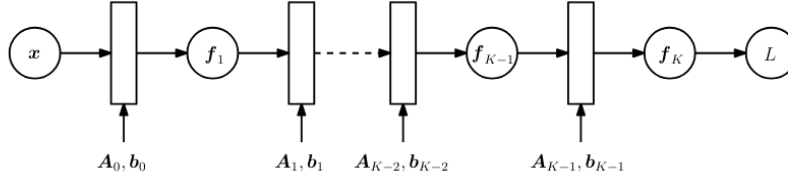


# Backpropagation

*Jack David Carson - February 1, 2025*

Finding good parameters for a machine learning model relies on the fact we can compute the gradient of a learning objective with respect to the parameters of a model. The `backpropagation` algorithm is an efficient way to compute the gradient of an error function with respect to the model parameters.

In deep learning, the chain rule is used to the extreme where the function  $\mathbf{y}$  is computed as a many level function composition  $\mathbf{y} = (f_K \circ f_{K-1} \circ \dots \circ f_1)(\mathbf{x})$  for  $\mathbf{x}$  inputs and  $\mathbf{y}$  observations.



In a neural network with multiple layers we have functions  $f_i(\mathbf{x}_{i-1}) = \sigma(A_{i-1}\mathbf{x}_{i-1} + \mathbf{b}_{i-1})$  in layer  $i$  for **sigmoid activation function** such as  $\sigma(x) = \frac{1}{1+e^{-x}}$ . For instance, we can have loss function

$$\mathcal{L}(\theta) = \|\mathbf{y} - \mathbf{f}_K(\theta, \mathbf{x})^2\|$$

where  $\theta = \{A_0, \mathbf{b}_0, \dots, A_{K-1}, \mathbf{b}_{K-1}\}$  contains the previous layers of the model. We apply the chain rule to say

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta_{K-1}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \theta_{K-1}} \\ \frac{\partial \mathcal{L}}{\partial \theta_{K-1}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \frac{\partial \mathbf{f}_{K-1}}{\partial \theta_{K-2}} \\ \frac{\partial \mathcal{L}}{\partial \theta_i} &= \frac{\partial \mathcal{L}}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \dots \frac{\partial \mathbf{f}_{i+2}}{\partial \mathbf{f}_{i+1}} \frac{\partial \mathbf{f}_{i+1}}{\partial \theta_i} \end{aligned}$$