\alpha$ is the learning rate.

$$\mathbf{W} := \mathbf{W} - \alpha \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}}$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} = ?$$

This will be a much harder problem, since now we have weights $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}$ for a 3-layer NN.

$$L(\mathbf{W}) = -\left[\sum_{n=1}^{N}\left(t_n \log(f(\mathbf{x}_n, \mathbf{W})) + (1 - t_n)\log(1 - f(\mathbf{x}_n, \mathbf{W}))\right)\right]$$

where each $t_n$ is 0 (if $\mathbf{x}_n$ is in class 0) or 1 (if $\mathbf{x}_n$ is in class 1)

In classification, the output neuron takes the sigmoid activation

Similarly we can use gradient descent to try to minimize the loss function as a function of parameters $\mathbf{W}$. $\alpha$ is the learning rate.

$$\mathbf{W} := \mathbf{W} - \alpha \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}}$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} = ?$$

This will be a much harder problem, since now we have weights $\mathbf{W}^{(1)}$, $\mathbf{W}^{(2)}$ for a 3-layer NN.
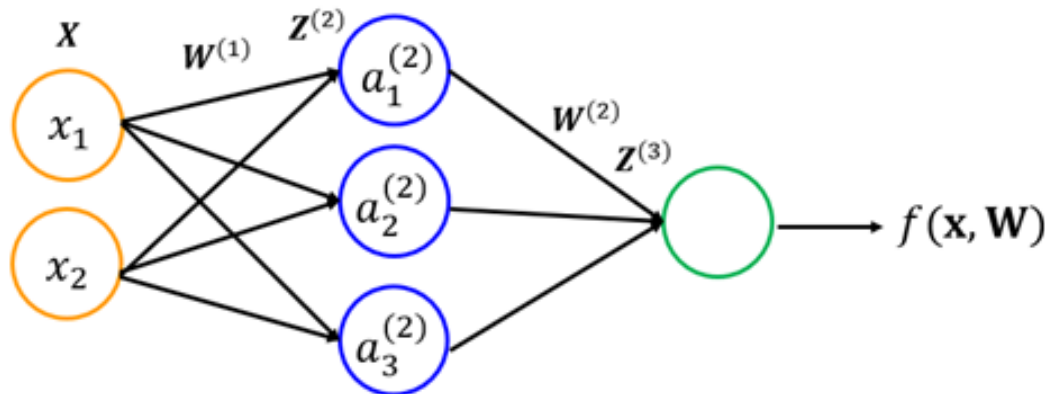
# Neural Network Training

# **Parameters to be estimated**

Target: minimize the loss function by changing the weights $W^{(1)}$, $W^{(2)}$

Let's still use 3-layer NN without bias units example that we had last week



$$\mathbf{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \end{bmatrix}$$

$$\mathbf{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} \\ w_{21}^{(2)} \\ w_{31}^{(2)} \end{bmatrix}$$

Poll 1

# Dimensionality

Why not simply try all the potential weights and see which combination produces the smallest loss?

Curse of dimensionality:

- 1 parameter: 1000 trials, 0.05 seconds
- 2 parameters: 1000 ×1000=1 million trials, 0.05*1000= 50 seconds
- 3 parameters: 1000 × 1000 × 1000=1 billion trials; (1,000,000*0.05)/(60*60)= 13.89 hours
- …
- 9 parameters: how long it will take?

Besides estimating the parameters of NN, other problems in neural network modelling:

- How to select the number of hidden layers?
- How to select the number of units in each hidden layer?
- How to perform feature selection?
- ...

# Backpropagation Intuition

# Backpropagation

- Backpropagation is a method used in artificial neural networks to calculate the error contribution of each neuron after a batch of data (in image recognition, multiple images) is processed. This is used by an enveloping optimization algorithm to adjust the weight of each neuron, completing the learning process for that case.

- Technically it calculates the gradient of the loss function. It is commonly used in the gradient descent optimization algorithm. It is also called **backward propagation of errors**, because the **error is calculated at the output and distributed back** through the network layers.

Poll 2

https://en.wikipedia.org/wiki/Backpropagation

- Phase 1: propagation. Each propagation involves the following steps:
  - Propagation forward through the network to generate the output value(s)
  - Calculation of the loss (error term)
  - Propagation of the output activations back through the network using the training pattern target in order to generate the **deltas ($\delta$)** of all output and hidden neurons.

- Phase 2: weight (parameter) update. For each weight, the following steps must be followed:
  - The weight's output **deltas** and input activation are multiplied to find the gradient of the weight.
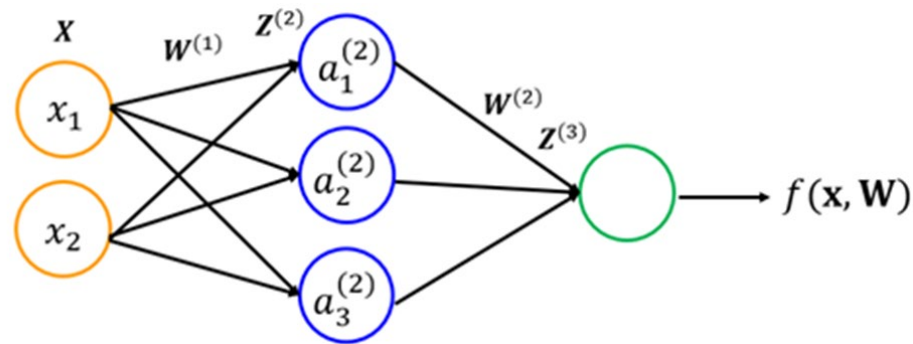  - A ratio (learning rate) of the weight's gradient is subtracted from the weight.

# Backpropagation

- Now we are going to use back propagation to train this NN for regression
- The example we had last week
- No bias unit for simplicity

$$L(\mathbf{W}) = \frac{1}{2} \sum_{n=1}^{N} (f(\mathbf{x}_n, \mathbf{W}) - t_n)^2$$

| $x_1$ | $x_2$ |
|-------|-------|
| 1 | 0.875 |
| 0.25 | 1 |
| 0.50 | 0.45 |
| 0.75 | 0.25 |



| expense |
|---------|
| 1 |
| 0.4 |
| 0.5 |
| 0.6 |

Consider 4 data

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} = ?$$

Our parameter set $\mathbf{W}$ contains $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$

Hidden layer size

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \end{bmatrix}$$

Input layer size.
2 features.

Output layer size

$$\mathbf{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} \\ w_{21}^{(2)} \\ w_{31}^{(2)} \end{bmatrix}$$

Hidden layer size

Hidden layer size

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \begin{bmatrix} \dfrac{\partial L(\mathbf{W})}{\partial w_{11}^{(1)}} & \dfrac{\partial L(\mathbf{W})}{\partial w_{12}^{(1)}} & \dfrac{\partial L(\mathbf{W})}{\partial w_{13}^{(1)}} \\ \dfrac{\partial L(\mathbf{W})}{\partial w_{21}^{(1)}} & \dfrac{\partial L(\mathbf{W})}{\partial w_{22}^{(1)}} & \dfrac{\partial L(\mathbf{W})}{\partial w_{23}^{(1)}} \end{bmatrix}$$

Input layer size

Output layer size

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \begin{bmatrix} \dfrac{\partial L(\mathbf{W})}{\partial w_{11}^{(2)}} \\ \dfrac{\partial L(\mathbf{W})}{\partial w_{21}^{(2)}} \\ \dfrac{\partial L(\mathbf{W})}{\partial w_{31}^{(2)}} \end{bmatrix}$$

Hidden layer size

- First, work on the partial derivatives of loss with respect to $\boldsymbol{W}^{(2)}$.
- Note that now we start from the **"right"** of NN, or an opposite direction compared to forward propagation

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \frac{\partial \frac{1}{2} \sum_{n=1}^{N} (f(\mathbf{x}_n, \mathbf{W}) - t_n)^2}{\partial \mathbf{W}^{(2)}}$$

Sum rules in differentiation: derivative of sum equals to the sum of derivatives

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \sum_{n=1}^{N} \frac{\partial \frac{1}{2} (f(\mathbf{x}_n, \mathbf{W}) - t_n)^2}{\partial \mathbf{W}^{(2)}}$$

Remove summation for simplicity, will look after it later.
~~Also remove the subscript $(n)$ in $f(\mathbf{x}_n, \mathbf{W})$ and $t_n$ for simplicity.~~

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \frac{\partial \frac{1}{2} (f(\mathbf{x}_n, \mathbf{W}) - t_n)^2}{\partial \mathbf{W}^{(2)}}$$

**Derivative chain rule**

$t_n$ is a constant
w.r.t. $\boldsymbol{W}^{(2)}$

$f(\mathbf{x}_n, \mathbf{W}) = \boldsymbol{\sigma}(\boldsymbol{z}_n^{(3)})$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \frac{\partial \frac{1}{2}(f(\mathbf{x}_n, \mathbf{W}) - t_n)^2}{\partial \mathbf{W}^{(2)}} = (f(\mathbf{x}_n, \mathbf{W}) - t_n)\frac{\partial f(\mathbf{x}_n, \mathbf{W})}{\partial \mathbf{W}^{(2)}}$$

$$= (f(\mathbf{x}_n, \mathbf{W}) - t_n)\boxed{\frac{\partial f(\mathbf{x}_n, \mathbf{W})}{\partial \mathbf{z}_n^{(3)}} * \frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(2)}}}$$

**Derivative chain rule.
Based on FP, $z_n^{(3)}$ is
between $f(\mathbf{x}_n, \mathbf{W})$ and $W^{(2)}$**

$$= (f(\mathbf{x}_n, \mathbf{W}) - t_n)\boxed{\frac{\partial \boldsymbol{\sigma}(\mathbf{z}_n^{(3)})}{\partial \mathbf{z}_n^{(3)}}}\frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(2)}}$$

**How to
calculate this?**

$$\frac{\partial \boldsymbol{\sigma}(\mathbf{z}_n^{(3)})}{\partial \mathbf{z}_n^{(3)}} = \sigma'\left(\mathbf{z}_n^{(3)}\right) = \sigma\left(\mathbf{z}_n^{(3)}\right).*(1 - \sigma\left(\mathbf{z}_n^{(3)}\right))$$
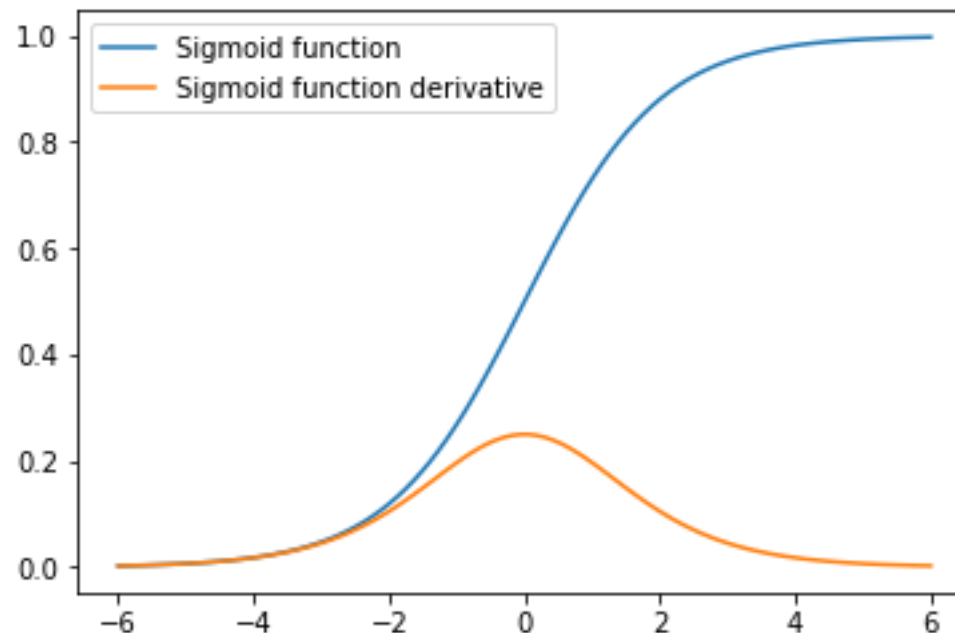
Poll 3 & 4

For single data

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial \boldsymbol{\sigma}(\mathbf{z}_n^{(3)})}{\partial \mathbf{z}_n^{(3)}} = ?$$

$$\sigma'(z) = \frac{d}{dz}\left(\frac{1}{1+e^{-z}}\right) = \frac{e^{-z}}{(1+e^{-z})^2} = \sigma(z)(1 - \sigma(z))$$

**For a single data $\mathbf{x}_n$, further we have**

$$\frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(2)}} = \frac{\partial}{\partial \mathbf{W}^{(2)}} \left( w_{11}^{(2)} a_{n1}^{(2)} + w_{21}^{(2)} a_{n2}^{(2)} + w_{31}^{(2)} a_{n3}^{(2)} \right) = [a_{n1}^{(2)}, \ a_{n2}^{(2)}, a_{n3}^{(2)}]^T$$

**Hence the derivative for a single data is**

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = (f(\mathbf{x}_n, \mathbf{W}) - t_n) \, \sigma' \left( \mathbf{z}_n^{(3)} \right) [a_{n1}^{(2)}, \ a_{n2}^{(2)}, a_{n3}^{(2)}]^T$$

**We add them together**

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \sum_{n=1}^{N} (f(\mathbf{x}_n, \mathbf{W}) - t_n) \, \sigma' \left( \mathbf{z}_n^{(3)} \right) \begin{bmatrix} a_{n1}^{(2)} \\ a_{n2}^{(2)} \\ a_{n3}^{(2)} \end{bmatrix}$$

**Next, we show how this can be done in matrix operation**

Now we write it for all the data $\mathbf{X}$ and $\mathbf{t}$ in matrix form

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \boxed{\left(f(\mathbf{X}, \mathbf{W}) - \mathbf{t}\right)\sigma'\left(\mathbf{Z}^{(3)}\right)} \circledast \boxed{\frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{W}^{(2)}}}$$

**Known**

**How to organize this?**

here $N = 4$:

$$\mathbf{Z}^{(3)} = \begin{bmatrix} z_{11}^{(3)} \\ z_{21}^{(3)} \\ z_{31}^{(3)} \\ z_{41}^{(3)} \end{bmatrix} \qquad \boldsymbol{a}^{(2)} = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} & a_{13}^{(2)} \\ a_{21}^{(2)} & a_{22}^{(2)} & a_{23}^{(2)} \\ a_{31}^{(2)} & a_{32}^{(2)} & a_{33}^{(2)} \\ a_{41}^{(2)} & a_{42}^{(2)} & a_{43}^{(2)} \end{bmatrix} \qquad \mathbf{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} \\ w_{21}^{(2)} \\ w_{31}^{(2)} \end{bmatrix}$$

$$\frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(2)}} = [a_{n1}^{(2)}, \ a_{n2}^{(2)}, \ a_{n3}^{(2)}]^T$$

$$\mathbf{Z}^{(3)} = \boxed{a^{(2)}}\mathbf{W}^{(2)}$$

$4 \times 1 \qquad\qquad 3 \times 1$

$4 \times 3$

$\mathbf{Z}^{(3)}$ is actually a linear combination of $\boldsymbol{W}^{(2)}$ with $\boldsymbol{a}^{(2)}$, so:

$$\frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{W}^{(2)}} = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} & a_{13}^{(2)} \\ a_{21}^{(2)} & a_{22}^{(2)} & a_{23}^{(2)} \\ a_{31}^{(2)} & a_{32}^{(2)} & a_{33}^{(2)} \\ a_{41}^{(2)} & a_{42}^{(2)} & a_{43}^{(2)} \end{bmatrix}$$

# Backpropagation- Step 1

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \overbrace{\left(f(\mathbf{X}, \mathbf{W}) - \mathbf{t}\right) . * \sigma'\left(\mathbf{Z}^{(3)}\right)}^{4 \times 1 \qquad 4 \times 1} \quad \circledast \quad \overbrace{\frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{W}^{(2)}}}^{4 \times 3}$$

Here we use matrix element wise product

$$\begin{bmatrix} f(\mathbf{x}_1, \mathbf{W}) - t_1 \\ f(\mathbf{x}_2, \mathbf{W}) - t_2 \\ f(\mathbf{x}_3, \mathbf{W}) - t_3 \\ f(\mathbf{x}_4, \mathbf{W}) - t_4 \end{bmatrix} .* \begin{bmatrix} \sigma'\left(z_{11}^{(3)}\right) \\ \sigma'\left(z_{21}^{(3)}\right) \\ \sigma'\left(z_{31}^{(3)}\right) \\ \sigma'\left(z_{41}^{(3)}\right) \end{bmatrix}$$

**Backpropagation error**

$$\begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \\ \delta_3^{(3)} \\ \delta_4^{(3)} \end{bmatrix}$$

$$a^{(2)} = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} & a_{13}^{(2)} \\ a_{21}^{(2)} & a_{22}^{(2)} & a_{23}^{(2)} \\ a_{31}^{(2)} & a_{32}^{(2)} & a_{33}^{(2)} \\ a_{41}^{(2)} & a_{42}^{(2)} & a_{43}^{(2)} \end{bmatrix}$$

Now we need to calculate $\delta^{(3)}$ times $a^{(2)}$

# Backpropagation- Step 1

$$4 \times 1 \qquad 4 \times 1 \qquad 4 \times 3$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \left( f(\mathbf{X}, \mathbf{W}) - \mathbf{t} \right) .* \; \sigma'\left(\mathbf{Z}^{(3)}\right) \;\; \circledast \;\; \frac{\partial \mathbf{Z}^{(3)}}{\partial \mathbf{W}^{(2)}}$$

Here we use matrix element wise product

$$\begin{bmatrix} f(\mathbf{x}_1, \mathbf{W}) - t_1 \\ f(\mathbf{x}_2, \mathbf{W}) - t_2 \\ f(\mathbf{x}_3, \mathbf{W}) - t_3 \\ f(\mathbf{x}_4, \mathbf{W}) - t_4 \end{bmatrix} .* \begin{bmatrix} \sigma'\left(z_{11}^{(3)}\right) \\ \sigma'\left(z_{21}^{(3)}\right) \\ \sigma'\left(z_{31}^{(3)}\right) \\ \sigma'\left(z_{41}^{(3)}\right) \end{bmatrix}$$

How to connect them?

**Backpropagation error**

$$\begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \\ \delta_3^{(3)} \\ \delta_4^{(3)} \end{bmatrix}$$

$$a^{(2)} = \begin{bmatrix} a_{11}^{(2)} & a_{12}^{(2)} & a_{13}^{(2)} \\ a_{21}^{(2)} & a_{22}^{(2)} & a_{23}^{(2)} \\ a_{31}^{(2)} & a_{32}^{(2)} & a_{33}^{(2)} \\ a_{41}^{(2)} & a_{42}^{(2)} & a_{43}^{(2)} \end{bmatrix}$$

Taking transpose of $\boldsymbol{a}^{(2)}$, then times $\left(\boldsymbol{a}^{(2)}\right)^T$ by $\delta^{(3)}$

$$\left(\boldsymbol{a}^{(2)}\right)^T = \begin{bmatrix} a_{11}^{(2)} & a_{21}^{(2)} & a_{31}^{(2)} & a_{41}^{(2)} \\ a_{12}^{(2)} & a_{22}^{(2)} & a_{32}^{(2)} & a_{42}^{(2)} \\ a_{13}^{(2)} & a_{23}^{(2)} & a_{33}^{(2)} & a_{43}^{(2)} \end{bmatrix}$$

$$\delta^{(3)} = \begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \\ \delta_3^{(3)} \\ \delta_4^{(3)} \end{bmatrix}$$

$3 \times 4$

$4 \times 1$

**Step 1: summary**

$4 \times 1$

1st equation uses matrix product

2nd equation uses matrix **element wise** product

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \boxed{\left(\boldsymbol{a}^{(2)}\right)^T \delta^{(3)}}$$
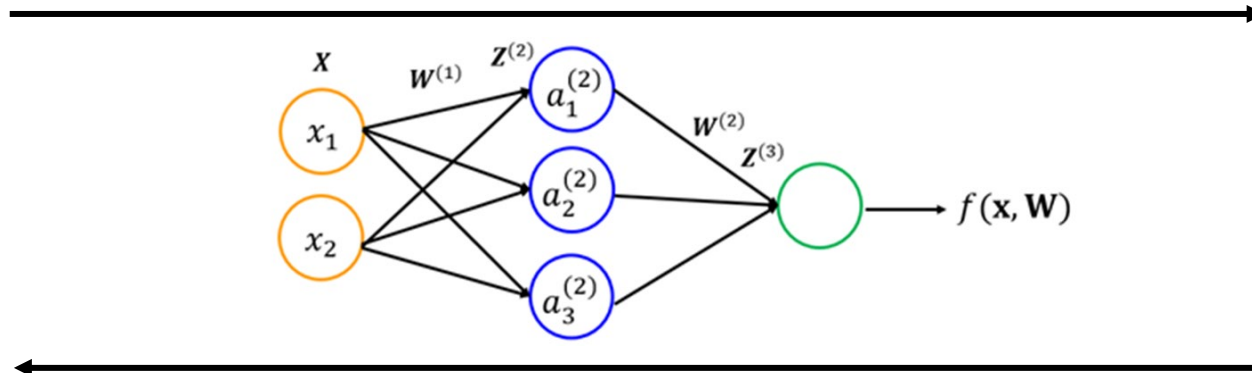
$$\delta^{(3)} = (f(\mathbf{X}, \mathbf{W}) - \mathbf{t}) .* \sigma'\left(\mathbf{Z}^{(3)}\right)$$

- Now matrix multiplication does the sum job we need on slide 19.

- In the real implementation, first we calculate $\delta^{(3)}$, then calculate $\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}}$

Forward propagation



Backpropagation

$$\delta^{(3)} = \begin{bmatrix} \delta_1^{(3)} \\ \delta_2^{(3)} \\ \delta_3^{(3)} \\ \delta_4^{(3)} \end{bmatrix}$$

- Given training examples, first run a "forward propagation" to compute all the activations throughout the network, including the output value of the NN $f(\mathbf{X}, \mathbf{W})$

- Then, for each node $i$ in layer $l$, compute an "error term" $\delta_i^{(l)}$ that measures how much that node was "**responsible**" for any errors in our output

Hidden layer size

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \end{bmatrix}$$

Input unit size

$$\frac{\partial \mathrm{L}(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \begin{bmatrix} \dfrac{\partial \mathrm{L}(\mathbf{W})}{\partial w_{11}^{(1)}} & \dfrac{\partial \mathrm{L}(\mathbf{W})}{\partial w_{12}^{(1)}} & \dfrac{\partial \mathrm{L}(\mathbf{W})}{\partial w_{13}^{(1)}} \\ \dfrac{\partial \mathrm{L}(\mathbf{W})}{\partial w_{21}^{(1)}} & \dfrac{\partial \mathrm{L}(\mathbf{W})}{\partial w_{22}^{(1)}} & \dfrac{\partial \mathrm{L}(\mathbf{W})}{\partial w_{23}^{(1)}} \end{bmatrix}$$

**How to calculate this?**

Let's follow the similar process as step 1

Each ⊛ has its own meaning

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \frac{\partial \frac{1}{2}(f(\mathbf{x}_n, \mathbf{W}) - t_n)^2}{\partial \mathbf{W}^{(1)}} = (f(\mathbf{x}_n, \mathbf{W}) - t_n) \circledast \frac{\partial f(\mathbf{x}_n, \mathbf{W})}{\partial \mathbf{W}^{(1)}}$$

$$= (f(\mathbf{x}_n, \mathbf{W}) - t_n) \circledast \frac{\partial f(\mathbf{x}_n, \mathbf{W})}{\partial \mathbf{z}_n^{(3)}} \circledast \frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(1)}}$$

**We calculated this in step 1** →

$$= \boxed{(f(\mathbf{x}_n, \mathbf{W}) - t_n) \circledast \frac{\partial \boldsymbol{\sigma}(\mathbf{z}_n^{(3)})}{\partial \mathbf{z}_n^{(3)}}} \circledast \frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(1)}}$$

**Note the difference compared to step 1**

$$= (f(\mathbf{x}_n, \mathbf{W}) - t_n) .* \boldsymbol{\sigma}'\left(\mathbf{z}_n^{(3)}\right) \circledast \frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(1)}}$$

$$= \boxed{\delta_n^{(3)}} \circledast \boxed{\frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(1)}}}$$

**How to calculate this?**

Poll 5

Continued…

For simplicity we skip those ⊛



$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}}$$

$$= \delta_n^{(3)} \circledast \boxed{\frac{\partial \mathbf{z}_n^{(3)}}{\partial \mathbf{W}^{(1)}}}$$

**Derivative chain rule.**

**Based on FP, $W^{(2)}$, $a_n^{(2)}$, $z_n^{(2)}$ are between $z_n^{(3)}$ and $W^{(1)}$**

$$\frac{\partial \mathbf{z}_n^{(3)}}{\partial a_n^{(2)}} = \frac{\partial a_n^{(2)} \mathbf{W}^{(2)}}{\partial a_n^{(2)}} = \mathbf{W}^{(2)} \otimes \mathbf{I}$$

$$= \delta_n^{(3)} \circledast \frac{\partial \mathbf{z}_n^{(3)}}{\partial a_n^{(2)}} \circledast \frac{\partial a_n^{(2)}}{\partial \mathbf{W}^{(1)}} = \delta_n^{(3)} \circledast \frac{\partial \mathbf{z}_n^{(3)}}{\partial a_n^{(2)}} \circledast \frac{\partial a_n^{(2)}}{\partial \mathbf{z}_n^{(2)}} \circledast \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(1)}}$$

$$= \delta_n^{(3)} \left(\mathbf{W}^{(2)}\right)^T \circledast \boxed{\frac{\partial a_n^{(2)}}{\partial \mathbf{z}_n^{(2)}}} \circledast \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(1)}} = \delta_n^{(3)} \left(\mathbf{W}^{(2)}\right)^T \circledast \frac{\partial \sigma(\mathbf{z}_n^{(2)})}{\partial \mathbf{z}_n^{(2)}} \circledast \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(1)}}$$

$$\frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(1)}} = \frac{\partial a_n^{(1)} \mathbf{W}^{(1)}}{\partial \mathbf{W}^{(1)}} = I \otimes a_n^{(1)}$$

$$= \delta_n^{(3)} \left(\mathbf{W}^{(2)}\right)^T .* \sigma'(\mathbf{z}_n^{(2)}) \circledast \frac{\partial \mathbf{z}_n^{(2)}}{\partial \mathbf{W}^{(1)}}$$

**Use the same strategy as step 1 to calculate this.**

For one data
$a_n^{(1)} = \mathbf{x}_n$ (row of $\mathbf{X}$)

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \left(a_n^{(1)}\right)^T \left(\delta_n^{(3)}(\mathbf{W}^{(2)})^T\right) .* \sigma'\left(\mathbf{z}_n^{(2)}\right)$$

$$= \boxed{(\mathbf{x}_n)^T} \left(\delta_n^{(3)}(\mathbf{W}^{(2)})^T\right) .* \sigma'\left(\mathbf{z}_n^{(2)}\right)$$

For All data:

$$N \times 3$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = (\mathbf{X})^T \underbrace{\left(\delta^{(3)}(\mathbf{W}^{(2)})^T\right) .* \sigma'(\mathbf{Z}^{(2)})}_{\delta^{(2)}}$$

Poll 6

**Step 2: summary for all data (e.g. N=3)**

$$2 \times 3$$
$$2 \times 4 \quad 4 \times 3$$
$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \boxed{\mathbf{X}^T \delta^{(2)}}$$

- We removed summation from loss function calculation previously, but now matrix multiplication also does the sum job.

$$\delta^{(2)} = \left(\delta^{(3)}(\mathbf{W}^{(2)})^T\right) .* \sigma'(\mathbf{Z}^{(2)})$$
$$4 \times 3 \quad 4 \times 1 \ 1 \times 3 \quad\quad 4 \times 3$$

- In the real implementation, first we calculate $\delta^{(2)}$, then calculate $\frac{\partial L(W)}{\partial W^{(1)}}$

**Step 1: summary**

$$\delta^{(3)} = (f(\mathbf{X}, \mathbf{W}) - \mathbf{t}).* \sigma'\big(\mathbf{Z}^{(3)}\big)$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \big(\boldsymbol{a}^{(2)}\big)^T \delta^{(3)}$$

Note the where the matrix element wise product " .* " is used
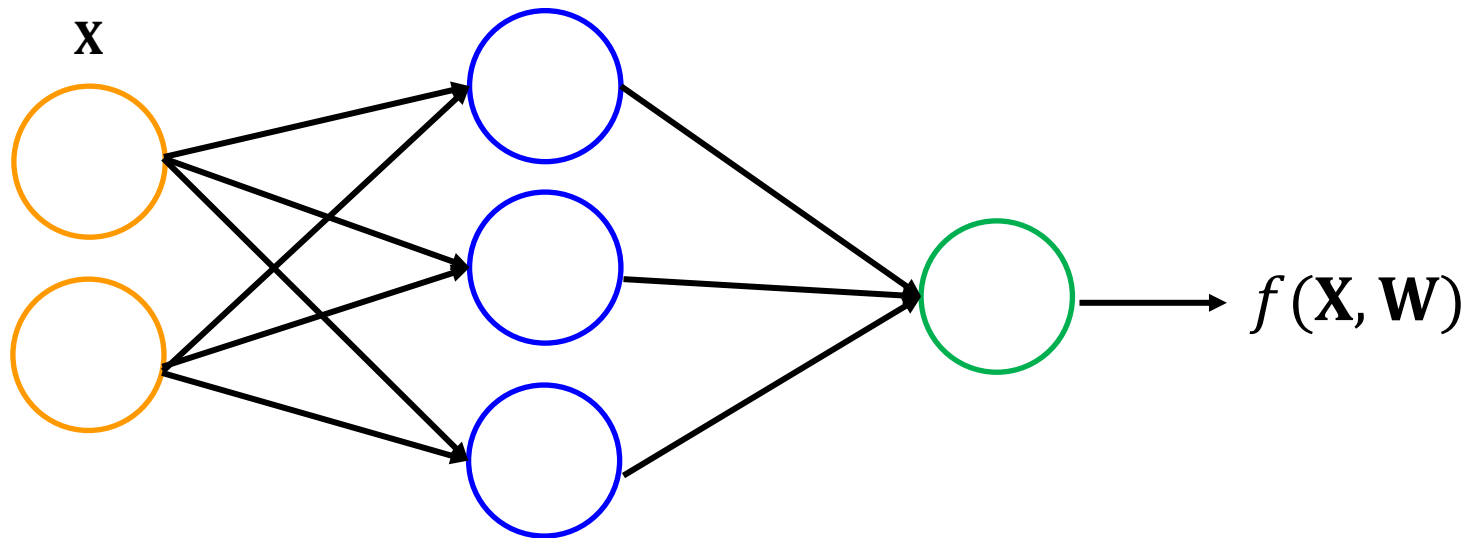
**Step 2: summary**

$$\delta^{(2)} = \Big(\delta^{(3)}\big(\mathbf{W}^{(2)}\big)^T\Big).* \sigma'\big(\mathbf{Z}^{(2)}\big)$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \mathbf{X}^T \delta^{(2)} = \big(\boldsymbol{a}^{(1)}\big)^T \delta^{(2)}$$

Based on the layer index, we can see that this process can be easily repeated with more layers.

Forward propagation



$\delta^{(2)}$ is a $4 \times 3$ matrix

Backpropagation
`Example04_Example01.py`

Calculate gradient

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \begin{bmatrix} \frac{\partial L(\mathbf{W})}{\partial w_{11}^{(1)}} & \frac{\partial L(\mathbf{W})}{\partial w_{12}^{(1)}} & \frac{\partial L(\mathbf{W})}{\partial w_{13}^{(1)}} \\ \frac{\partial L(\mathbf{W})}{\partial w_{21}^{(1)}} & \frac{\partial L(\mathbf{W})}{\partial w_{22}^{(1)}} & \frac{\partial L(\mathbf{W})}{\partial w_{23}^{(1)}} \end{bmatrix}$$

**$2 \times 3$**

```
In [253]: dl_dW1
Out[253]:
array([[-0.01194708,  0.00497217,  0.00505324],
       [-0.00600021,  0.00238656,  0.00351687]])
```

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \begin{bmatrix} \frac{\partial L(\mathbf{W})}{\partial w_{11}^{(2)}} \\ \frac{\partial L(\mathbf{W})}{\partial w_{21}^{(2)}} \\ \frac{\partial L(\mathbf{W})}{\partial w_{31}^{(2)}} \end{bmatrix}$$

**$3 \times 1$**

## Forward Propagation

$$a^{(1)} = \mathbf{X}$$

$$\mathbf{Z}^{(2)} = \mathbf{X}\mathbf{W}^{(1)}$$

$$a^{(2)} = \sigma(\mathbf{Z}^{(2)})$$

$$\mathbf{Z}^{(3)} = a^{(2)}\mathbf{W}^{(2)}$$

$$f(\mathbf{X}, \mathbf{W}) = \mathbf{a}^{(3)} = \sigma(\mathbf{Z}^{(3)})$$

Identity

## Backward Propagation

$$\delta^{(3)} = (f(\mathbf{X}, \mathbf{W}) - \mathbf{t}).* \boxed{\sigma'(\mathbf{Z}^{(3)})}$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = (a^{(2)})^T \delta^{(3)}$$

=1

$$\delta^{(2)} = \left(\delta^{(3)}(\mathbf{W}^{(2)})^T\right).* \sigma'(\mathbf{Z}^{(2)})$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \mathbf{X}^T \delta^{(2)} = (a^{(1)})^T \delta^{(2)}$$

## Calculate loss function

$$L(\mathbf{W}) = \frac{1}{2}\sum_{n=1}^{N}(f(\mathbf{x}_n, \mathbf{W}) - t_n)^2 = \frac{1}{2}(f(\mathbf{X}, \mathbf{W}) - \mathbf{t})^T(f(\mathbf{X}, \mathbf{W}) - \mathbf{t})$$

Use " .* " to denote the element-wise product operator

`Lecture04_Example02.py`

36

## Forward Propagation

$$a^{(1)} = \mathbf{X}$$

$$\mathbf{Z}^{(2)} = \mathbf{X}\mathbf{W}^{(1)}$$

$$a^{(2)} = \sigma(\mathbf{Z}^{(2)})$$

$$\mathbf{Z}^{(3)} = a^{(2)}\mathbf{W}^{(2)}$$

$$f(\mathbf{X}, \mathbf{W}) = a^{(3)} = \sigma(\mathbf{Z}^{(3)})$$

sigmoid

## Backward Propagation

$$\delta^{(3)} = (f(\mathbf{X}, \mathbf{W}) - \mathbf{t}) \quad \text{[t in 1 or 0 code]}$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(2)}} = \left(a^{(2)}\right)^T \delta^{(3)}$$

$$\delta^{(2)} = \left(\delta^{(3)}\left(\mathbf{W}^{(2)}\right)^T\right) .* \sigma'\left(\mathbf{Z}^{(2)}\right)$$

$$\frac{\partial L(\mathbf{W})}{\partial \mathbf{W}^{(1)}} = \mathbf{X}^T \delta^{(2)} = \left(a^{(1)}\right)^T \delta^{(2)}$$
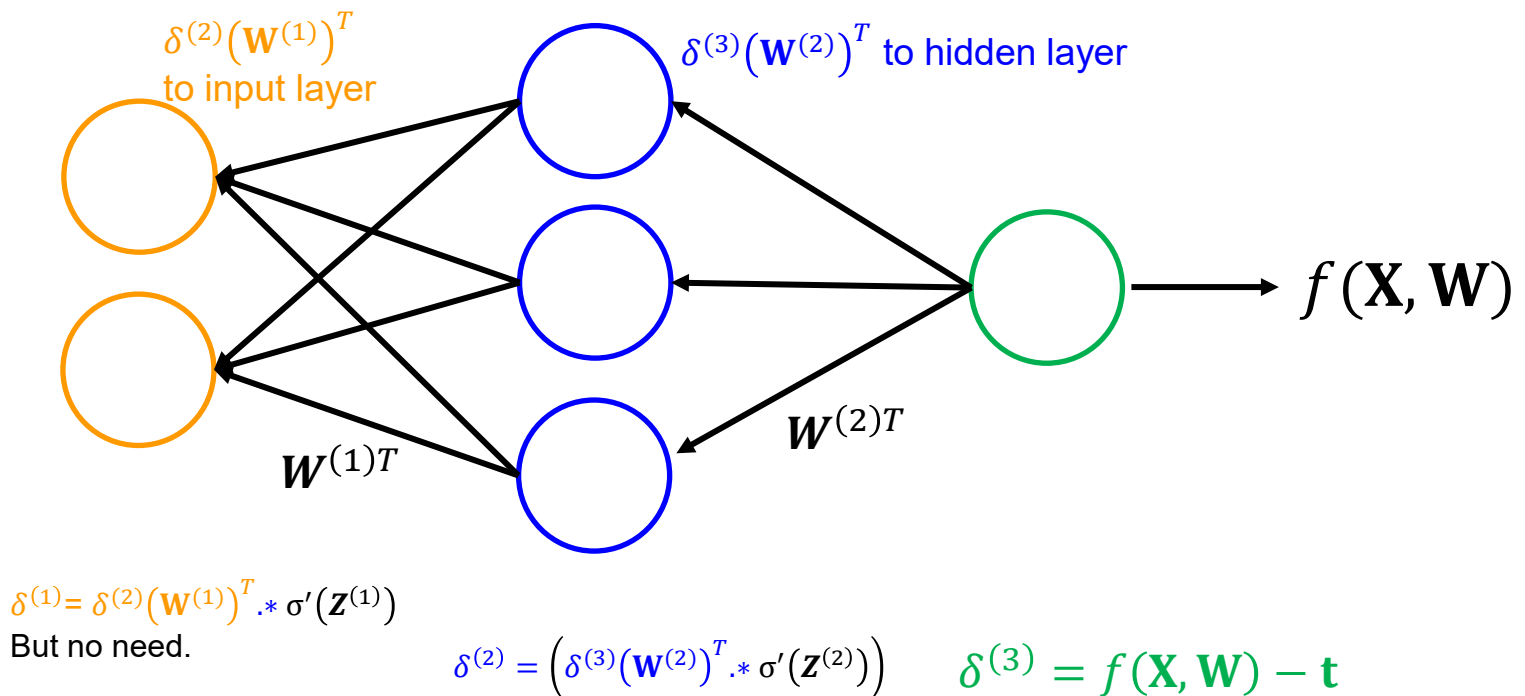
Calculate (logistic) loss function (1/N removed)

$$L(\mathbf{W}) = -\left[\sum_{n=1}^{N} \left(t_n \log(f(\mathbf{x}_n, \mathbf{W})) + (1 - t_n) \log(1 - f(\mathbf{x}_n, \mathbf{W}))\right)\right]$$

Use " .* " to denote the element-wise product operator

$\delta^{(2)}\big(\mathbf{W}^{(1)}\big)^{T}$
to input layer

$\delta^{(3)}\big(\mathbf{W}^{(2)}\big)^{T}$ to hidden layer

$\boldsymbol{W}^{(2)T}$

$\boldsymbol{W}^{(1)T}$

$f(\mathbf{X}, \mathbf{W})$

$\delta^{(1)} = \delta^{(2)}\big(\mathbf{W}^{(1)}\big)^{T} .* \sigma'(\boldsymbol{Z}^{(1)})$
But no need.

$\delta^{(2)} = \Big( \delta^{(3)}\big(\mathbf{W}^{(2)}\big)^{T} .* \sigma'(\boldsymbol{Z}^{(2)}) \Big)$     $\delta^{(3)} = f(\mathbf{X}, \mathbf{W}) - \mathbf{t}$

$\delta$ is backpropagated on the network

Two types of forward operations
for All Data:

$$\mathbf{Z}^{(l+1)} \rightarrow \sigma(\mathbf{Z}^{(l+1)}) \rightarrow \boldsymbol{a}^{(l+1)}$$

$$\boldsymbol{a}^{(l)} \rightarrow \mathbf{Z}^{(l+1)} = \boldsymbol{a}^{(l)}\mathbf{W}^{(l)} \rightarrow \mathbf{Z}^{(l+1)}$$

$$\mathbf{W}^{(l)}$$

Their BP:

$$\frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} = \sigma'(\mathbf{Z}^{(l+1)}).* \frac{\partial L}{\partial \mathbf{a}^{(l+1)}}$$

$$\frac{\partial L}{\partial \mathbf{a}^{(l+1)}}$$

$$\sigma'(\mathbf{Z}^{(l+1)})$$

Denoted by $\delta^{(l+1)}$

$$\frac{\partial L}{\partial \mathbf{a}^{(l)}} = \frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} \mathbf{W}^{(l)T}$$

$$\frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{W}^{(l)}} = \mathbf{I} \otimes \mathbf{a}^{(l)}$$

$$\frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{a}^{(l)}} = \mathbf{W}^{(l)} \otimes \mathbf{I}$$

$$\frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} = \delta^{(l+1)}$$

$$\frac{\partial L}{\partial \mathbf{W}^{(l)}} = (\mathbf{a}^{(l)})^T \frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} = (\mathbf{a}^{(l)})^T \delta^{(l+1)}$$