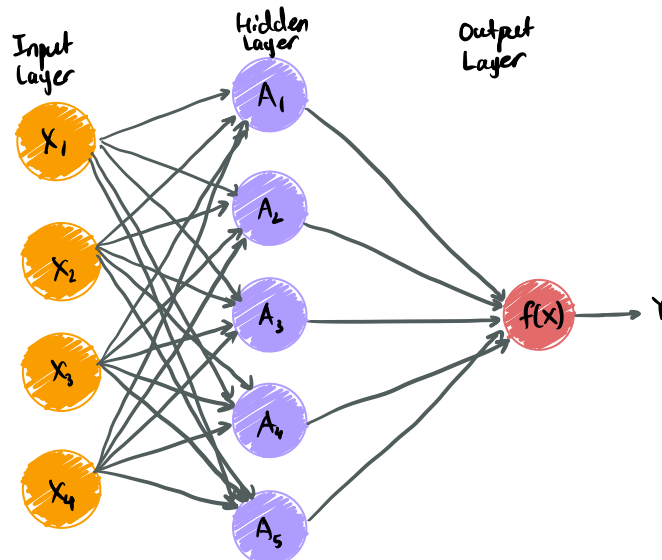## Neural Learning

Recap: The Perceptron



inputs        weights
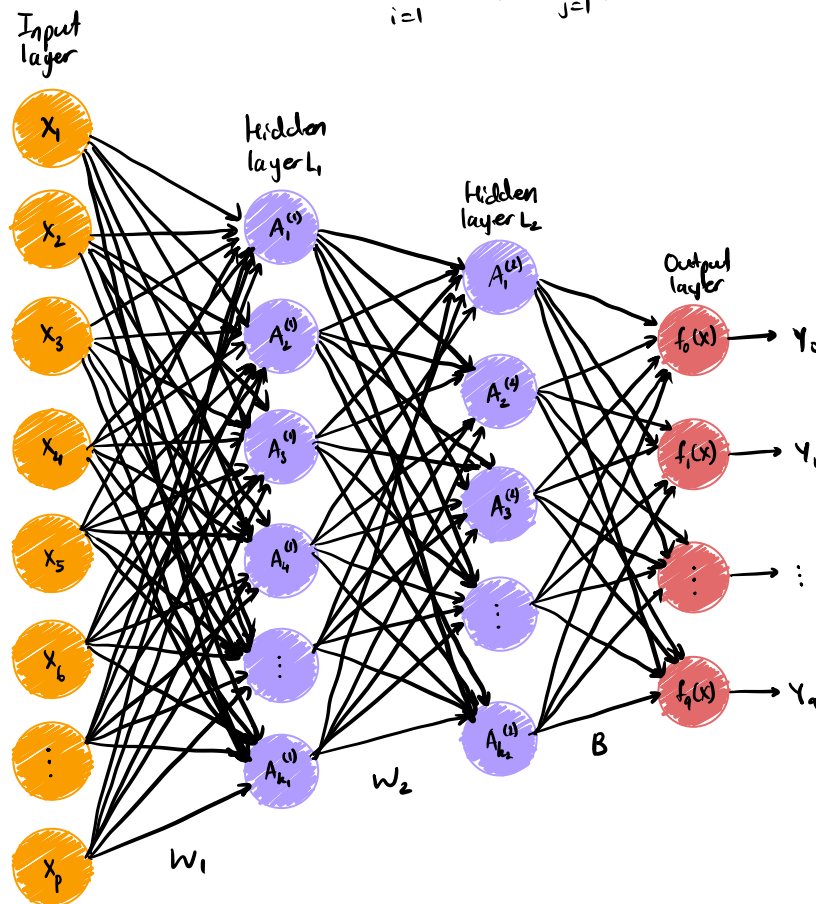
Activation function

## Multi-layer Perceptron

- A multi-layer perceptron is where we chain the perceptrons to learn non-linear patterns



If we define the activation function used for the hidden layer as $g$ and the weights for

input features as $\beta$:

$$f(X) = w_0 + \sum_{i=1}^{n} w_i A_i$$

$$= w_0 + \sum_{i=1}^{n} w_i g(X_i)$$

$$= w_0 + \sum_{i=1}^{n} w_i g\left(\beta_0 + \sum_{j=1}^{p} \beta_j X_i\right)$$

Input layer

Hidden layer $L_1$

Hidden layer $L_2$

Output layer

$X_1$ $X_2$ $X_3$ $X_4$ $X_5$ $X_6$ $\vdots$ $X_p$

$A_1^{(1)}$ $A_2^{(1)}$ $A_3^{(1)}$ $A_4^{(1)}$ $\vdots$ $A_{k_1}^{(1)}$

$A_1^{(2)}$ $A_2^{(2)}$ $A_3^{(2)}$ $\vdots$ $A_{k_2}^{(2)}$

$f_0(x)$ $f_1(x)$ $\vdots$ $f_q(x)$

$Y_0$ $Y_1$ $\vdots$ $Y_q$

$W_1$ $W_2$ $B$

## Back — propagation

The main problem now becomes: How do we learn this large number of weights?
As always, we define an appropriate loss function and optimise it. Due to the
complexity of the function which is the neural network, we'll need to perform gradient
descent to gradually improve our model over time
But how do we calculate the gradient?
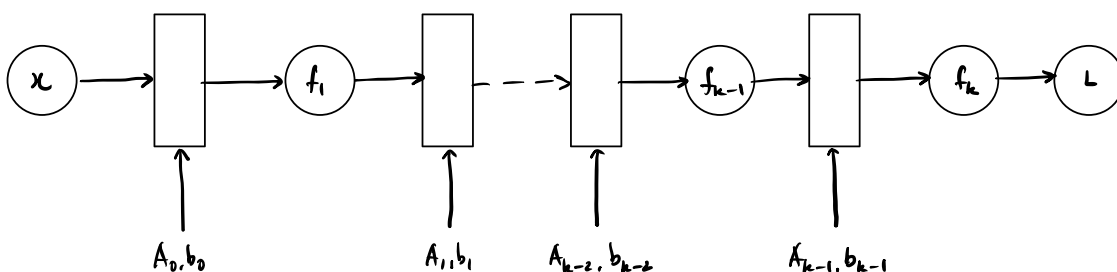
The loss function of the activation:

$$L(a, y)$$

The activation is a function of the inputs, the weights and the bias:

$$a(x_1, \ldots, x_n, w_0, \ldots, w_n)$$

What we want to optimise the loss function is:

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial w_i}$$

Say we have the following network architecture, where A represents the weights and b the bias:



If we say that $\theta_k = \{A_k, b_k\}$ for a layer k. The gradient of our coefficients looks like this:

$$\frac{\partial L}{\partial \theta_{k-1}} = \frac{\partial L}{\partial f_k} \frac{\partial f_k}{\partial f_{k-1}} \frac{\partial f_{k-1}}{\partial f_{k-2}} \frac{\partial f_{k-2}}{\partial \theta_{k-3}}$$

Using back-propagation, we can therefore calculate:

$$\nabla L(\theta, y) = \left[ \frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \cdots, \frac{\partial L}{\partial \theta_k} \right]$$

We can then all of our parameters (in the basic case):

$$\theta^{(t)} = \theta^{(t-1)} - \nabla L(\theta, y)$$

Typically, the optimiser will be some form of stochastic gradient descent (minibatch in some cases) as classic gradient descent is expensive for a large number of parameters and data points