

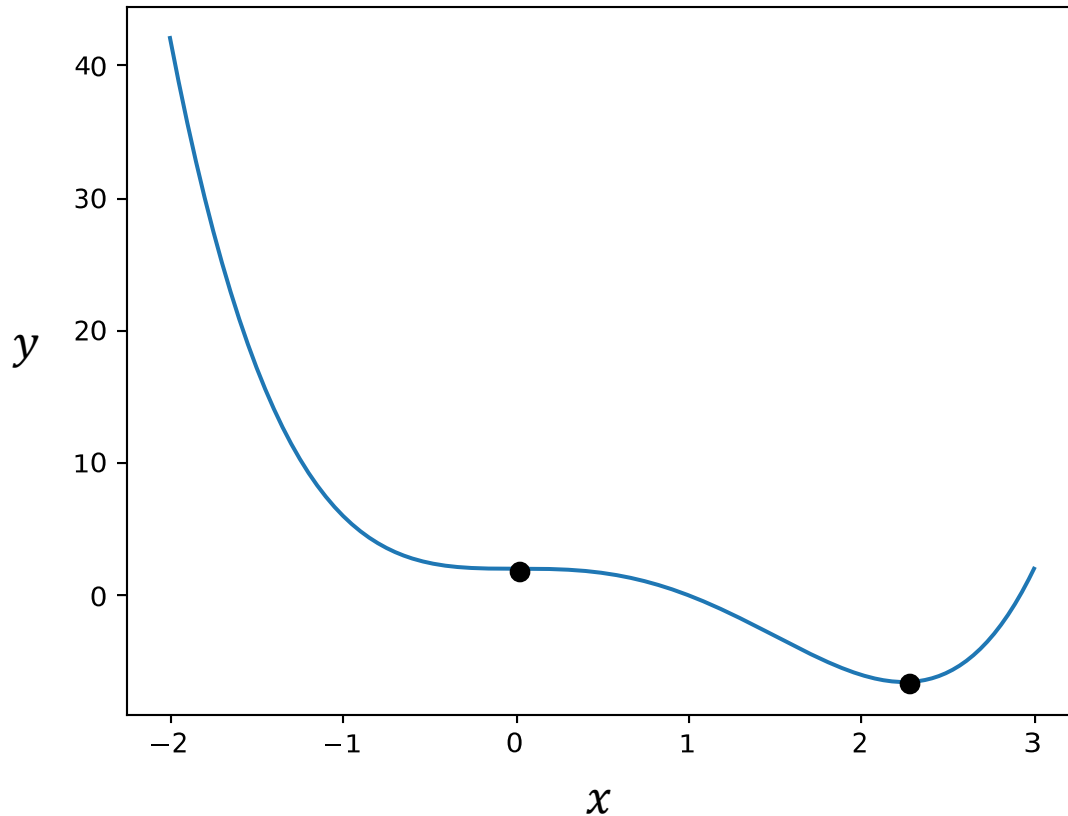
CLIMATE 405

Lecture 10 – Backpropagation and Multilayer Perceptron

Mohammed Ombadi
Assistant Professor
Climate and Space Sciences and Engineering
University of Michigan

Optimization: selecting the best value of variables that will minimize (or maximize) a specific function.

$$y = f(x) = x^4 - 3x^3 + 2$$



For the function on the left, we can find the optimum values of x that minimize the function **analytically**.

$$f'(x) = 4x^3 - 9x^2$$

$$0 = 4x^3 - 9x^2$$

$$4x^3 = 9x^2$$

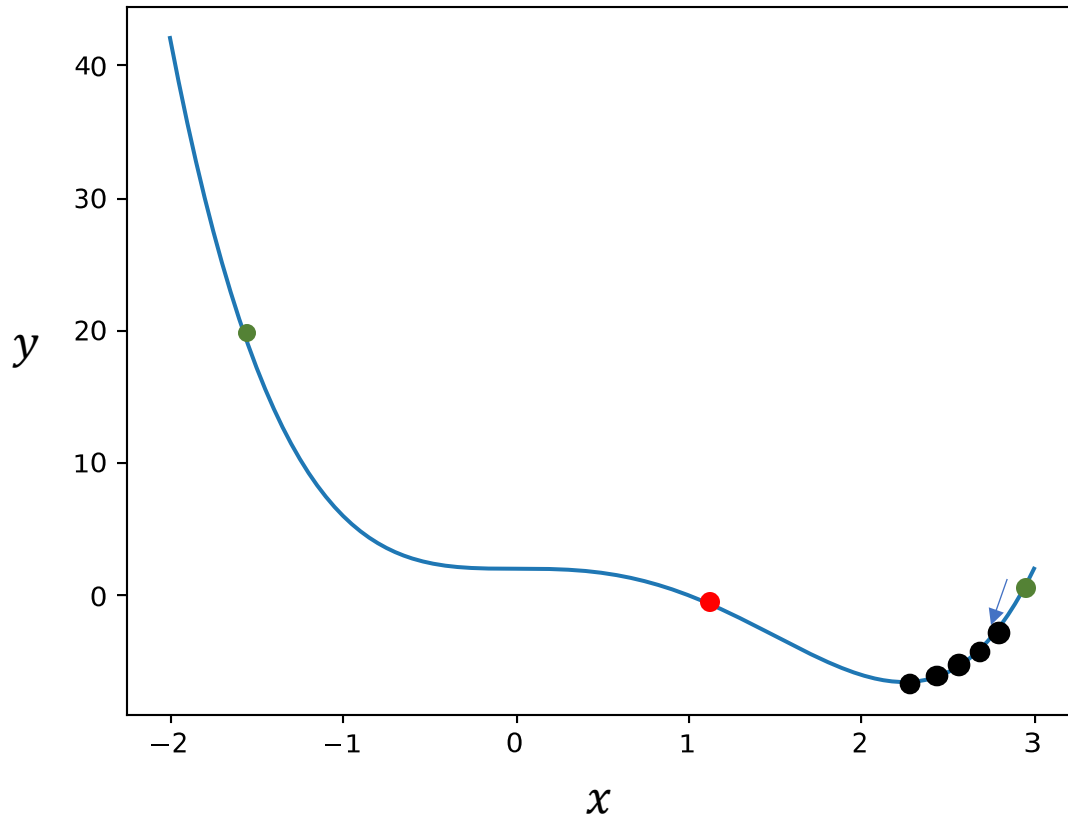
$$\frac{x^3}{x^2} = x = \frac{9}{4} = 2.25$$

$x = 0$ is a local minimum.

Gradient Descent

Gradient Descent: an algorithm that uses the gradient (derivative of a function) to reach a global (local) minimum.

$$y = f(x) = x^4 - 3x^3 + 2$$



Gradient Descent

Inputs: x_{init} , learning rate, num_iterations

Algorithm steps:

$x = x_{\text{init}}$

for i in range(num_iterations):

 calculate $f'(x)$

$x = x - [\text{learning rate} * f'(x)]$

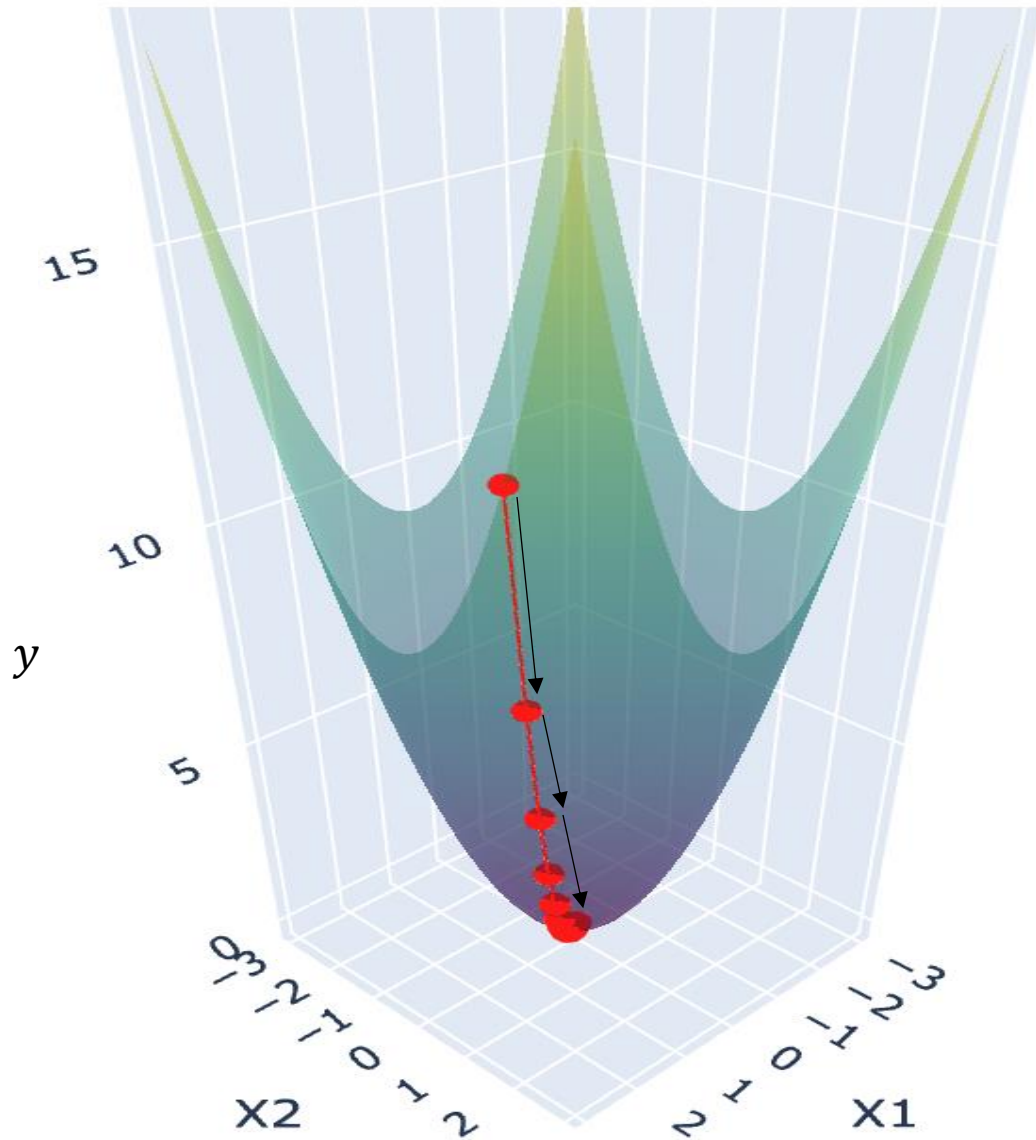
$x_{\text{opt}} = x_{\text{new}}$

In this case, the optimization problem is:

$$\arg \min_{x \in \mathbb{R}} f(x)$$

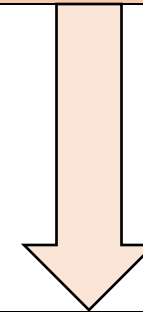
Gradient Descent in 3D

$$y = f(x_1, x_2) = x_1^2 + x_2^2$$



In this case, the optimization problem is:

$$\arg \min_{x_1, x_2 \in \mathbb{R}} f(x_1, x_2)$$



In the context of machine learning, the optimization problem is:

$$\arg \min_{W, b \in \mathbb{R}} f^*(W, b, x, y)$$

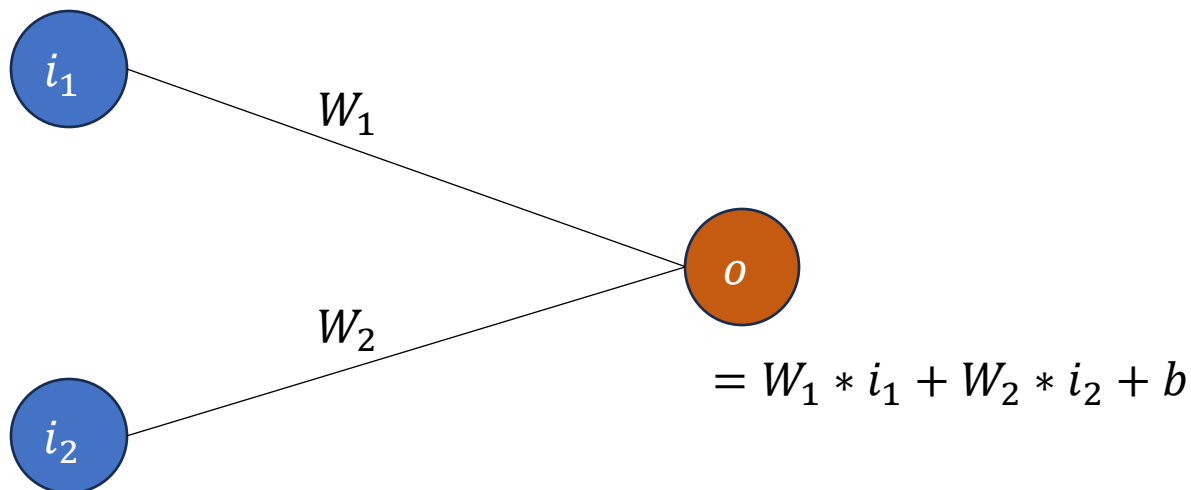
Where $f^*(W, b, X, y) = \sum [y - f(W, b, X)]^2$

In matrix format:

$$[o] = [w_1 \quad w_2] \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} + b$$

Input Layer

Output Layer

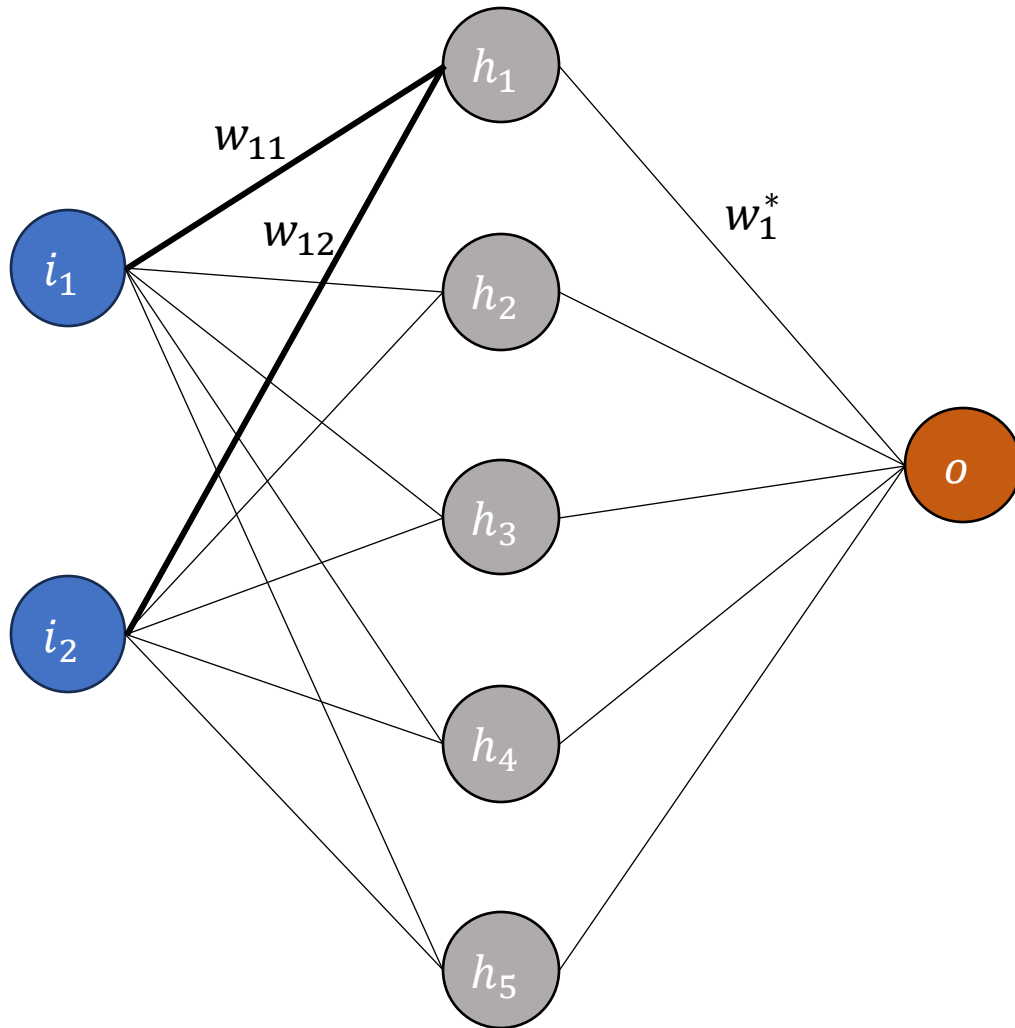


Architecture of a Neural Network: Multilayer Perceptron (MLP)

Input Layer

Hidden Layer

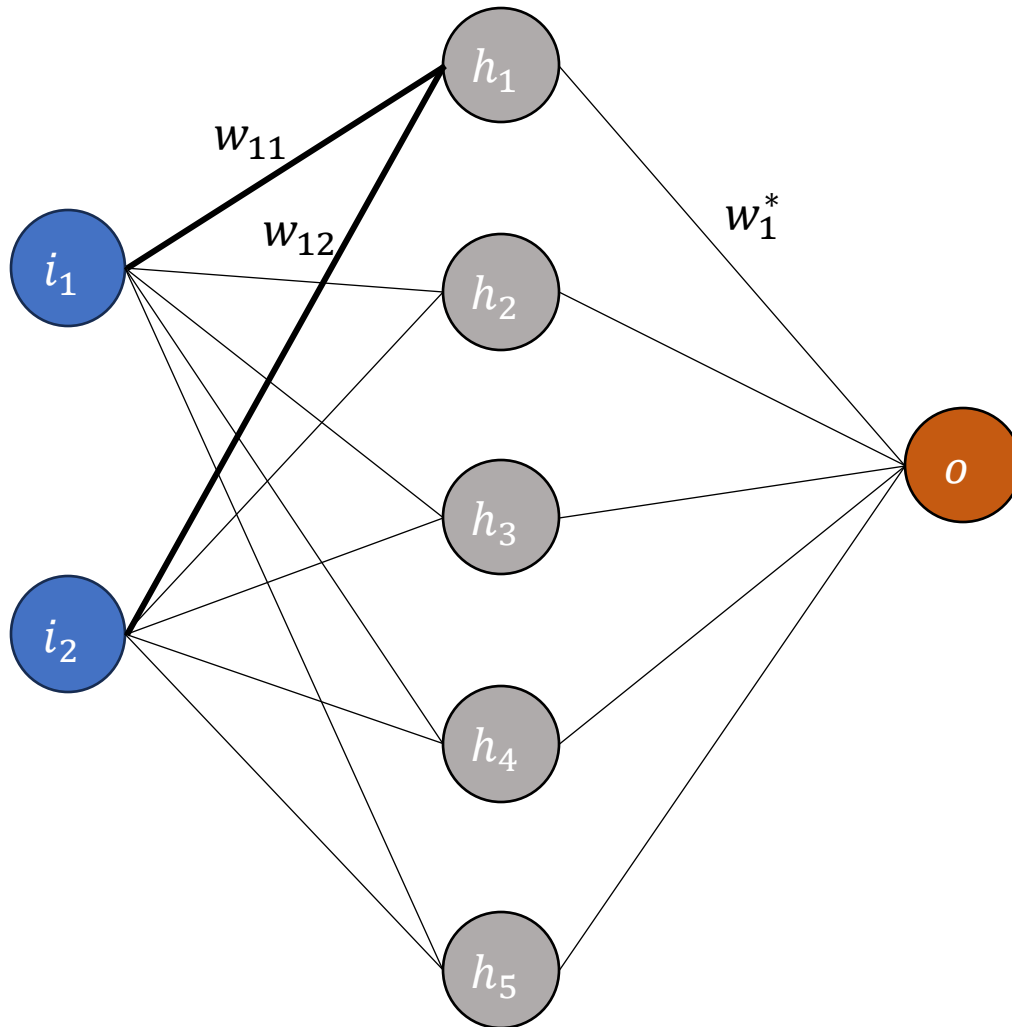
Output Layer



In matrix format:

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \\ w_{41} & w_{42} \\ w_{51} & w_{52} \end{bmatrix} \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{bmatrix}$$

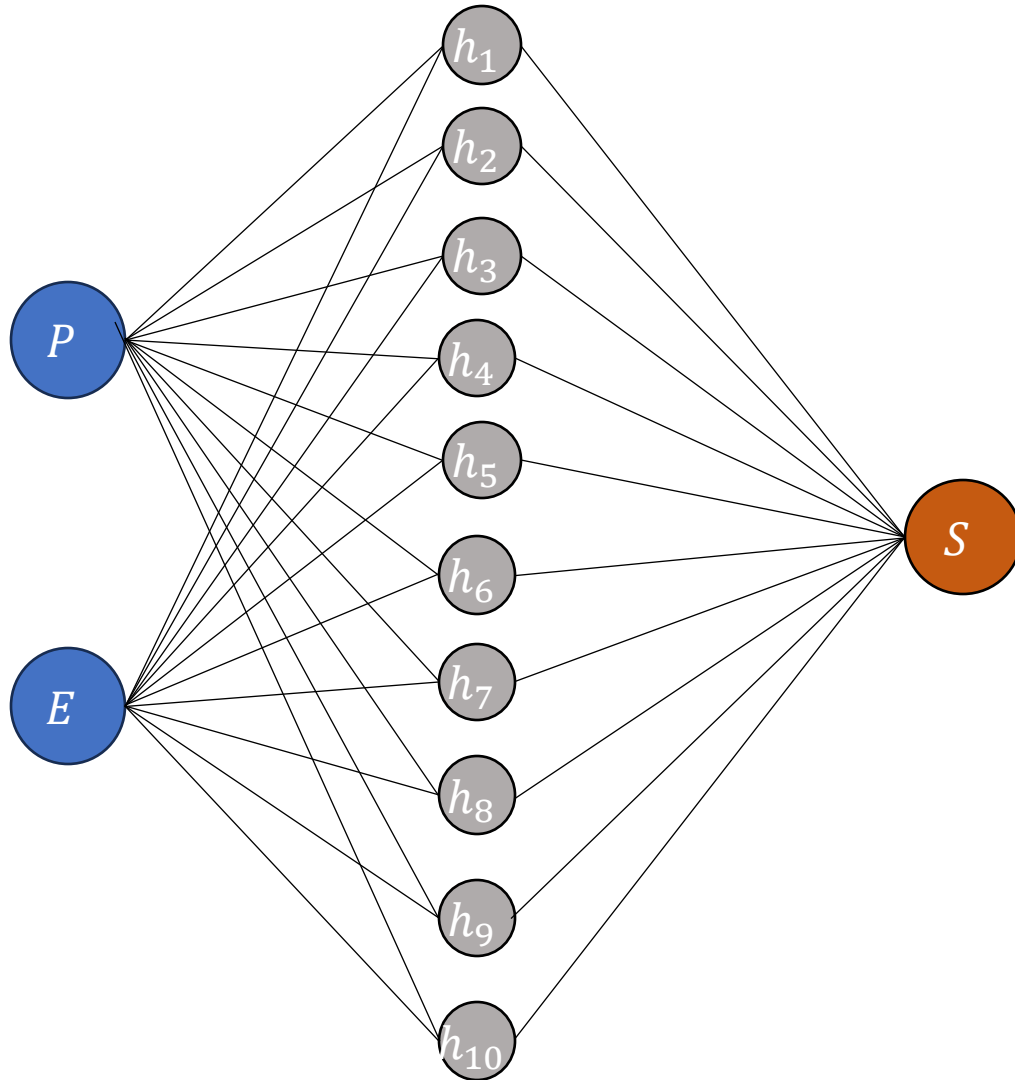
$$[o] = [w_1^* \quad w_2^* \quad w_3^* \quad w_4^* \quad w_5^*] \begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \end{bmatrix} + b^*$$

Input LayerHidden LayerOutput Layer

Number of parameters
 $= (\#Inputs \times \#Nodes_{hidden-layer})$
 $+ (\#Nodes_{hidden-layer} \times \#Outputs)$
 $+ \#Nodes_{hidden-layer} + \#Outputs$

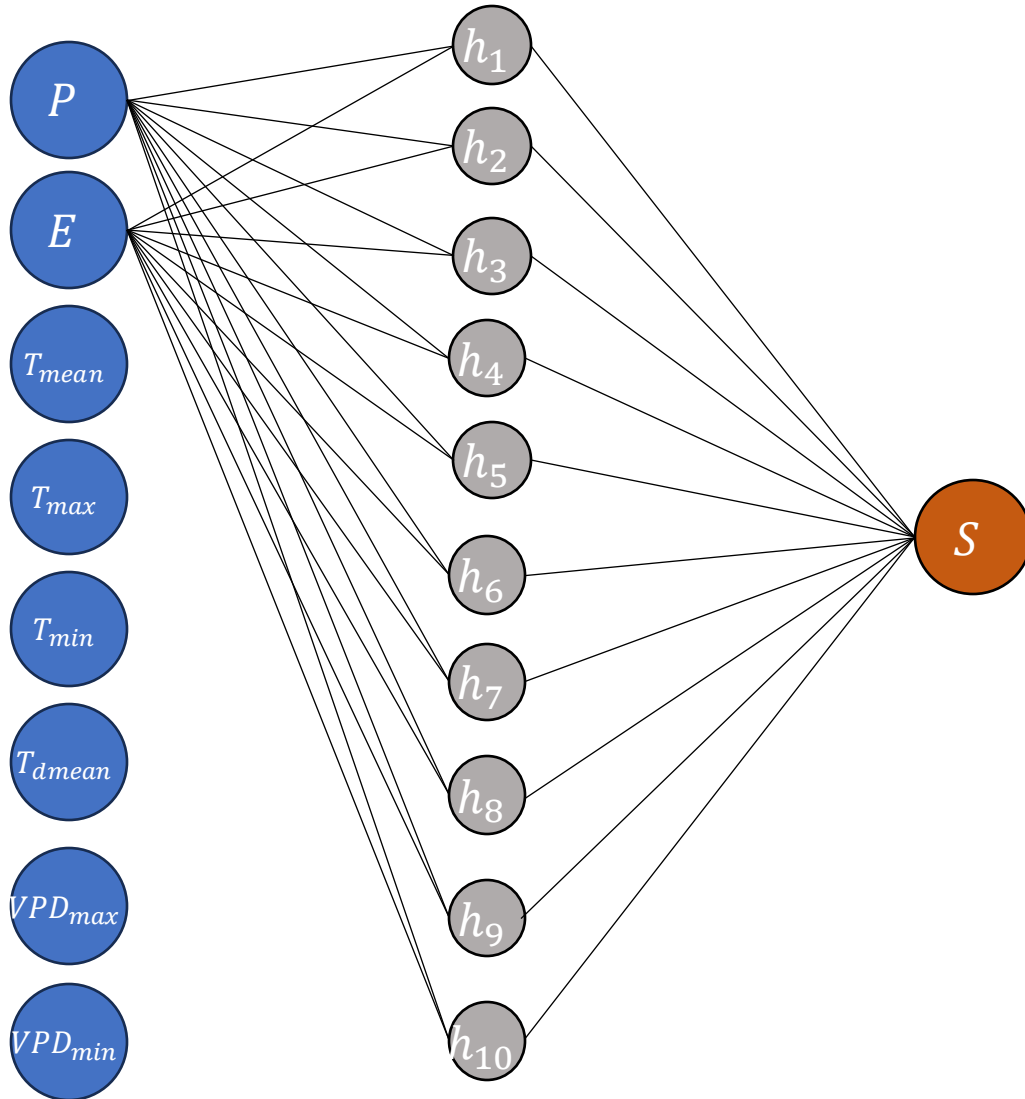
Number of parameters
 $= (2 \times 5) + (5 \times 1) + 5 + 1 = 21$

partial_correlation_data_annual.csv



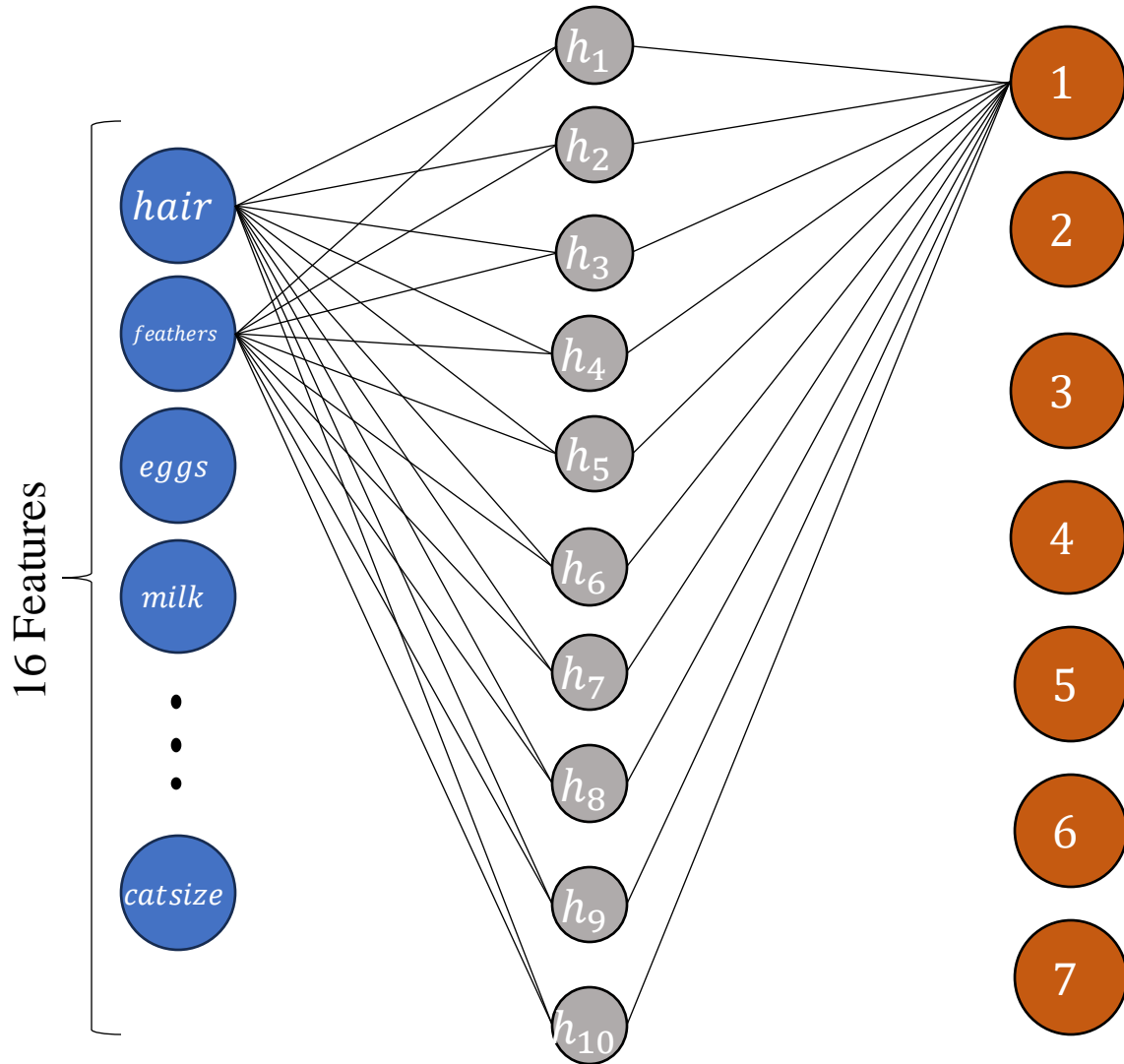
Number of parameters
 $= (2 \times 10) + (10 \times 1) + 10 + 1 = 41$

monthly_meteo_streamflow.csv

*Number of parameters*

$$= (8 \times 10) + (10 \times 1) + 10 + 1 = 101$$

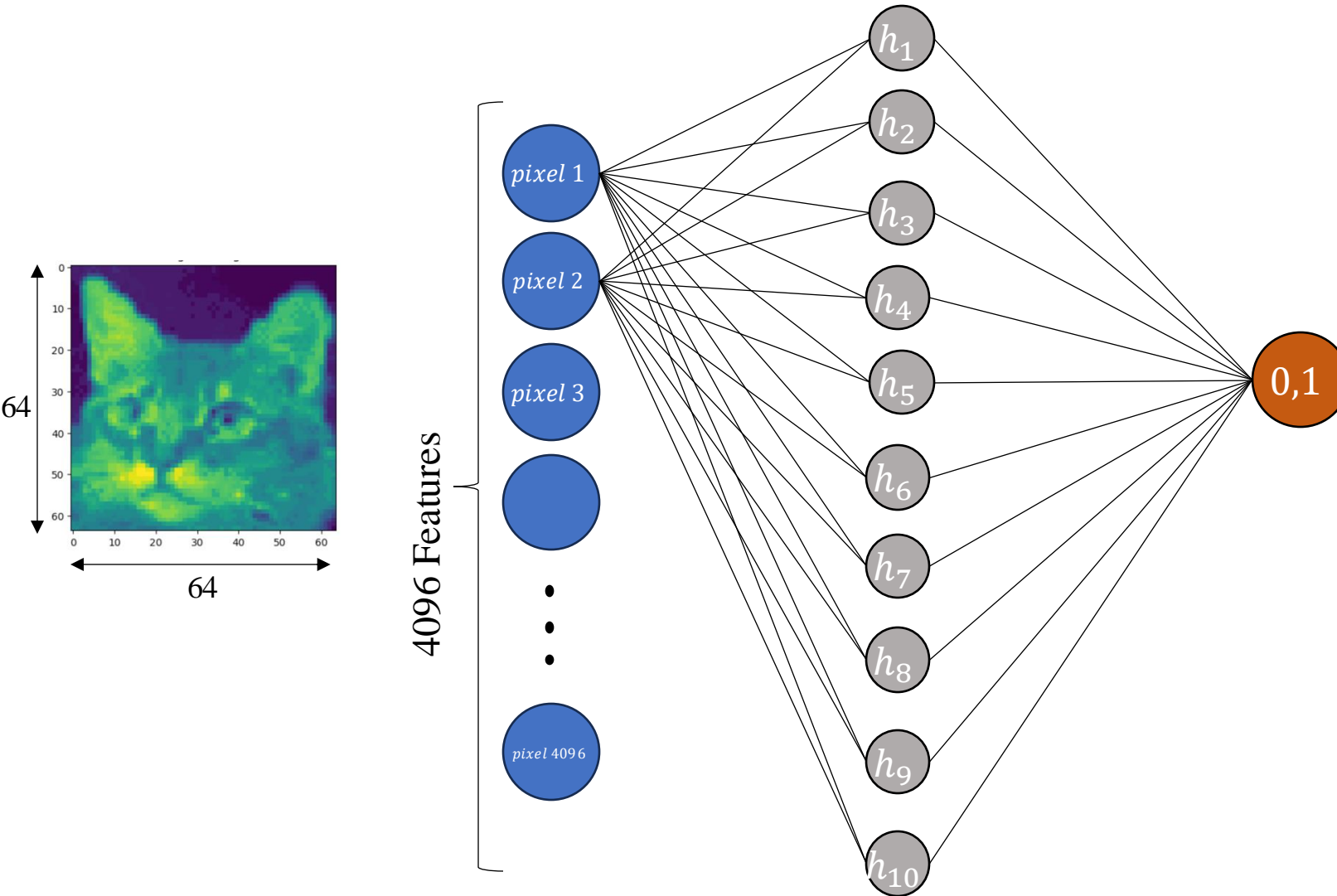
zoo_data.csv



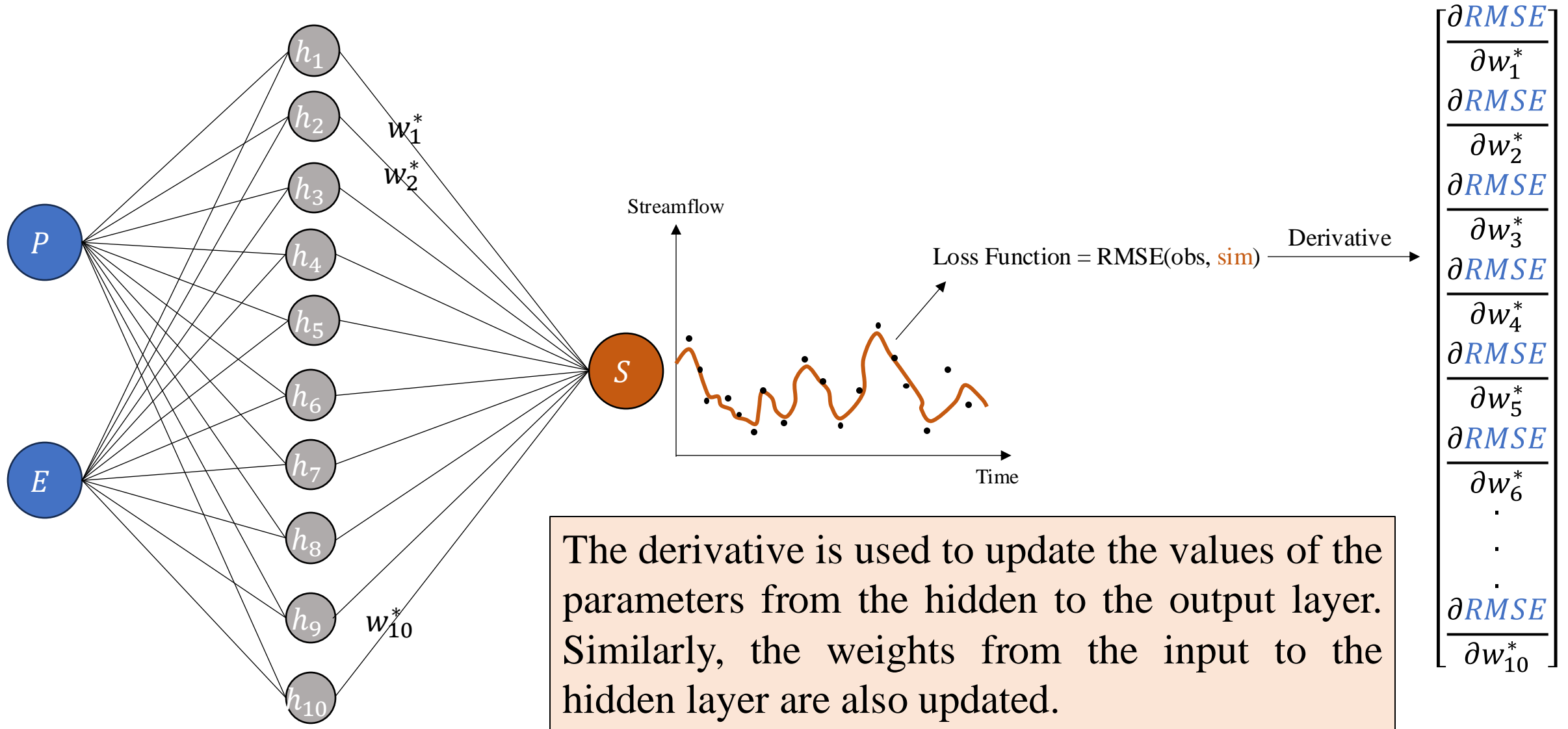
Number of parameters

$$= (16 \times 10) + (10 \times 7) + 10 + 7 = 247$$

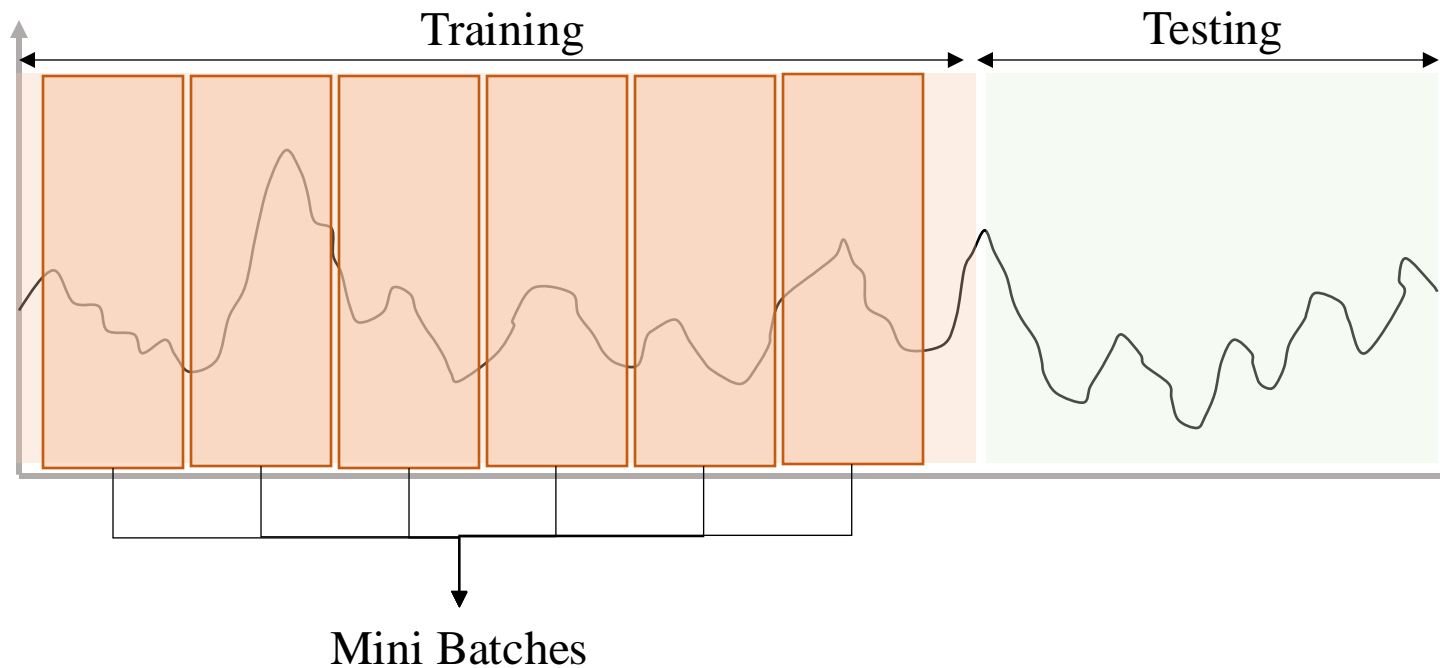
partial_correlation_data_annual.csv

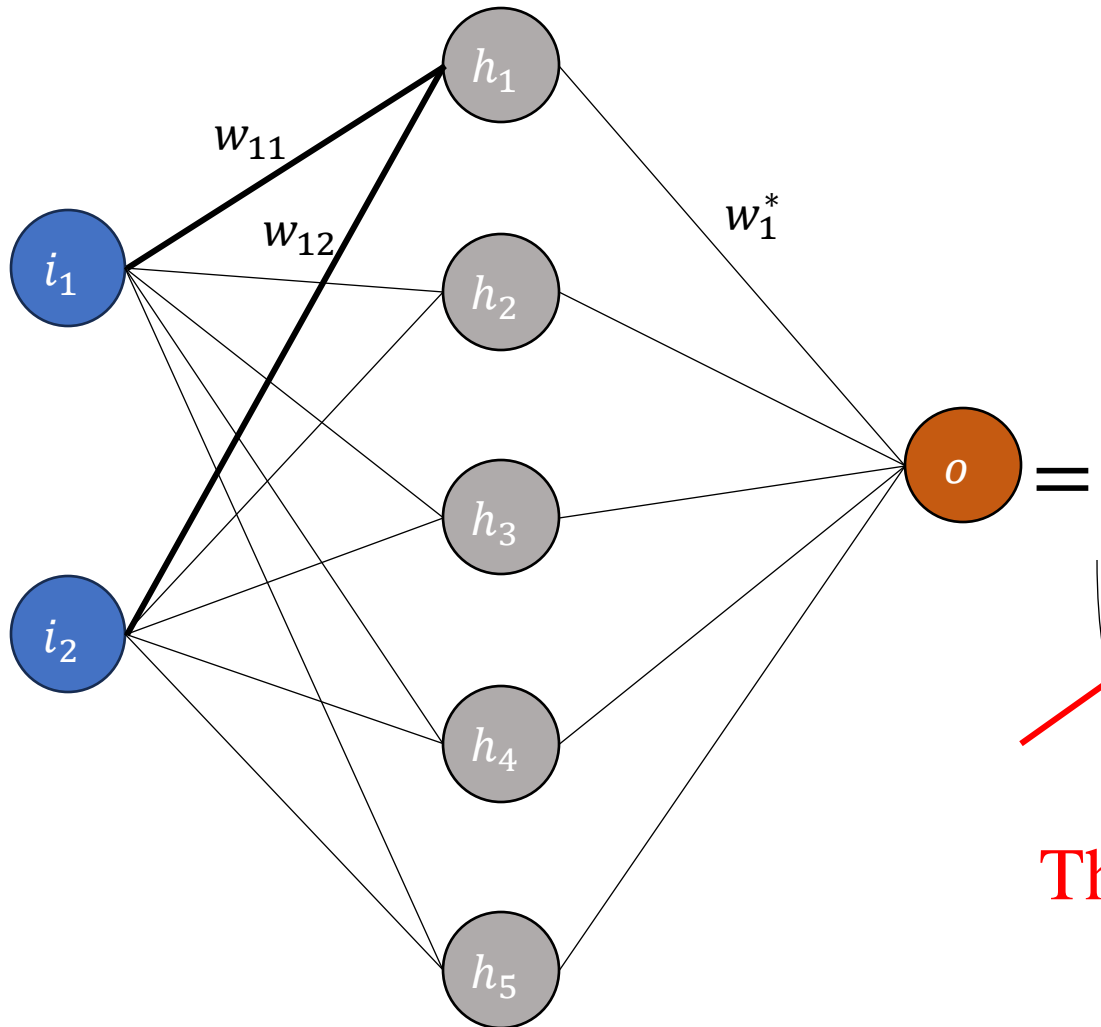


$$\begin{aligned} \text{Number of parameters} &= (4096 * 10) + (10 * 1) \\ &+ 10 + 1 = 40,981 \end{aligned}$$



Stochastic Gradient Descent: a variant of the stochastic gradient optimization algorithm, where the gradient is calculated from *a subset of the training data* instead of the entire training dataset.



Input LayerHidden LayerOutput Layer

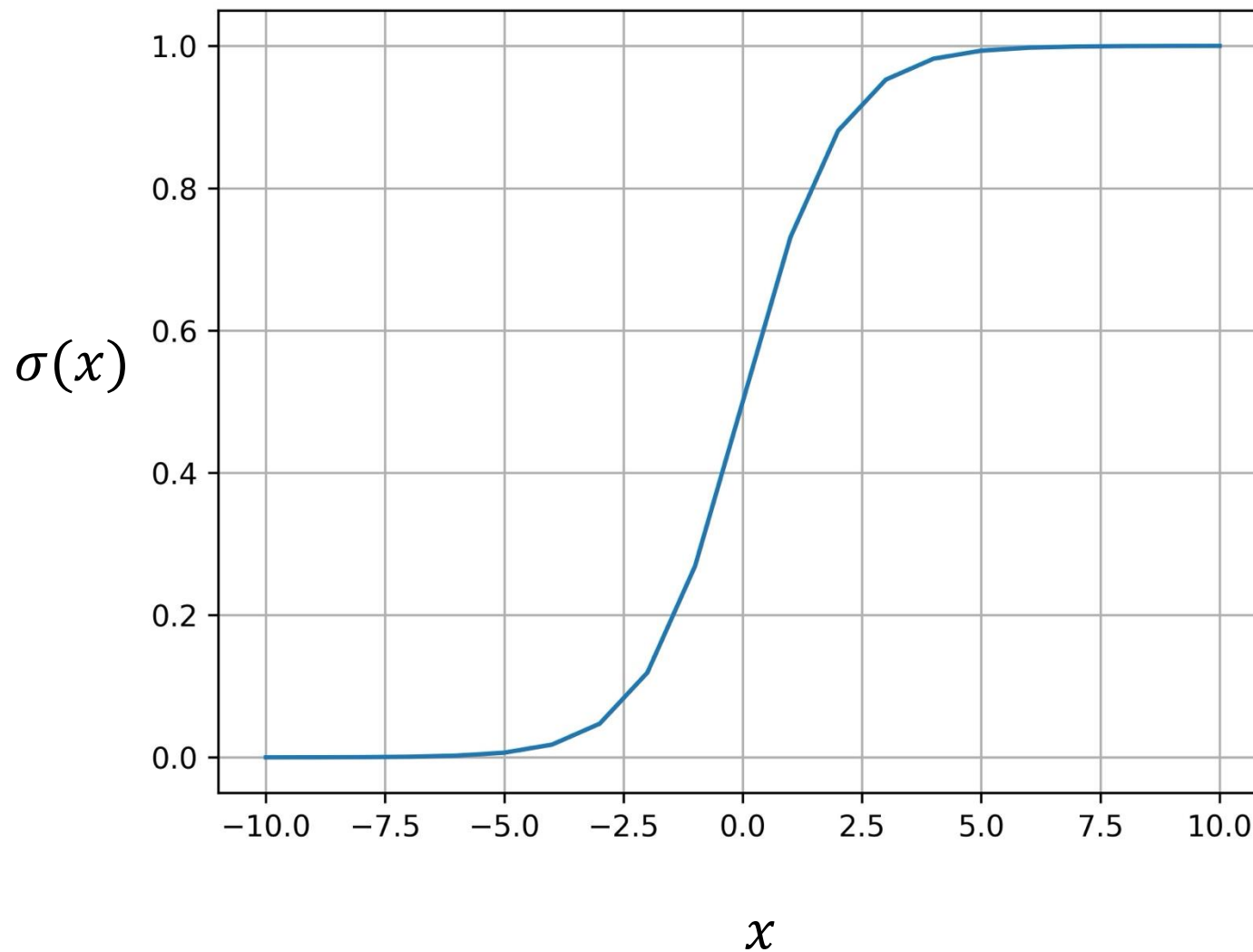
$$= f \left(\sum_{i=1}^n w_i^* h_i + b \right)$$

$f \equiv \text{activation function}$

$$= \left(\sum_{i=1}^n w_i^* h_i + b \right)$$

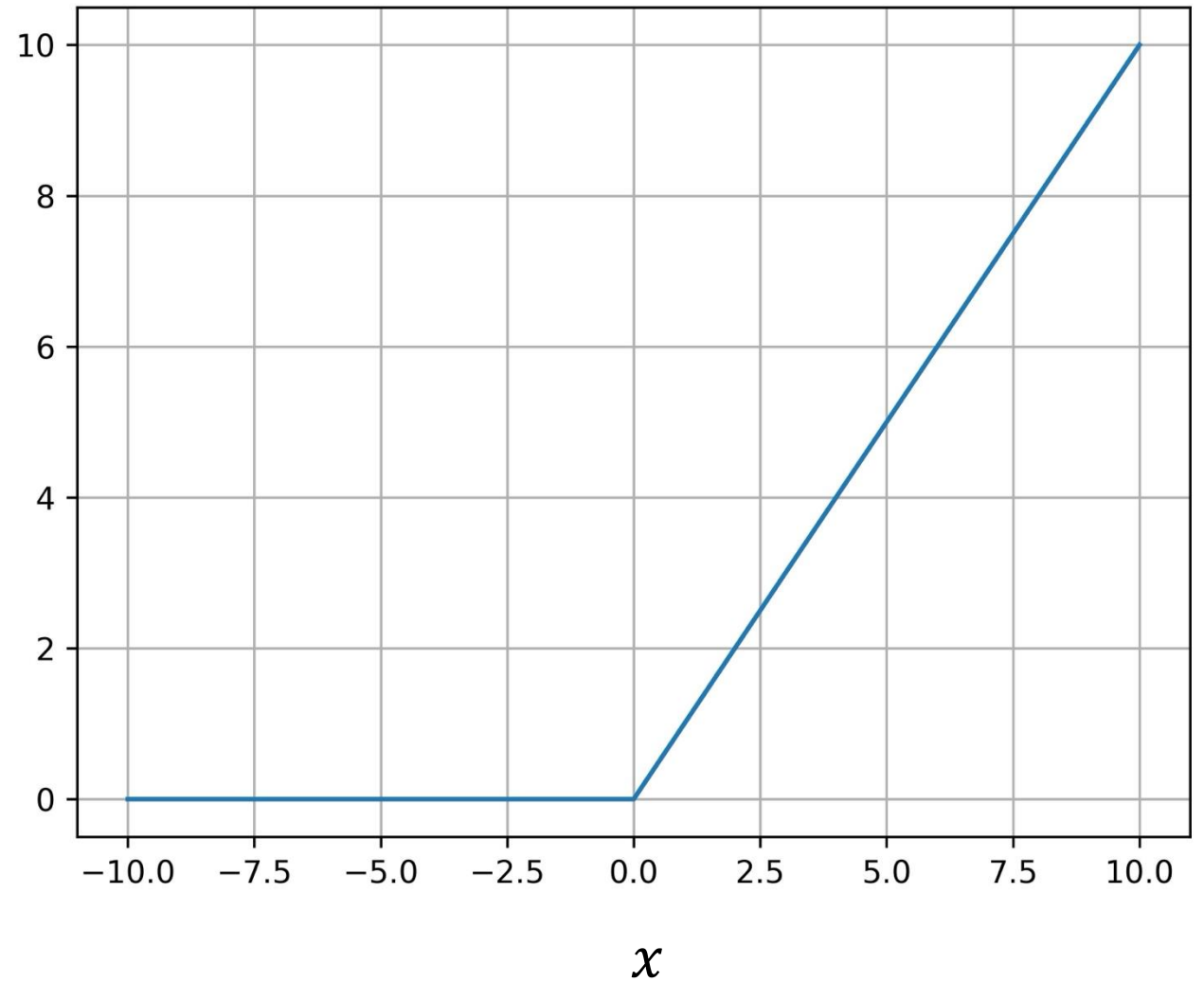
This Sum can be a very large number!

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



$$\text{ReLU}(x) = \max(0, x)$$

$\text{ReLU}(x)$



For an input vector of $x = [x_1, x_2, \dots, x_n]$

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{i=0}^n e^{x_i}}$$

