



DEEP LEARNING

Amrita Vishwa Vidyapeetham
Amritapuri Campus

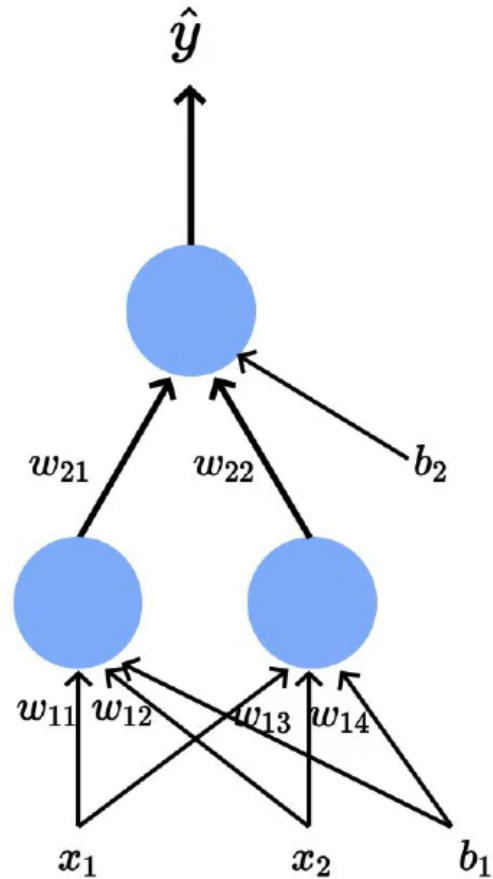




Feed Forward Neural Networks- Introduction

Citation Note: The content, of this presentation were inspired by the awesome lectures and the material offered by Prof. [Mitesh M. Khapra](#) on [NPTEL's Deep Learning](#) course

Multiple sigmoid Neurons



$$h_1 = f_1(x_1, x_2)$$

$$h_2 = f_2(x_1, x_2)$$

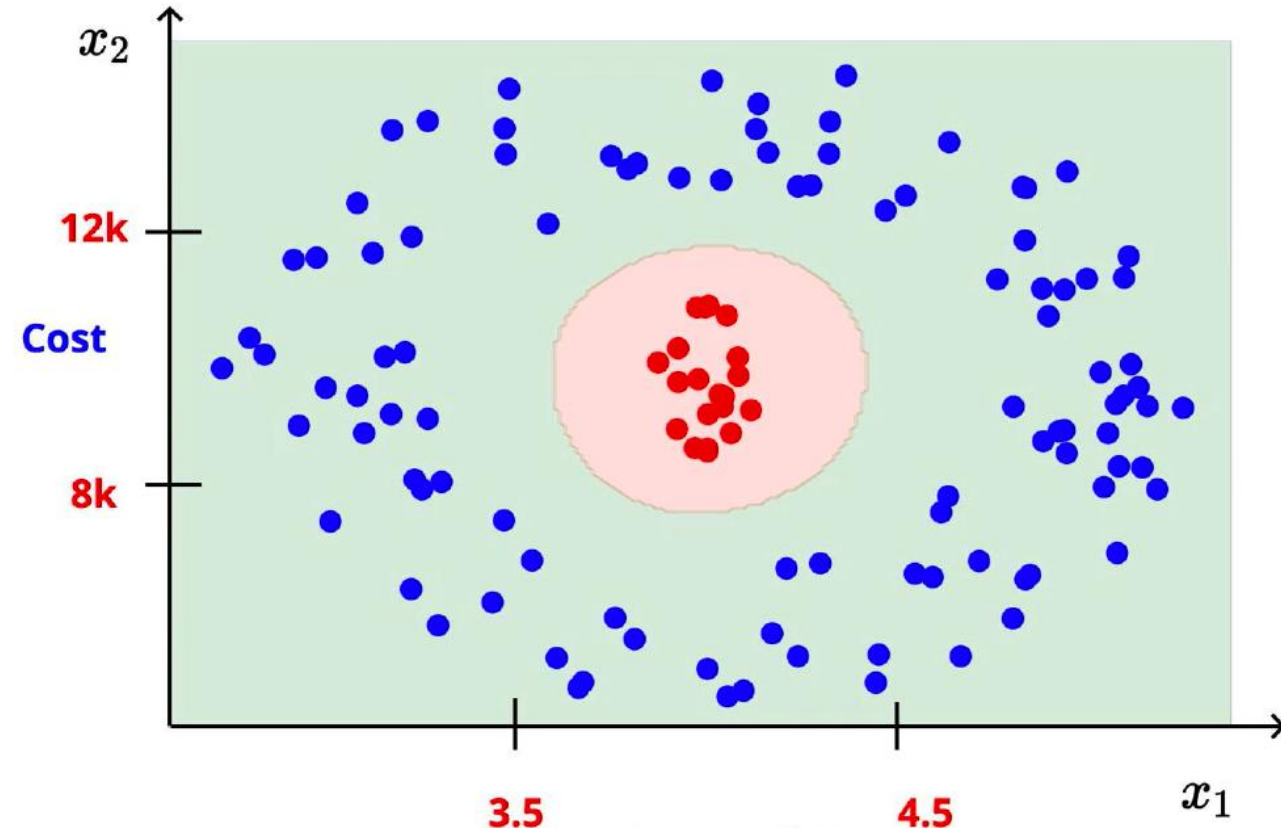
$$\hat{y} = g(h_1, h_2)$$

$$h_1 = \frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}}$$

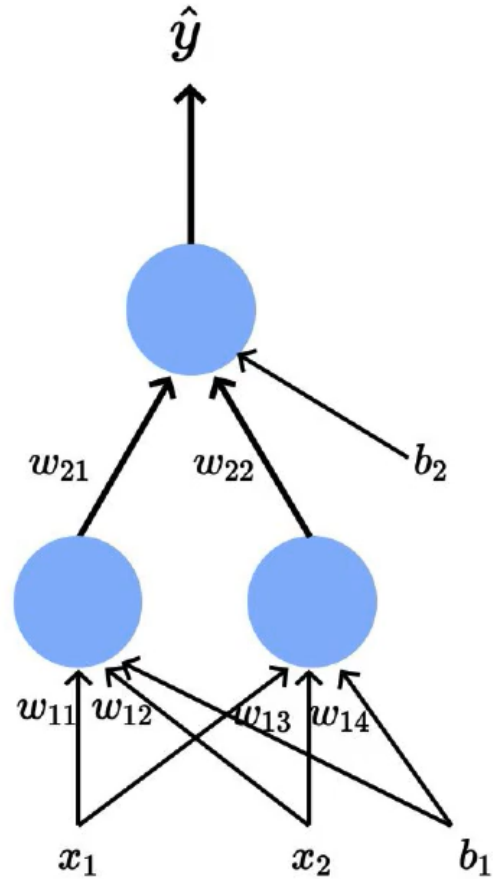
$$h_2 = \frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}}$$

$$\hat{y} = \frac{1}{1+e^{-(w_{21}*h_1+w_{22}*h_2+b_2)}}$$

$$= \frac{1}{1+e^{-\left(w_{21}*\left(\frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}}\right)+w_{22}*\left(\frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}}\right)+b_2\right)}}$$



Multiple sigmoid Neurons



$$h_1 = f_1(x_1, x_2)$$

$$h_2 = f_2(x_1, x_2)$$

$$\hat{y} = g(h_1, h_2)$$

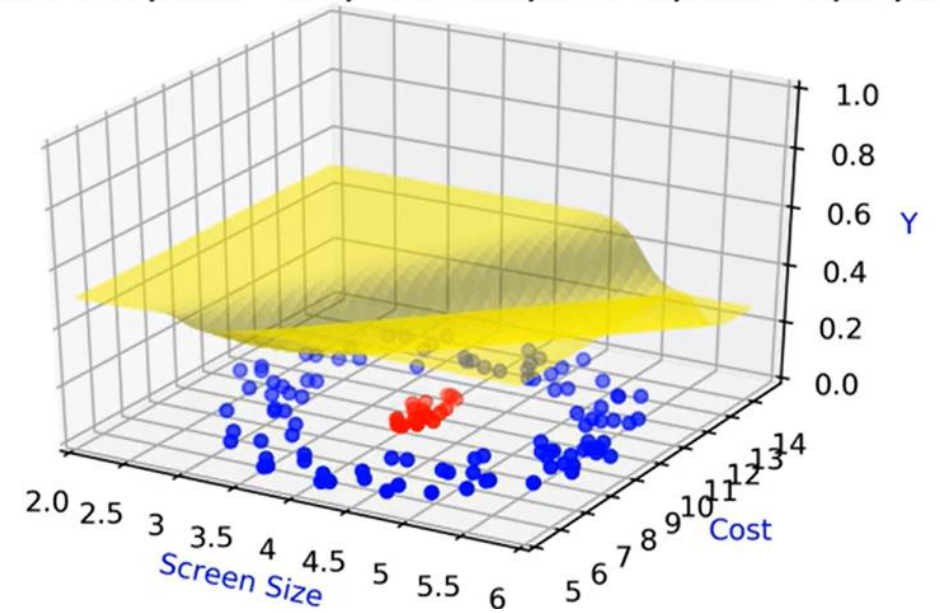
$$h_1 = \frac{1}{1 + e^{-(w_{11} * x_1 + w_{12} * x_2 + b_1)}}$$

$$h_2 = \frac{1}{1 + e^{-(w_{13} * x_1 + w_{14} * x_2 + b_1)}}$$

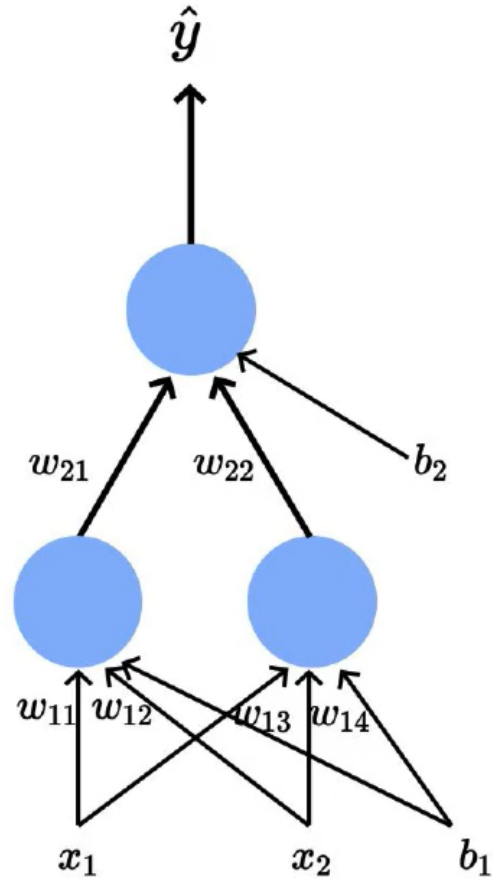
$$\hat{y} = \frac{1}{1 + e^{-(w_{21} * h_1 + w_{22} * h_2 + b_2)}}$$

$$= \frac{1}{1 + e^{-(w_{21} * (\frac{1}{1 + e^{-(w_{11} * x_1 + w_{12} * x_2 + b_1)}}) + w_{22} * (\frac{1}{1 + e^{-(w_{13} * x_1 + w_{14} * x_2 + b_1)}}) + b_2)}}$$

$$w_{11}=-2, w_{12}=1.0, w_{13}=-1.5, w_{14}=1.5, w_{21}=1, w_{22}=-1, b_1, b_2=0$$



Multiple sigmoid Neurons



$$h_1 = f_1(x_1, x_2)$$

$$h_2 = f_2(x_1, x_2)$$

$$\hat{y} = g(h_1, h_2)$$

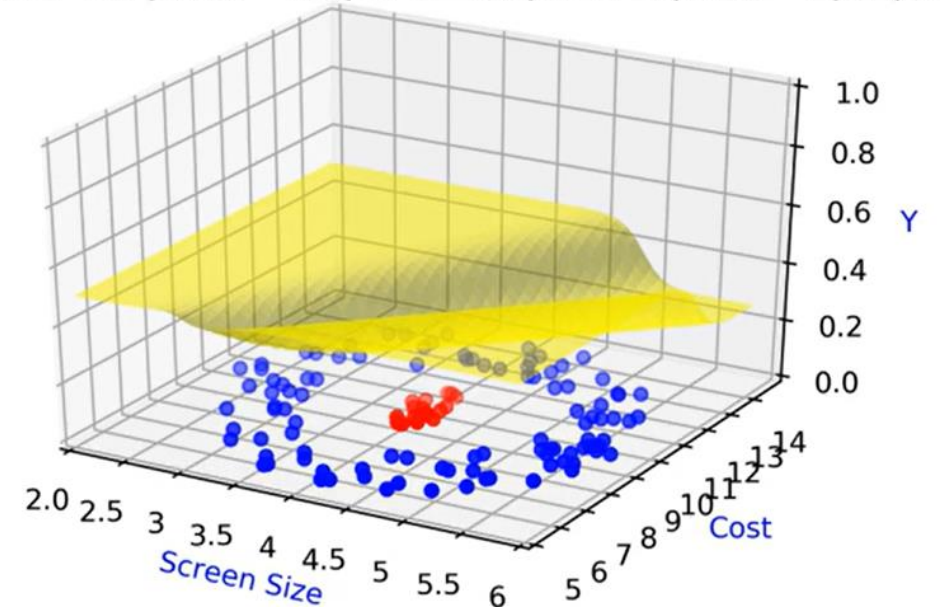
$$h_1 = \frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}}$$

$$h_2 = \frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}}$$

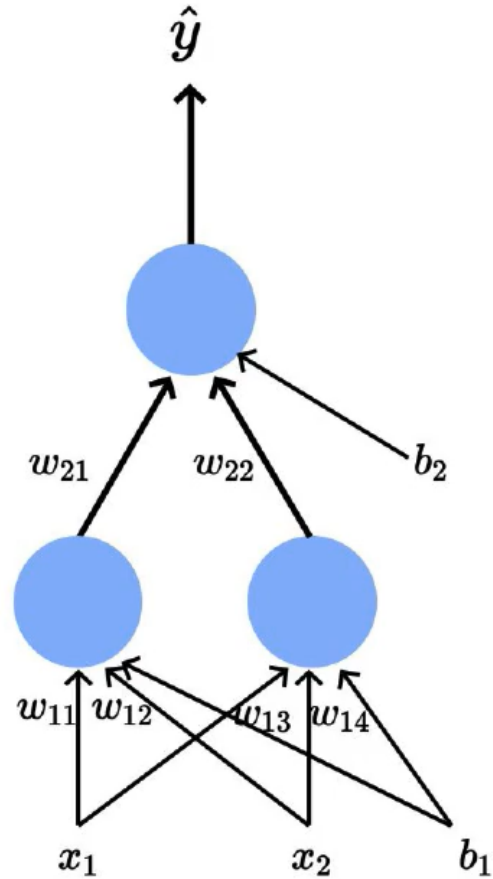
$$\hat{y} = \frac{1}{1+e^{-(w_{21}*h_1+w_{22}*h_2+b_2)}}$$

$$= \frac{1}{1+e^{-(w_{21}*(\frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}})+w_{22}*(\frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}})+b_2)}}$$

$$w_{11}=-2, w_{12}=1.0, w_{13}=-1.5, w_{14}=1.5, w_{21}=1, w_{22}=-1, b_1, b_2=0$$



Multiple sigmoid Neurons



$$h_1 = f_1(x_1, x_2)$$

$$h_2 = f_2(x_1, x_2)$$

$$\hat{y} = g(h_1, h_2)$$

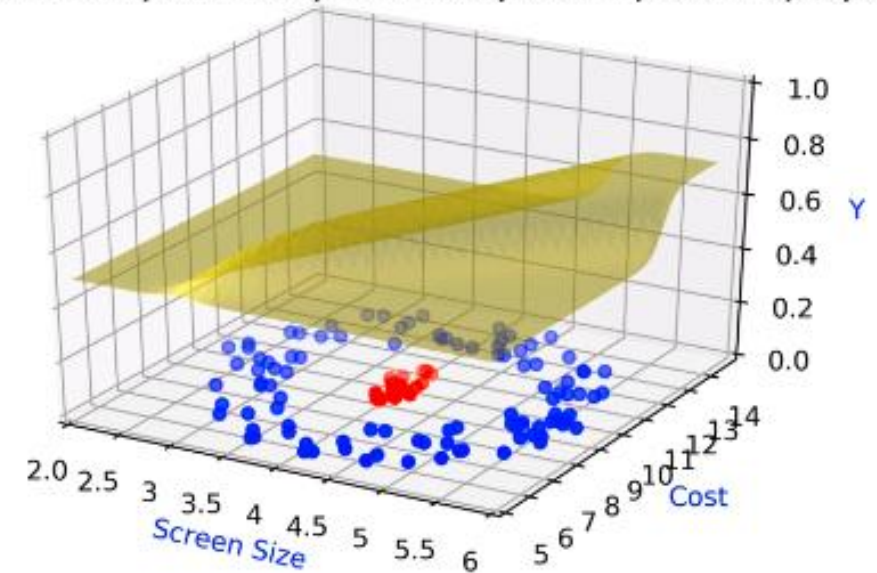
$$h_1 = \frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}}$$

$$h_2 = \frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}}$$

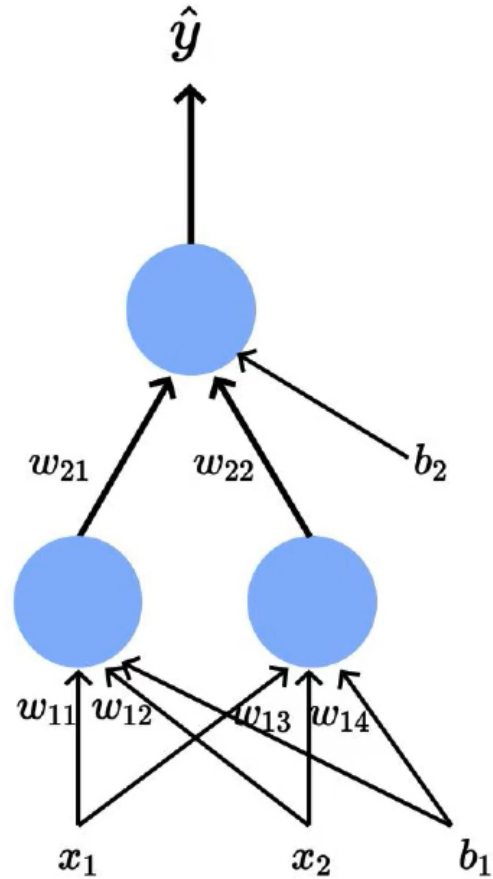
$$\hat{y} = \frac{1}{1+e^{-(w_{21}*h_1+w_{22}*h_2+b_2)}}$$

$$= \frac{1}{1+e^{-(w_{21}*(\frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}})+w_{22}*(\frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}})+b_2)}}$$

$$w_{11}=2, w_{12}=-1.0, w_{13}=2.0, w_{14}=-2.0, w_{21}=1, w_{22}=-1, b_1, b_2=0$$



Multiple sigmoid Neurons



$$h_1 = f_1(x_1, x_2)$$

$$h_2 = f_2(x_1, x_2)$$

$$\hat{y} = g(h_1, h_2)$$

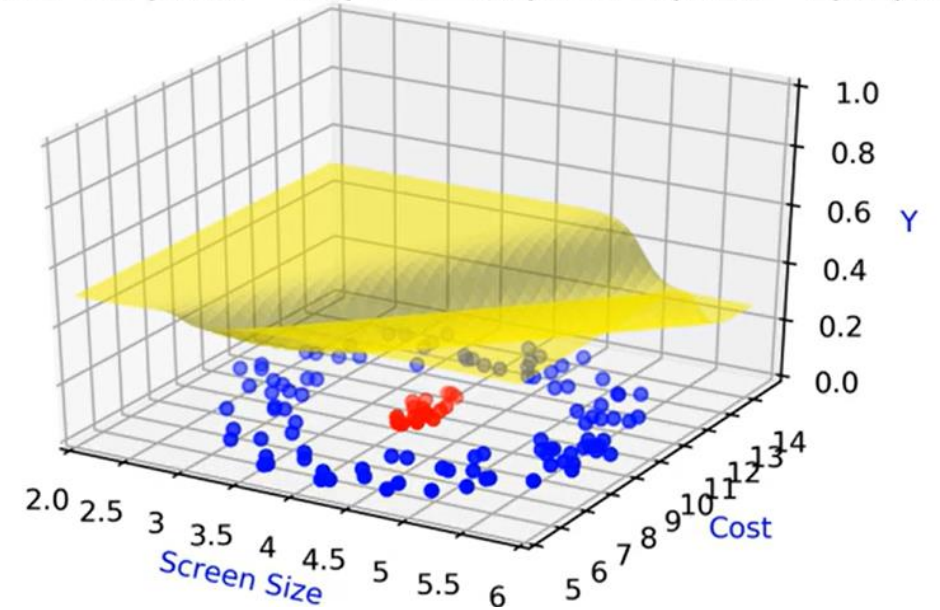
$$h_1 = \frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}}$$

$$h_2 = \frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}}$$

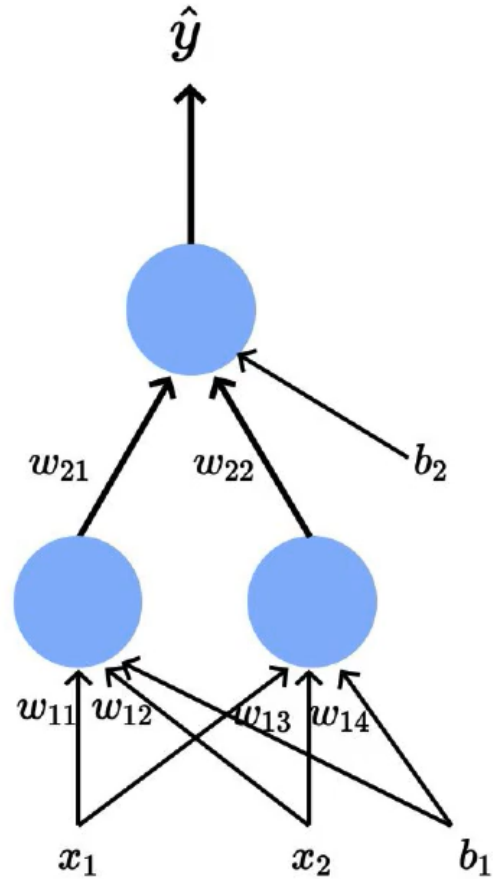
$$\hat{y} = \frac{1}{1+e^{-(w_{21}*h_1+w_{22}*h_2+b_2)}}$$

$$= \frac{1}{1+e^{-(w_{21}*(\frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}})+w_{22}*(\frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}})+b_2)}}$$

$$w_{11}=-2, w_{12}=1.0, w_{13}=-1.5, w_{14}=1.5, w_{21}=1, w_{22}=-1, b_1, b_2=0$$



Multiple sigmoid Neurons



$$h_1 = f_1(x_1, x_2)$$

$$h_2 = f_2(x_1, x_2)$$

$$\hat{y} = g(h_1, h_2)$$

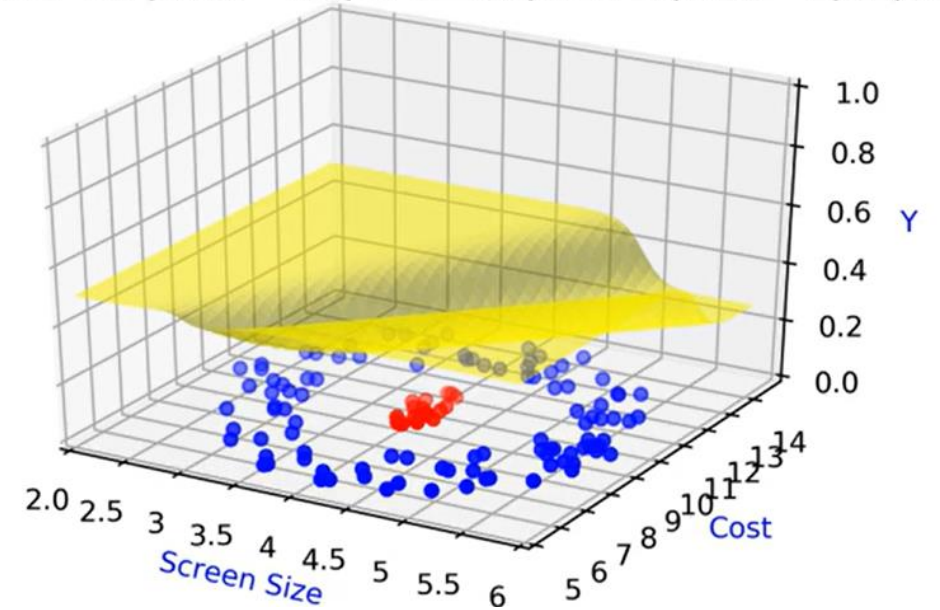
$$h_1 = \frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}}$$

$$h_2 = \frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}}$$

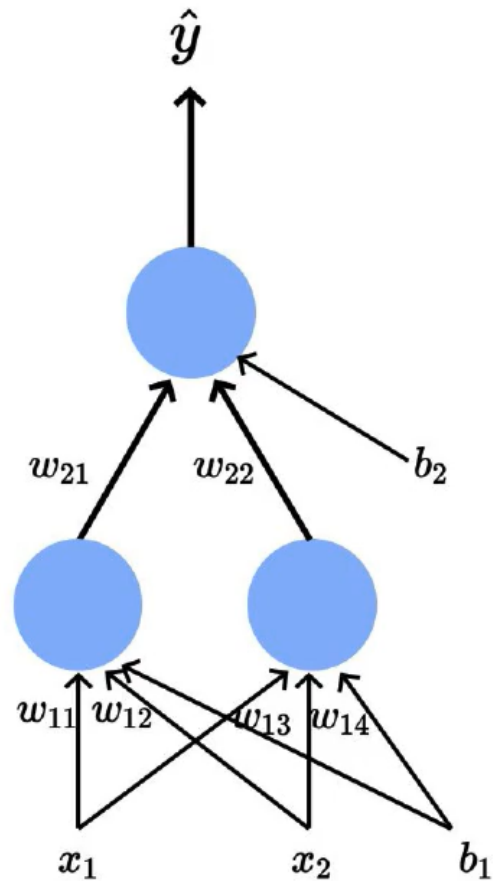
$$\hat{y} = \frac{1}{1+e^{-(w_{21}*h_1+w_{22}*h_2+b_2)}}$$

$$= \frac{1}{1+e^{-(w_{21}*(\frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}})+w_{22}*(\frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}})+b_2)}}$$

$$w_{11}=-2, w_{12}=1.0, w_{13}=-1.5, w_{14}=1.5, w_{21}=1, w_{22}=-1, b_1, b_2=0$$



Multi Layer Neural Networks



$$h_1 = f_1(x_1, x_2)$$

$$h_2 = f_2(x_1, x_2)$$

$$\hat{y} = g(h_1, h_2)$$

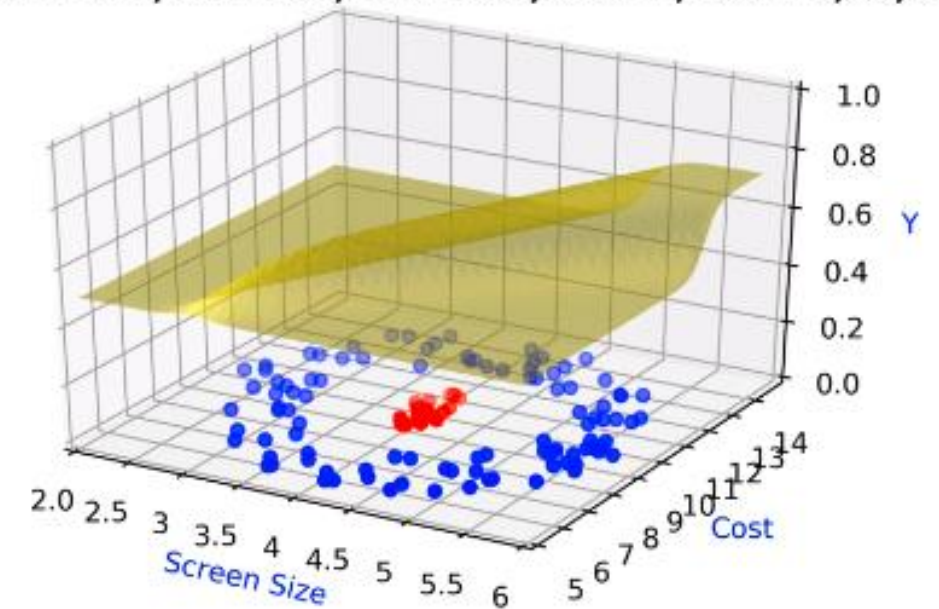
$$h_1 = \frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}}$$

$$h_2 = \frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}}$$

$$\hat{y} = \frac{1}{1+e^{-(w_{21}*h_1+w_{22}*h_2+b_2)}}$$

$$= \frac{1}{1+e^{-(w_{21}*(\frac{1}{1+e^{-(w_{11}*x_1+w_{12}*x_2+b_1)}})+w_{22}*(\frac{1}{1+e^{-(w_{13}*x_1+w_{14}*x_2+b_1)}})+b_2)}}$$

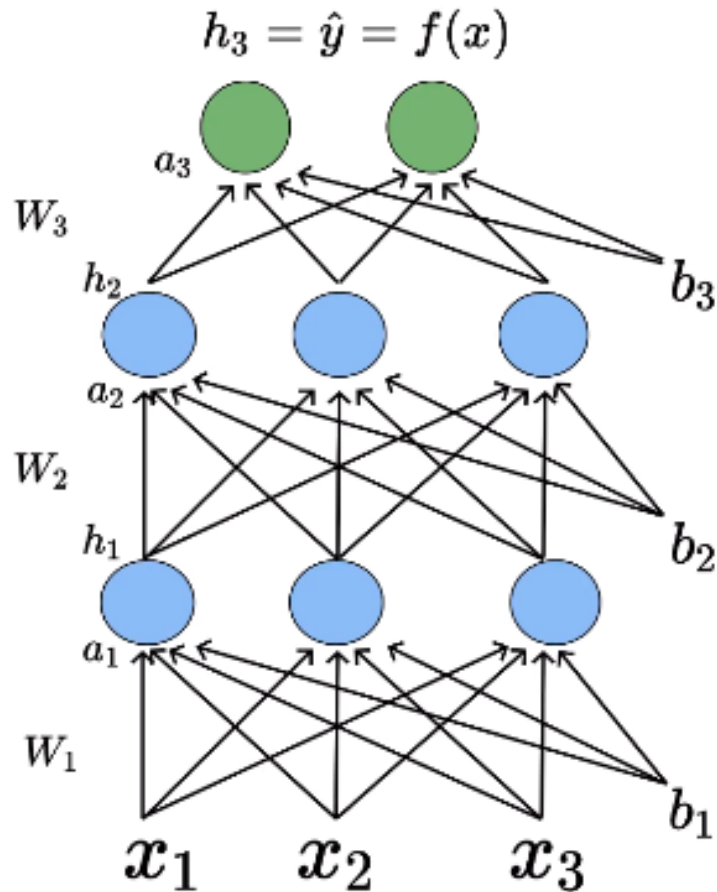
$$w_{11}=2, w_{12}=-1.0, w_{13}=2.0, w_{14}=-2.0, w_{21}=1, w_{22}=-1, b_1, b_2=0$$



Feed Forward Neural Networks

Understanding the computation

$W_{111} = W$ Layer no, Neuron in the next layer ,Input Neuron



$$W_1 = \begin{bmatrix} w_{111} & w_{112} & \dots & w_{1199} & w_{11100} \\ w_{121} & w_{122} & \dots & w_{1299} & w_{12100} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ w_{1101} & w_{1102} & \dots & w_{11099} & w_{110100} \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{100} \end{bmatrix}$$

$$a_{11} = w_{111} * x_1 + w_{112} * x_2 + w_{113} * x_3 + \dots + w_{11100} * x_{100} + b_{11}$$

$$a_{12} = w_{121} * x_1 + w_{122} * x_2 + w_{123} * x_3 + \dots + w_{12100} * x_{100} + b_{12}$$

$$\vdots$$

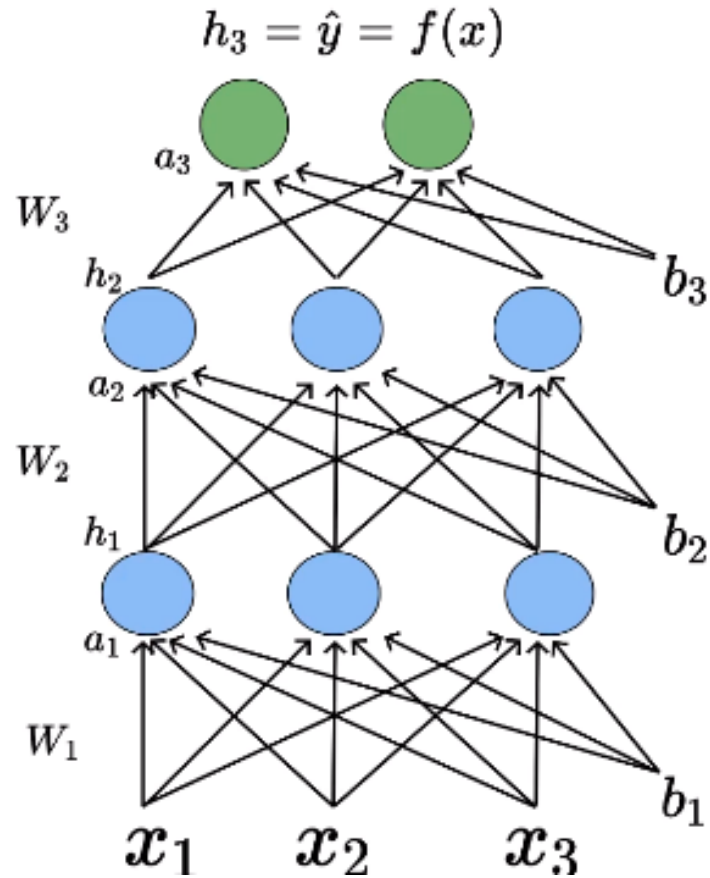
$$a_{110} = w_{1101} * x_1 + w_{1102} * x_2 + w_{1103} * x_3 + \dots + w_{110100} * x_{100} + b_{1,10}$$

$$a_1 = W_1 * x + b$$

(c) Or

$$h_{11} = g(a_{11}) \quad h_{12} = g(a_{12}) \quad \dots \quad h_{110} = g(a_{110})$$

Feed Forward Neural Networks



$$W_1 = \begin{bmatrix} w_{111} & w_{112} & \cdot & \cdot & \cdot & w_{1199} & w_{11100} \\ w_{121} & w_{122} & \cdot & \cdot & \cdot & w_{1299} & w_{12100} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ w_{1101} & w_{1102} & \cdot & \cdot & \cdot & w_{11099} & w_{110100} \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ x_{100} \end{bmatrix}$$

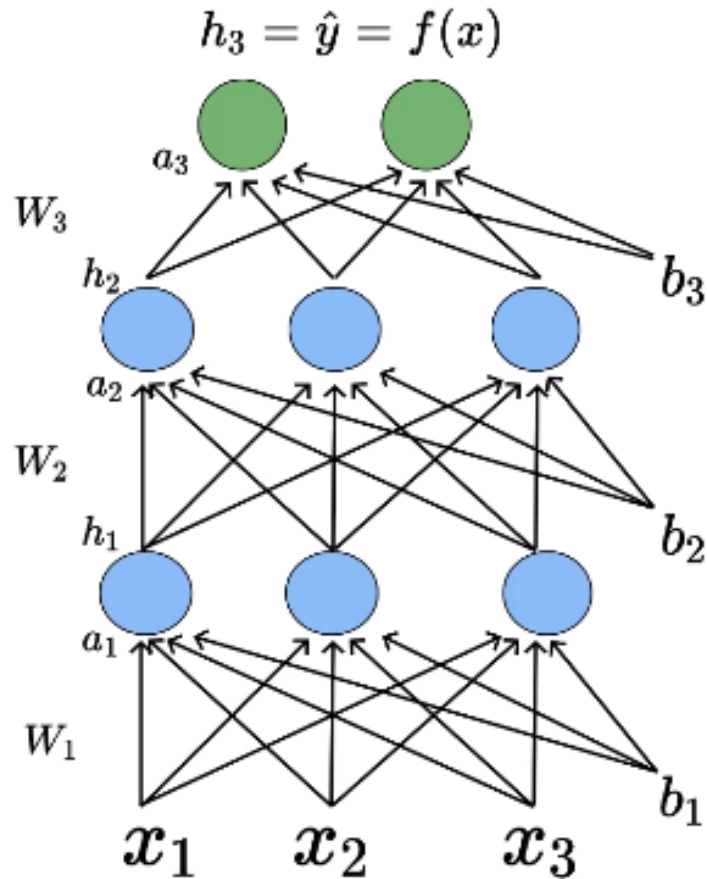
$$a_1 = W_1 * x + b$$

$$h_{11} = g(a_{11}) \quad h_{12} = g(a_{12}) \quad \cdot \quad \cdot \quad \cdot \quad \cdot \quad h_{110} = g(a_{110})$$

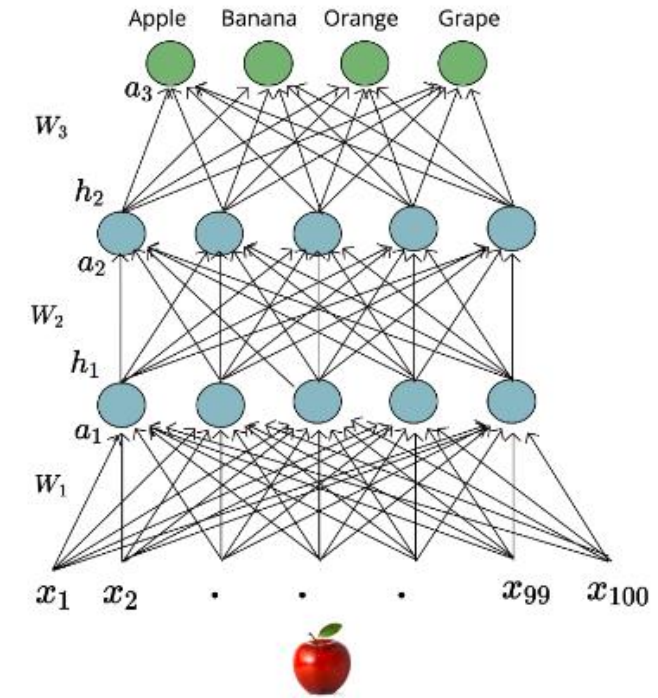
$$h_1 = g(a_1)$$

$$\hat{y} = f(x) = O(W_3 g(W_2 g(W_1 x + b_1) + b_2) + b_3)$$

Feed Forward Neural Networks

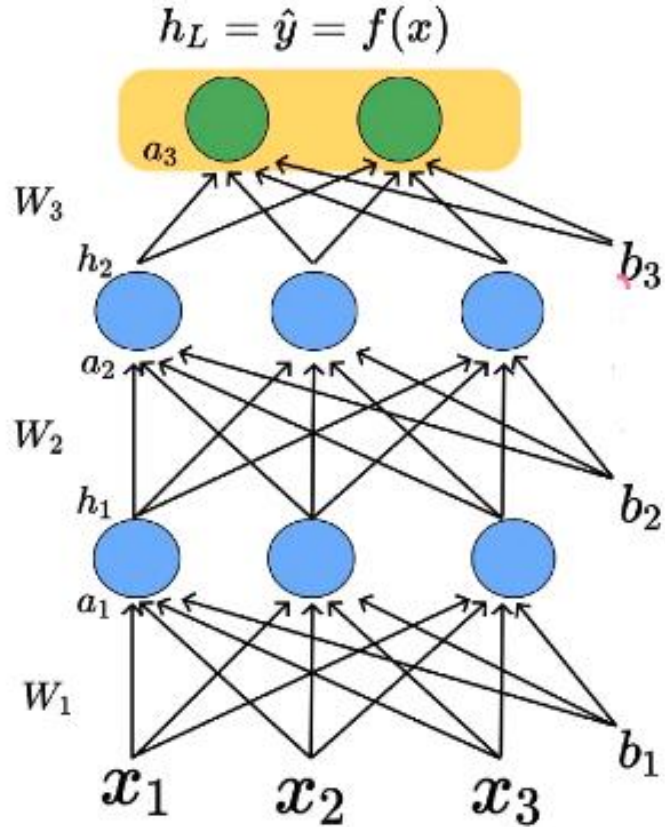


- The pre-activation at layer 'i' is given by
$$a_i(x) = W_i h_{i-1}(x) + b_i$$
- The activation at layer 'i' is given by
$$h_i(x) = g(a_i(x))$$
where 'g' is called as the activation function
- The activation at output layer 'L' is given by
$$f(x) = h_L = O(a_L)$$
where 'O' is called as the output activation function



AMRITA VISHWA VIDYAPEETHAM

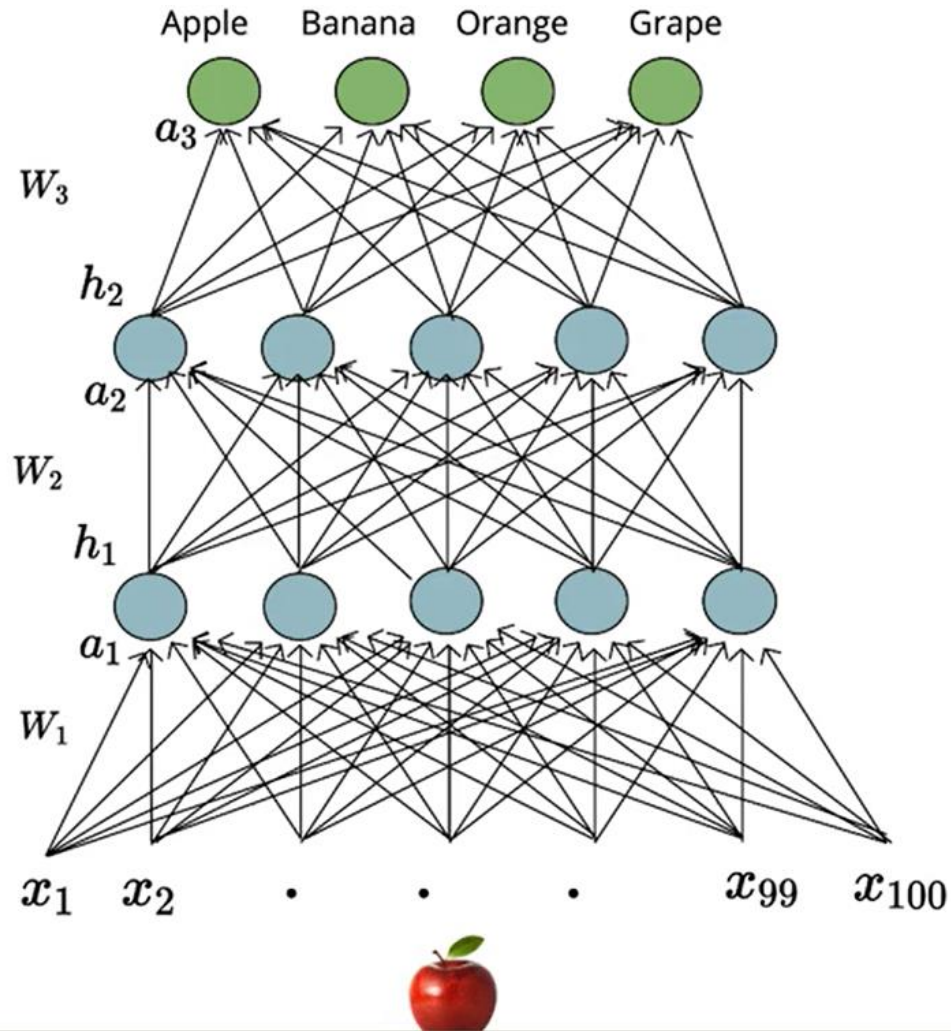
Feed Forward Neural Networks



$$\hat{y} = f(x) = O(W_3 g(W_2 g(W_1 x + b_1) + b_2) + b_3)$$

Output Activation function is chosen depending on the task at hand (can be a softmax, linear)

Is this ok?



$$\text{Say } a_3 = [3 \ 4 \ 10 \ 3]$$

Output Activation Function has to be chosen such that output is probability

$$\hat{y}_1 \Rightarrow \frac{3}{(3 + 4 + 10 + 3)} = 0.15$$

$$\hat{y}_2 = \frac{4}{(3 + 4 + 10 + 3)} = 0.20$$

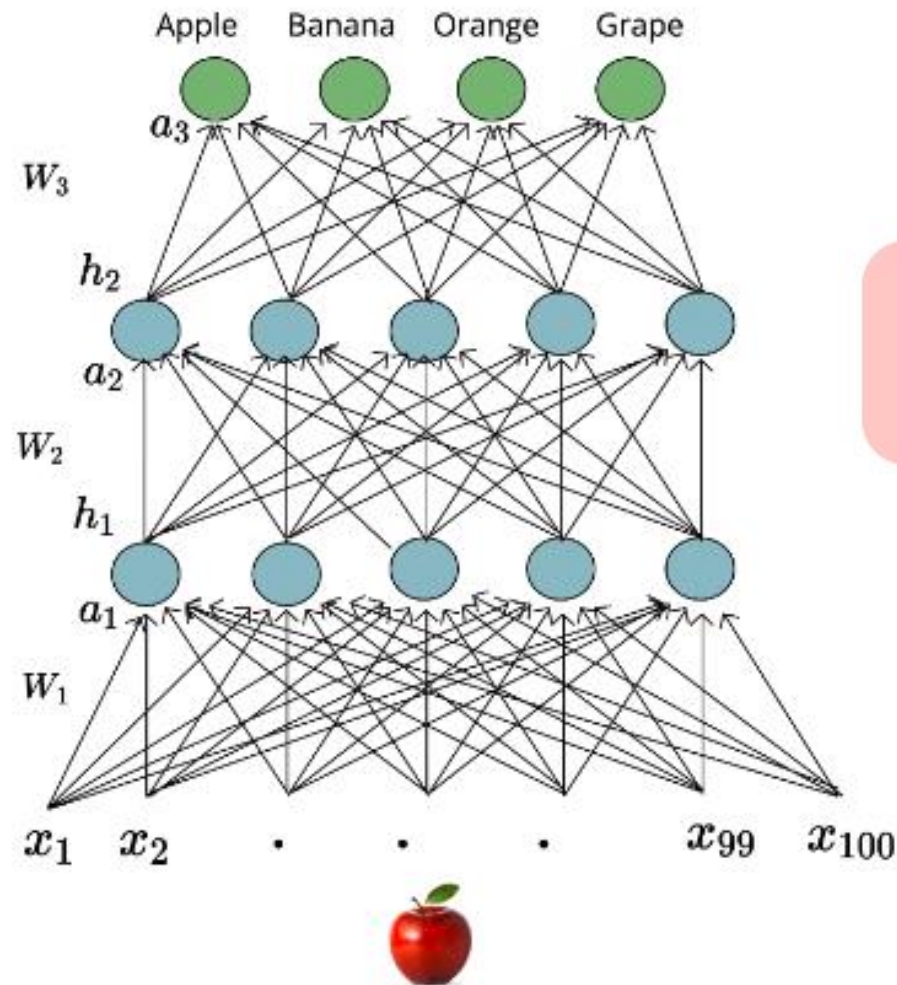
$$\hat{y}_3 = \frac{10}{(3 + 4 + 10 + 3)} = 0.50$$

$$\hat{y}_4 = \frac{3}{(3 + 4 + 10 + 3)} = 0.15$$

Take each entry and divide by the sum of all entries

Output Layer

Is this ok?-No



Say for other input $a_3 = [7 \ -2 \ 4 \ 1]$

Output Activation Function has to be chosen such that output is probability

$$\hat{y}_1 \Rightarrow \frac{7}{(7 + (-2) + 4 + 1)} = 0.70$$

$$\hat{y}_2 = \frac{-2}{(7 + (-2) + 4 + 1)} = -0.20 \quad \times$$

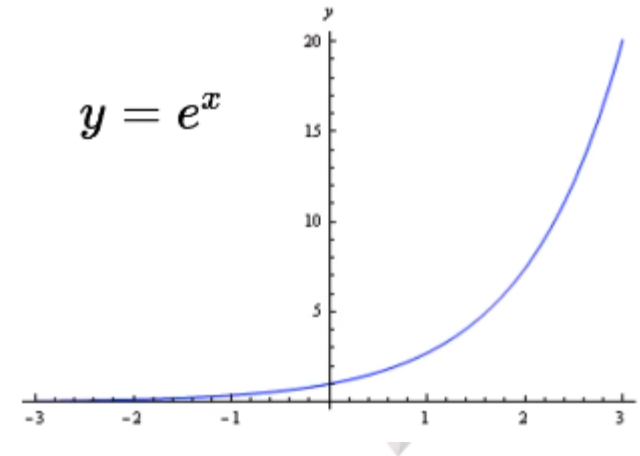
$$\hat{y}_3 = \frac{4}{(7 + (-2) + 4 + 1)} = 0.40$$

$$\hat{y}_4 = \frac{1}{(7 + (-2) + 4 + 1)} = 0.10$$

Softmax

Softmax is a kind of activation function with the speciality of output summing to 1.

$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \text{ for } i = 1, \dots, k$$



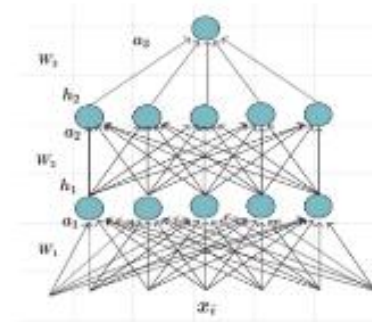
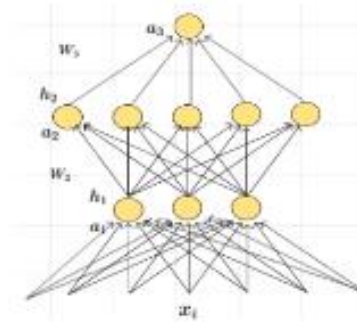
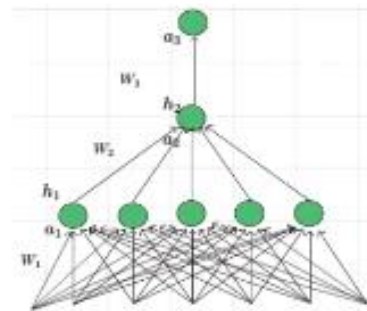
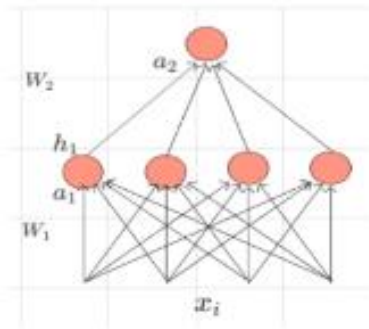
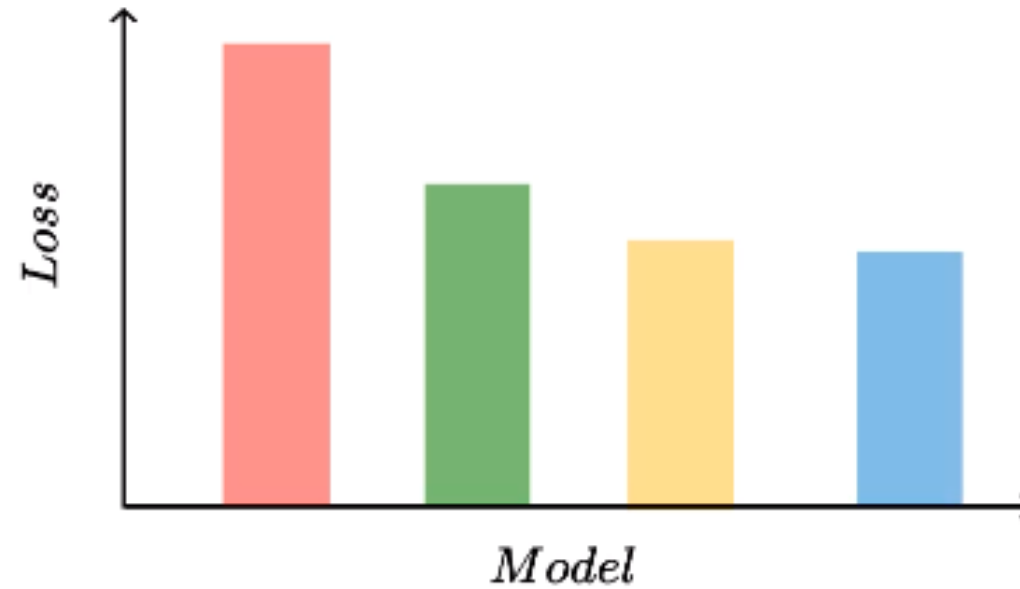
$$h = [h_1 \ h_2 \ h_3 \ h_4]$$

$$\text{softmax}(h) = [\text{softmax}(h_1) \ \text{softmax}(h_2) \ \text{softmax}(h_3) \ \text{softmax}(h_4)]$$

$$\text{softmax}(h) = \left[\frac{e^{h_1}}{\sum_{j=1}^4 e^{h_j}} \quad \frac{e^{h_2}}{\sum_{j=1}^4 e^{h_j}} \quad \frac{e^{h_3}}{\sum_{j=1}^4 e^{h_j}} \quad \frac{e^{h_4}}{\sum_{j=1}^4 e^{h_j}} \right]$$



Different network configurations



Learning Algorithm- backpropagation

Initialise w, b

Iterate over data:

compute \hat{y}

compute $\mathcal{L}(w, b)$

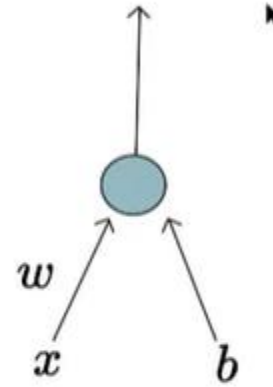
$$w_{111} = w_{111} - \eta \Delta w_{111}$$

$$w_{112} = w_{112} - \eta \Delta w_{112}$$

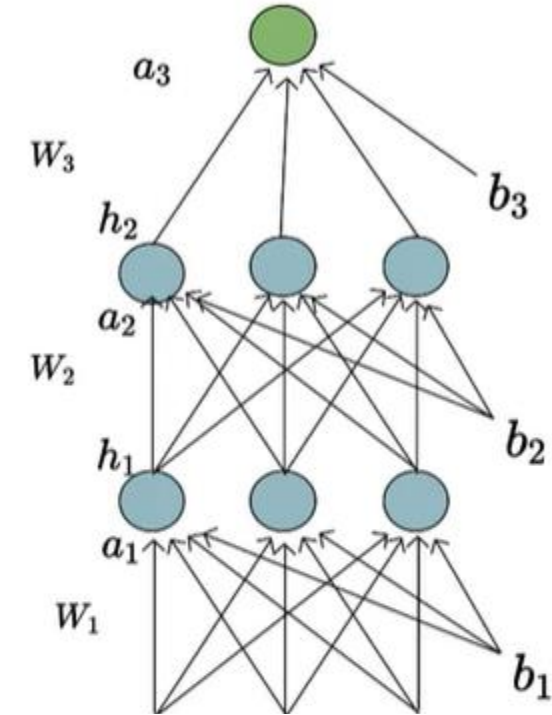
....

$$w_{313} = w_{313} - \eta \Delta w_{313}$$

till satisfied



$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$



$$\mathcal{L} = \frac{1}{5} \sum_{i=1}^{i=5} (f(x_i) - y_i)^2$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial}{\partial w} \left[\frac{1}{5} \sum_{i=1}^{i=5} (f(x_i) - y_i) \right]^2$$

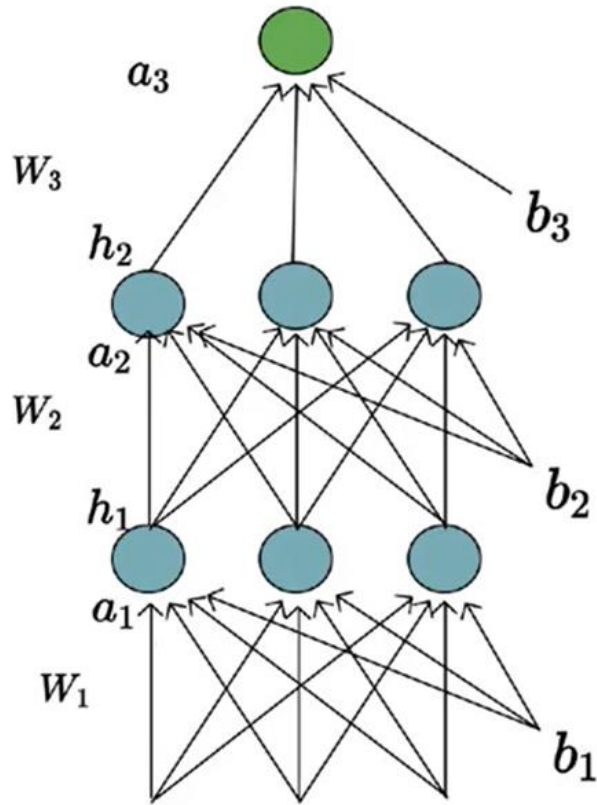
$$\Delta w = \frac{\partial \mathcal{L}}{\partial w} = \frac{1}{5} \sum_{i=1}^{i=5} \frac{\partial}{\partial w} (f(x_i) - y_i)^2$$

Calculation of Δw

$$\begin{aligned}\nabla w &= \frac{\partial}{\partial w} \left[\frac{1}{2} * (f(x) - y)^2 \right] \\&= \frac{1}{2} * [2 * (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y)] \\&= (f(x) - y) * \frac{\partial}{\partial w} (f(x)) \\&= (f(x) - y) * \frac{\partial}{\partial w} \left(\frac{1}{1+e^{-(wx+b)}} \right) \\&= (f(x) - y) * f(x) * (1 - f(x)) * x\end{aligned}$$

$$\begin{aligned}&\frac{\partial}{\partial w} \left(\frac{1}{1+e^{-(wx+b)}} \right) \\&= \frac{-1}{(1+e^{-(wx+b)})^2} \frac{\partial}{\partial w} (e^{-(wx+b)}) \\&= \frac{-1}{(1+e^{-(wx+b)})^2} * (e^{-(wx+b)}) \frac{\partial}{\partial w} (-(wx+b)) \\&= \frac{-1}{(1+e^{-(wx+b)})} * \frac{e^{-(wx+b)}}{(1+e^{-(wx+b)})} * (-x) \\&= \frac{1}{(1+e^{-(wx+b)})} * \frac{e^{-(wx+b)}}{(1+e^{-(wx+b)})} * (x) \\&= f(x) * (1 - f(x)) * x\end{aligned}$$

Partial derivative , Gradient



$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$

$$\begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial W_{111}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{11n}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{211}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{21n}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,11}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,1k}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,1k}} & \frac{\partial \mathcal{L}(\theta)}{\partial b_{11}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial b_{L1}} \\ \frac{\partial \mathcal{L}(\theta)}{\partial W_{121}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{12n}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{221}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{22n}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,21}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,2k}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,2k}} & \frac{\partial \mathcal{L}(\theta)}{\partial b_{12}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial b_{L2}} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial W_{1n1}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{1nn}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{2n1}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{2nn}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,n1}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,nk}} & \frac{\partial \mathcal{L}(\theta)}{\partial W_{L,nk}} & \frac{\partial \mathcal{L}(\theta)}{\partial b_{1n}} & \cdots & \frac{\partial \mathcal{L}(\theta)}{\partial b_{Lk}} \end{bmatrix}$$

The partial derivative notation is used to specify the derivative of a function of more than one variable with respect to one of its variables.

The gradient of a function f , denoted as ∇f , is the collection of all its partial derivatives into a vector.

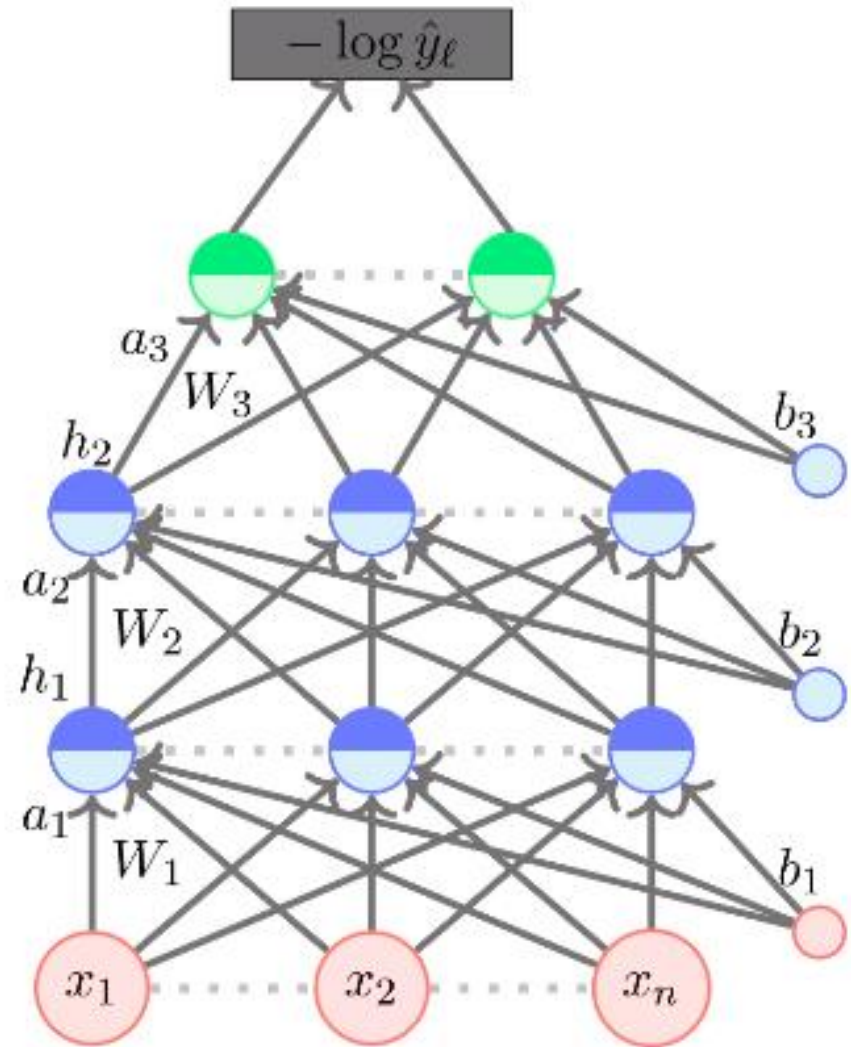
Backpropagation

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$

The partial derivative notation is used to **specify the derivative of a function of more than one variable with respect to one of its variables.**

The gradient of a function f , denoted as ∇f , is **the collection of all its partial derivatives into a vector.**



Calculus basics - Chain rule

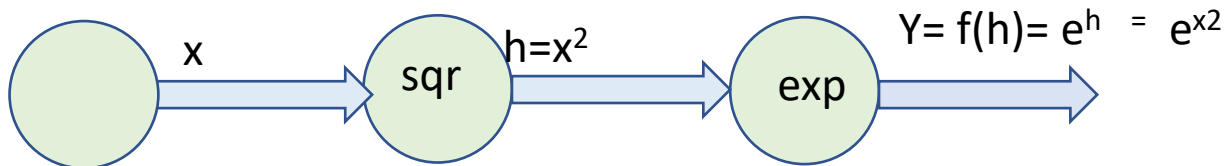
The chain rule is a method for finding the derivative of composite functions, or functions that are made by combining one or more functions.

$$\frac{de^x}{dx} = e^x$$

$$\frac{dx^2}{dx} = 2x$$

$$\frac{d(1/x)}{dx} = -\frac{1}{x^2}$$

$$\frac{de^{x^2}}{dx} = \frac{de^{x^2}}{dx^2} \cdot \frac{dx^2}{dx} = \frac{de^z}{dz} \cdot \frac{dx^2}{dx} = (e^z) \cdot (2x) = (e^{x^2}) \cdot (2x) = 2xe^{x^2}$$



$$h=f(x)$$

$$Y=f(h)= e^h$$

$$\frac{dy}{dx} = \frac{dy}{dh} \frac{dh}{dx} = \frac{de^h}{dh} \frac{dx^2}{dx} = e^h 2x = 2x e^{x^2}$$

Chain rule

$$\frac{de^{e^{x^2}}}{dx} = \frac{de^{e^{x^2}}}{de^{x^2}} \cdot \frac{de^{x^2}}{dx^2} \cdot \frac{dx^2}{dx}$$

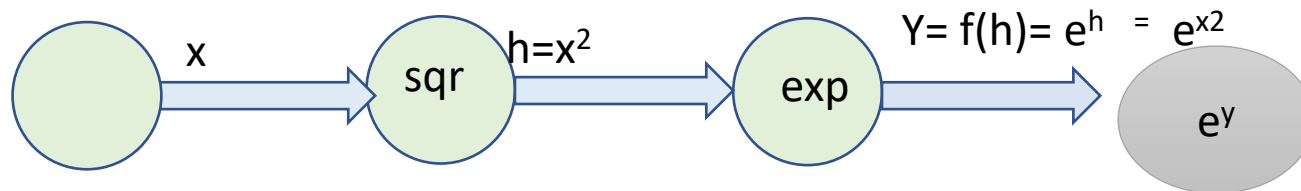
$$\frac{de^x}{dx} = e^x$$

$$\frac{dx^2}{dx} = 2x$$

$$\frac{d(1/x)}{dx} = -\frac{1}{x^2}$$

$$\frac{de^{x^2}}{dx} = \frac{de^{x^2}}{dx^2} \cdot \frac{dx^2}{dx} = \frac{de^z}{dz} \cdot \frac{dx^2}{dx} = (e^z) \cdot (2x) = (e^{x^2}) \cdot (2x) = 2xe^{x^2}$$

$$\frac{de^{e^{x^2}}}{dx} = \frac{de^{e^{x^2}}}{de^{x^2}} \cdot \frac{de^{x^2}}{dx} = \frac{de^z}{dz} \cdot \frac{de^{x^2}}{dx} = (e^z) \cdot (2xe^{x^2}) = (e^{e^{x^2}}) \cdot (2xe^{x^2}) = 2xe^{x^2} e^{e^{x^2}}$$



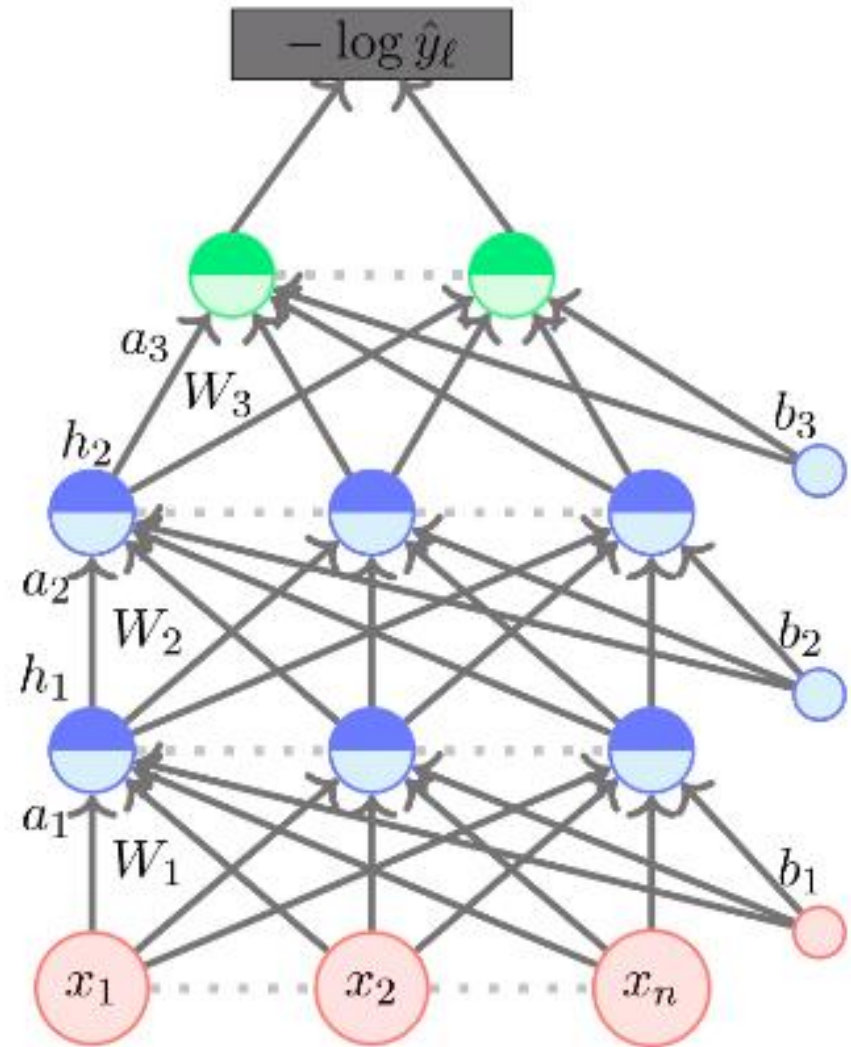
Backpropagation

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$

The partial derivative notation is used to **specify the derivative of a function of more than one variable with respect to one of its variables.**

The gradient of a function f , denoted as ∇f , is **the collection of all its partial derivatives into a vector.**



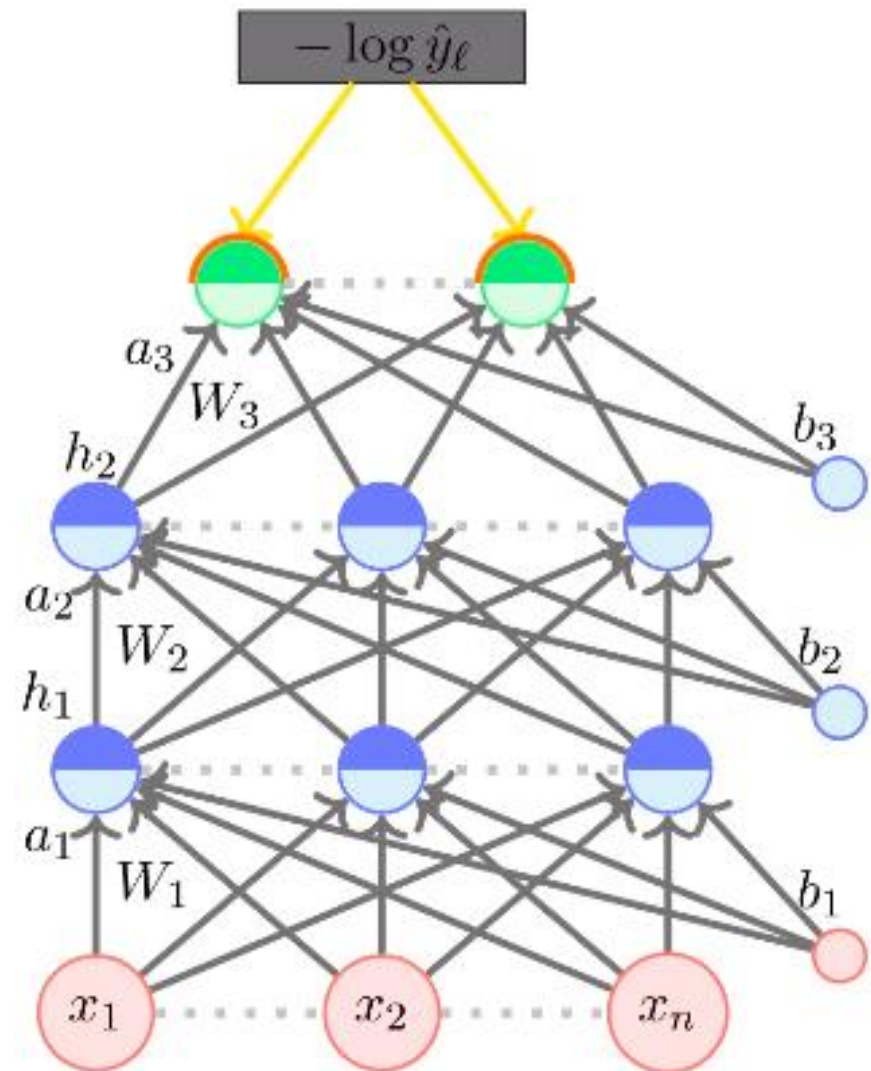
Backpropagation

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$

The partial derivative notation is used to **specify the derivative of a function of more than one variable with respect to one of its variables.**

The gradient of a function f , denoted as ∇f , is **the collection of all its partial derivatives into a vector.**



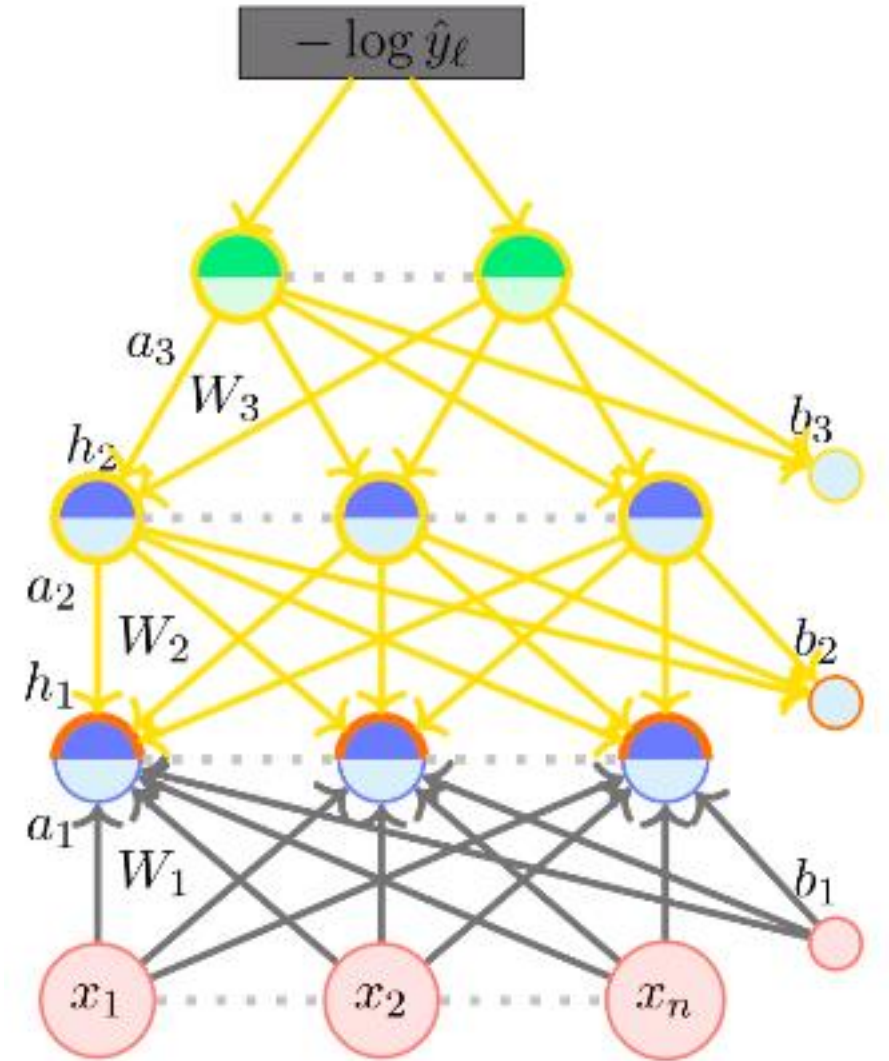
Backpropagation

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$

The partial derivative notation is used to **specify the derivative of a function of more than one variable with respect to one of its variables.**

The gradient of a function f , denoted as ∇f , is **the collection of all its partial derivatives into a vector.**



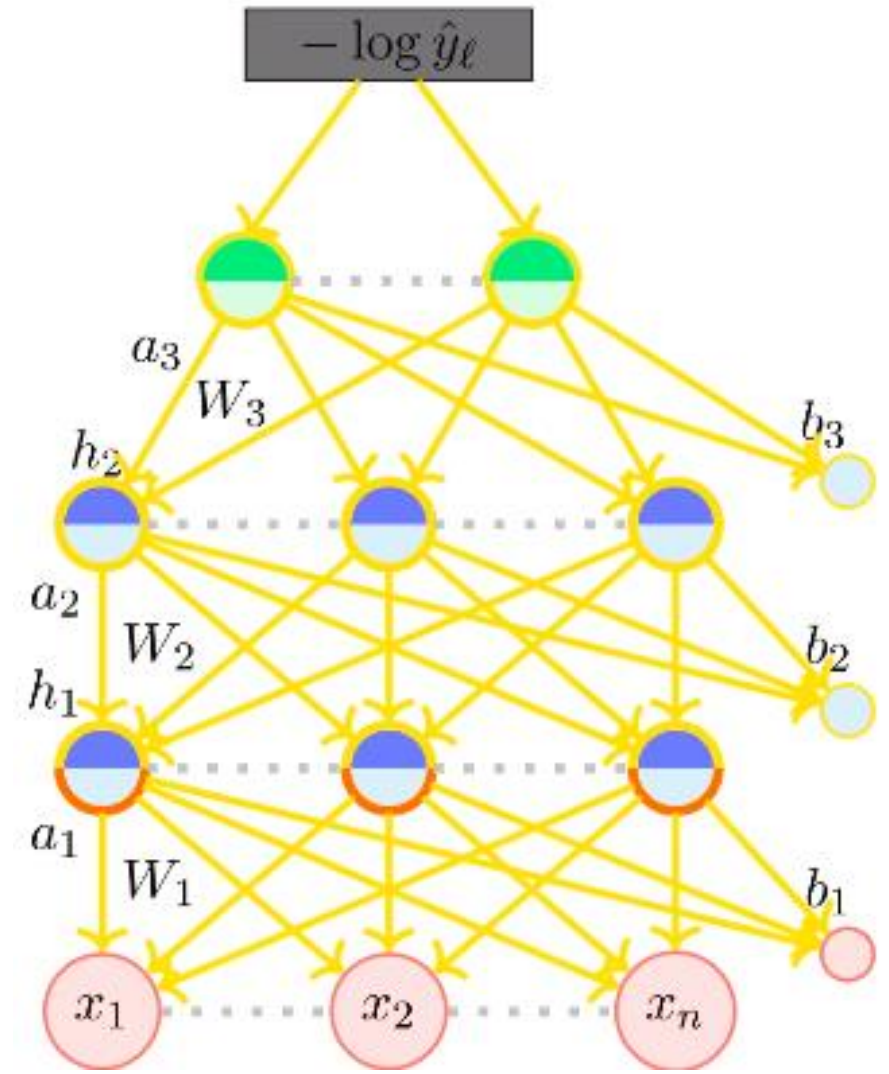
Backpropagation

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

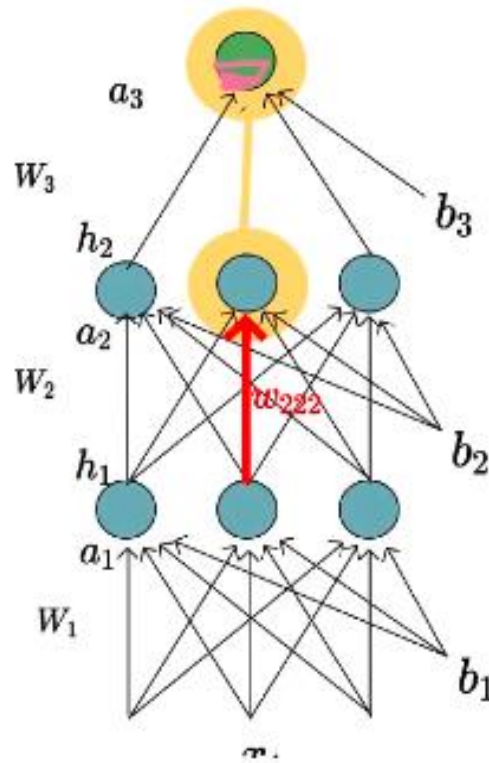
$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$

The partial derivative notation is used to **specify the derivative of a function of more than one variable with respect to one of its variables.**

The gradient of a function f , denoted as ∇f , is **the collection of all its partial derivatives into a vector.**



Learning algorithm- Back propagation

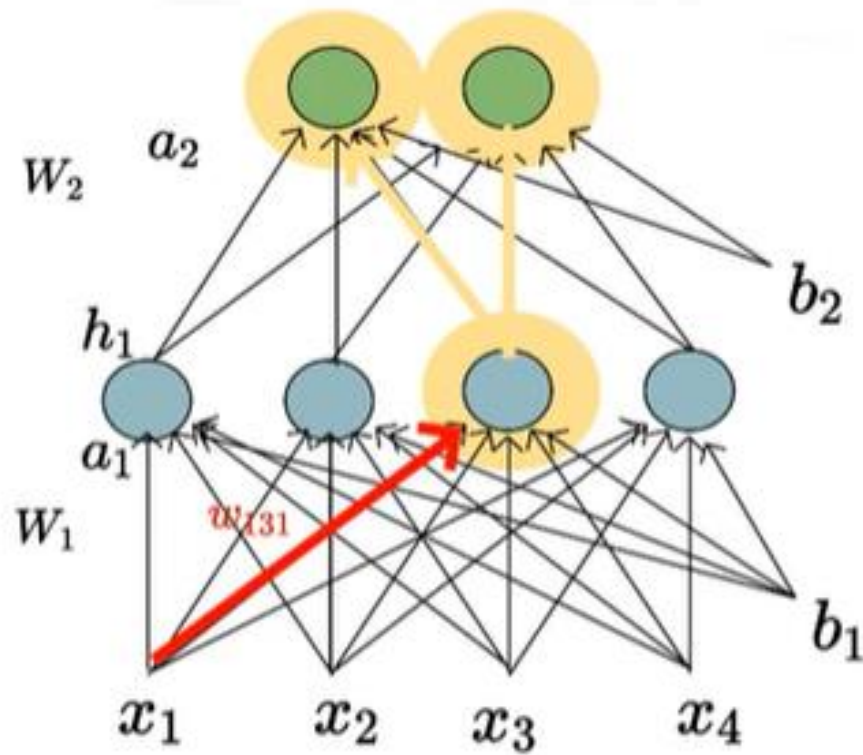


- Let us focus on the highlighted weight (w_{222})
- To learn this weight, we have to compute partial derivative w.r.t loss function

$$(w_{222})_{t+1} = (w_{222})_t - \eta * \left(\frac{\partial L}{\partial w_{222}} \right)$$

$$\begin{aligned} \frac{\partial L}{\partial w_{222}} &= \left(\frac{\partial L}{\partial a_{22}} \right) \cdot \left(\frac{\partial a_{22}}{\partial w_{222}} \right) \\ &= \left(\frac{\partial L}{\partial h_{22}} \right) \cdot \left(\frac{\partial h_{22}}{\partial a_{22}} \right) \cdot \left(\frac{\partial a_{22}}{\partial w_{222}} \right) \\ &= \left(\frac{\partial L}{\partial a_{31}} \right) \cdot \left(\frac{\partial a_{31}}{\partial h_{22}} \right) \cdot \left(\frac{\partial h_{22}}{\partial a_{22}} \right) \cdot \left(\frac{\partial a_{22}}{\partial w_{222}} \right) \\ &= \left(\frac{\partial L}{\partial \hat{y}} \right) \cdot \left(\frac{\partial \hat{y}}{\partial a_{31}} \right) \cdot \left(\frac{\partial a_{31}}{\partial h_{22}} \right) \cdot \left(\frac{\partial h_{22}}{\partial a_{22}} \right) \cdot \left(\frac{\partial a_{22}}{\partial w_{222}} \right) \end{aligned}$$

Multiple paths



- There are 2 different paths connecting w_{131} with loss function
- Consider all those paths through which gradient is flowing back
- Sum up the gradients along all these paths(2 here)
- ie Apply independent chain rules to all those multiple paths and sum up the derivatives across all these paths and get the total derivative of the loss function w.r.t w_{131}

Learning Algorithm- backpropagation

Initialise w, b

Iterate over data:

compute \hat{y}

compute $\mathcal{L}(w, b)$

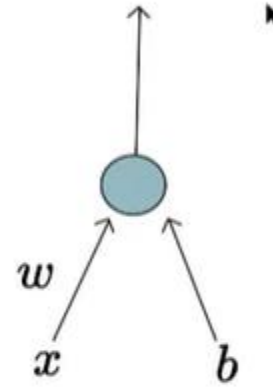
$$w_{111} = w_{111} - \eta \Delta w_{111}$$

$$w_{112} = w_{112} - \eta \Delta w_{112}$$

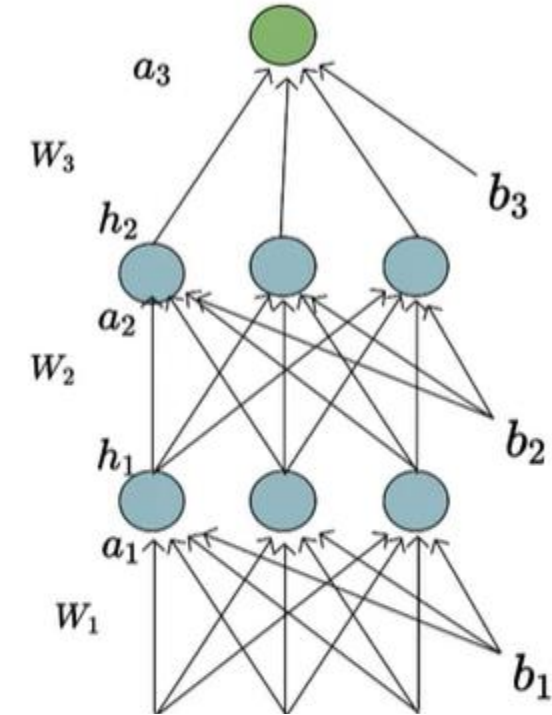
....

$$w_{313} = w_{313} - \eta \Delta w_{313}$$

till satisfied



$$\Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w}$$

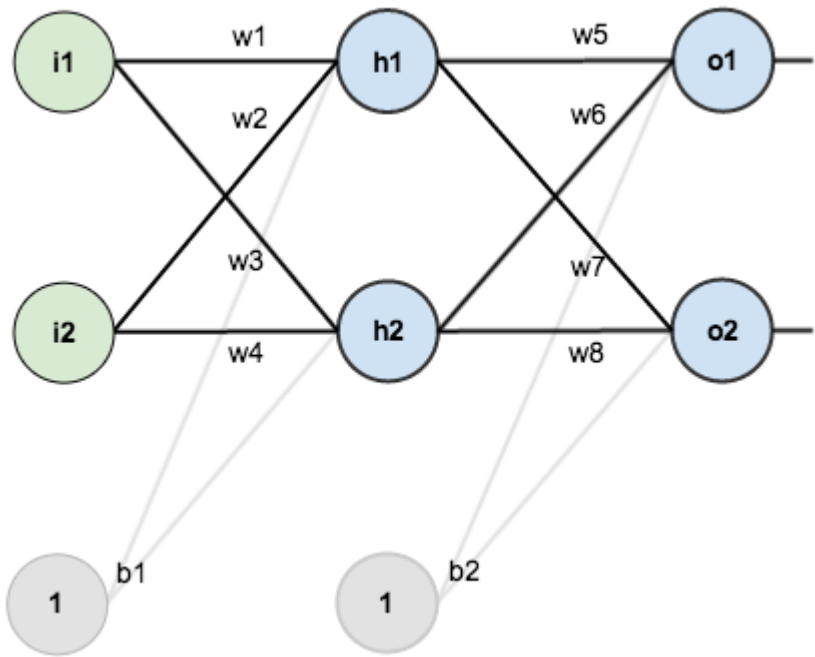


$$\mathcal{L} = \frac{1}{5} \sum_{i=1}^{i=5} (f(x_i) - y_i)^2$$

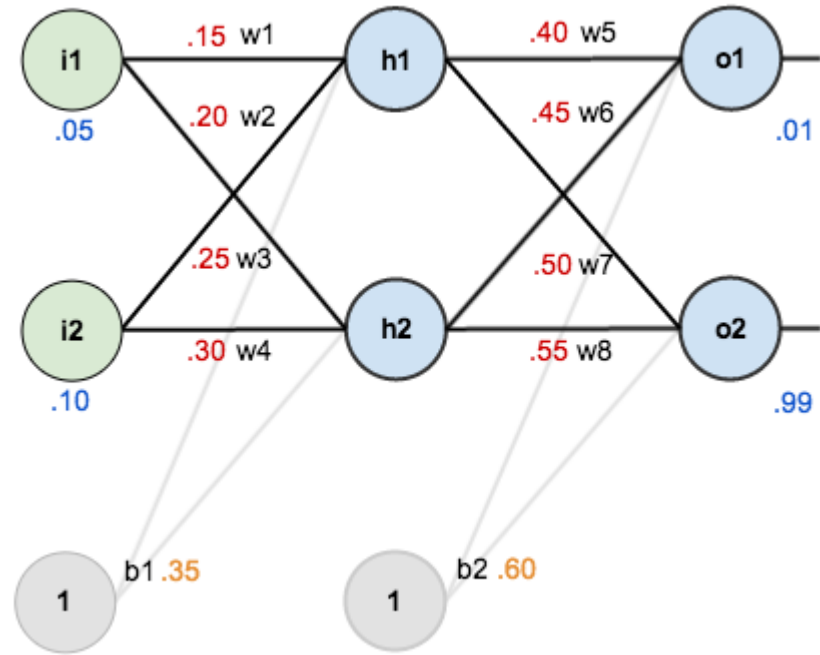
$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial}{\partial w} \left[\frac{1}{5} \sum_{i=1}^{i=5} (f(x_i) - y_i) \right]^2$$

$$\Delta w = \frac{\partial \mathcal{L}}{\partial w} = \frac{1}{5} \sum_{i=1}^{i=5} \frac{\partial}{\partial w} (f(x_i) - y_i)^2$$

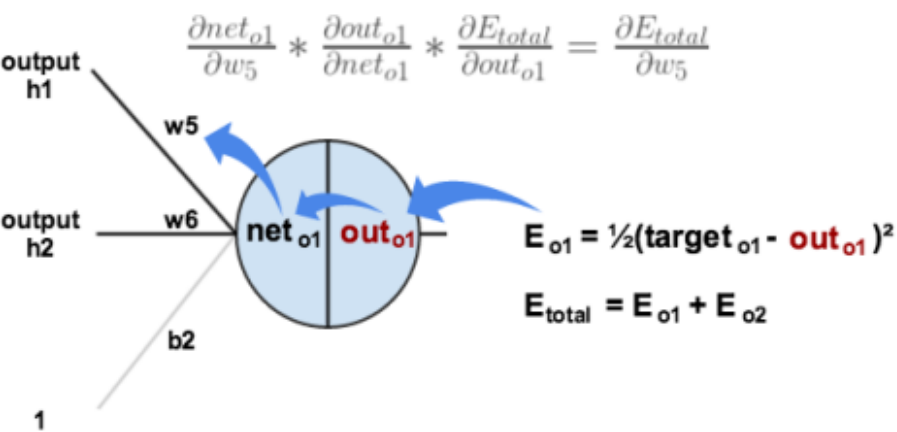
Practice Problem- Back Propagation



In order to have some numbers to work with, here are the **initial weights**, the **biases**, and **training inputs/outputs**:



The goal of backpropagation is to optimize the weights so that the neural network can learn how to correctly map arbitrary inputs to outputs.



Given a single training set: given inputs 0.05 and 0.10, we want the neural network to output 0.01 and 0.99.

Forward Pass

Carrying out the same process for h_2 we get:

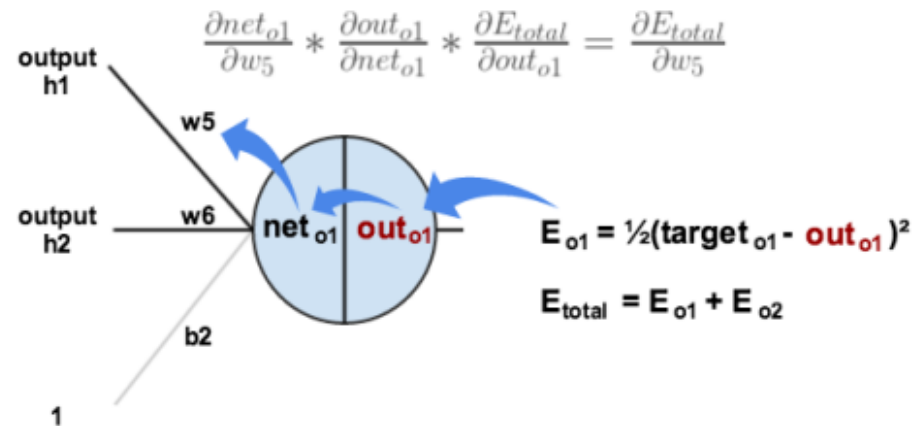
$$out_{h2} = 0.596884378$$

$$net_{h1} = w_1 * i_1 + w_2 * i_2 + b_1 * 1$$

$$net_{h1} = 0.15 * 0.05 + 0.2 * 0.1 + 0.35 * 1 = 0.3775$$

$$out_{h1} = \frac{1}{1+e^{-net_{h1}}} = \frac{1}{1+e^{-0.3775}} = 0.593269992$$

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$



Notation as per previous slides

Net O1= a1

Outo1=h1

- We repeat this process for the output layer neurons, using the output from the hidden layer neurons as inputs.
- Here's the output for o_1 :

Here's the output for o_1 :

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$net_{o1} = 0.4 * 0.593269992 + 0.45 * 0.596884378 + 0.6 * 1 = 1.105905967$$

$$out_{o1} = \frac{1}{1+e^{-net_{o1}}} = \frac{1}{1+e^{-1.105905967}} = 0.75136507$$

And carrying out the same process for o_2 we get:

$$out_{o2} = 0.772928465$$

Calculating the Total Error

- We can now calculate the error for each output neuron using the squared error function and sum them to get the total error:

$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

For example, the target output for o_1 is 0.01 but the neural network output 0.75136507, therefore its error is:

$$E_{o1} = \frac{1}{2}(target_{o1} - out_{o1})^2 = \frac{1}{2}(0.01 - 0.75136507)^2 = 0.274811083$$

Repeating this process for o_2 (remembering that the target is 0.99) we get:

$$E_{o2} = 0.023560026$$

The total error for the neural network is the sum of these errors:

$$E_{total} = E_{o1} + E_{o2} = 0.274811083 + 0.023560026 = 0.298371109$$

The Backwards Pass

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

$$out_{o1} = \frac{1}{1+e^{-net_{o1}}}$$

$$\frac{\partial out_{o1}}{\partial net_{o1}} = out_{o1}(1 - out_{o1}) = 0.75136507(1 - 0.75136507) = 0.186815602$$

Finally, how much does the total net input of $o1$ change with respect to w_5 ?

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\frac{\partial net_{o1}}{\partial w_5} = 1 * out_{h1} * w_5^{(1-1)} + 0 + 0 = out_{h1} = 0.593269992$$

Putting it all together:

$$\frac{\partial E_{total}}{\partial w_5} = -(target_{o1} - out_{o1}) * out_{o1}(1 - out_{o1}) * out_{h1}$$

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

$$\frac{\partial E_{total}}{\partial w_5} = 0.74136507 * 0.186815602 * 0.593269992 = 0.082167041$$

$$\frac{\partial E_{total}}{\partial w_5} = -(target_{o1} - out_{o1}) * out_{o1}(1 - out_{o1}) * out_{h1}$$

$$w_5^+ = w_5 - \eta * \frac{\partial E_{total}}{\partial w_5} = 0.4 - 0.5 * 0.082167041 = 0.35891648$$

$$w_6^+ = 0.408666186$$

$$w_7^+ = 0.511301270$$

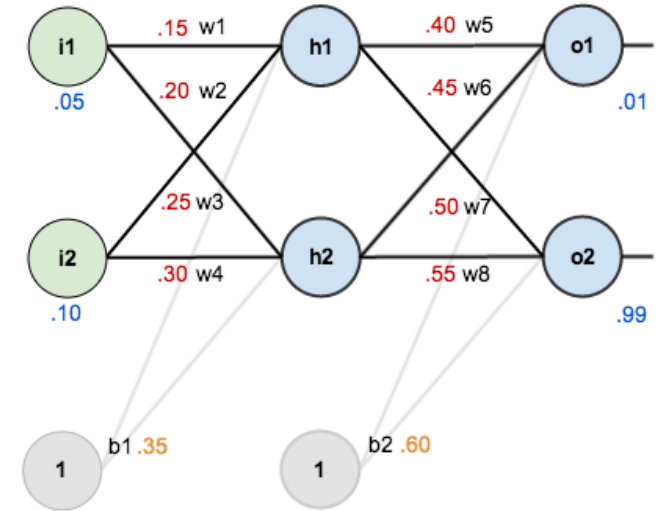
$$w_8^+ = 0.561370121$$

Hidden Layer

Next, we'll continue the backwards pass by calculating new values for w_1 , w_2 , w_3 , and w_4 .

Big picture, here's what we need to figure out:

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$



We can now update w_1 :

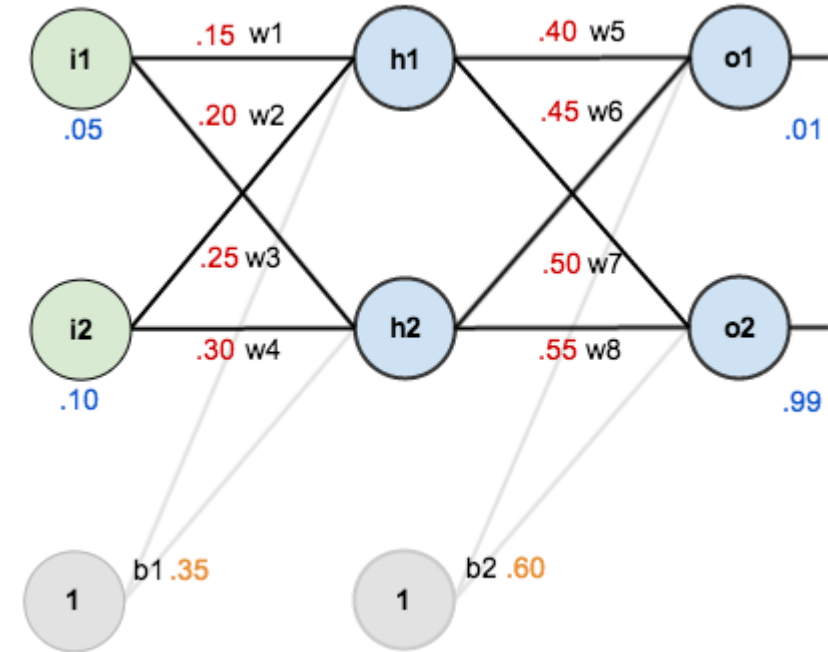
$$w_1^+ = w_1 - \eta * \frac{\partial E_{total}}{\partial w_1} = 0.15 - 0.5 * 0.000438568 = 0.149780716$$

Repeating this for w_2 , w_3 , and w_4

$$w_2^+ = 0.19956143$$

$$w_3^+ = 0.24975114$$

$$w_4^+ = 0.29950229$$



Cross Entropy Loss Function

Cross Entropy Loss- binary class

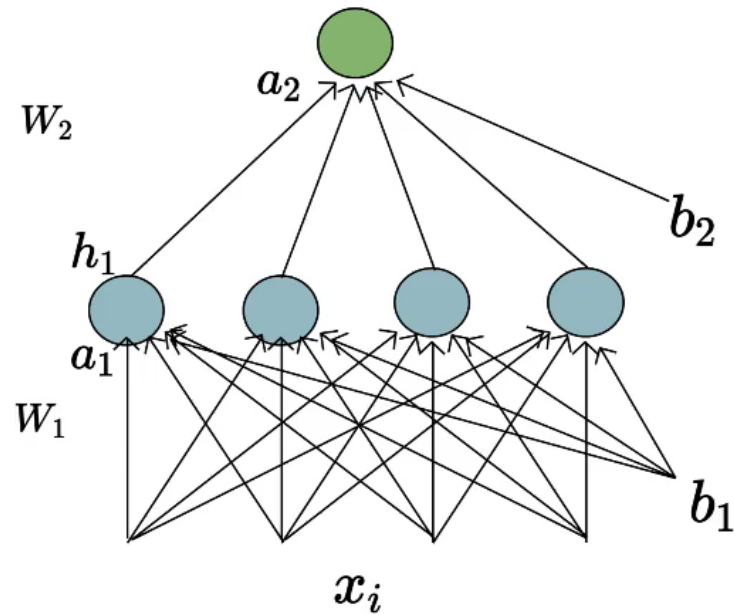
$$L(\Theta) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

Cross Entropy Loss: Multiclass

$$L(\Theta) = - \sum_{i=1}^k y_i \log(\hat{y}_i)$$

- Also called **logarithmic loss**, **log loss** or **logistic loss**.
- Each predicted class probability is compared to the actual class desired output 0 or 1 and a score/loss is calculated that **penalizes the probability based on how far it is from the actual expected value**.
- The **penalty is logarithmic in nature** yielding a large score for large differences close to 1 and small score for small differences tending to 0.
- Cross-entropy loss is used when adjusting model weights during training. **The aim is to minimize the loss**, i.e, the smaller the loss the better the model. A **perfect model has a cross-entropy loss of 0**.

Loss function for binary class classification



$$b = [0.5 \quad 0.3]$$

$$W_1 = \begin{bmatrix} 0.9 & 0.2 & 0.4 & 0.3 \\ -0.5 & 0.4 & 0.3 & 0.3 \\ 0.1 & 0.1 & -0.1 & 0.2 \\ -0.2 & 0.5 & 0.5 & 0.7 \end{bmatrix}$$

$$W_2 = [0.5 \quad 0.8 \quad -0.6 \quad 0.3]$$

$$x = [-0.6 \quad -0.6 \quad 0.2 \quad 0.3] \quad y = 0$$

Output :

$$a_1 = W_1 * x + b_1 = [0.01 \quad 0.71 \quad 0.42 \quad 0.63]$$

$$h_1 = \text{sigmoid}(a_1) = [0.50 \quad 0.67 \quad 0.60 \quad 0.65]$$

$$a_2 = W_2 * h_1 + b_2 = 0.921$$

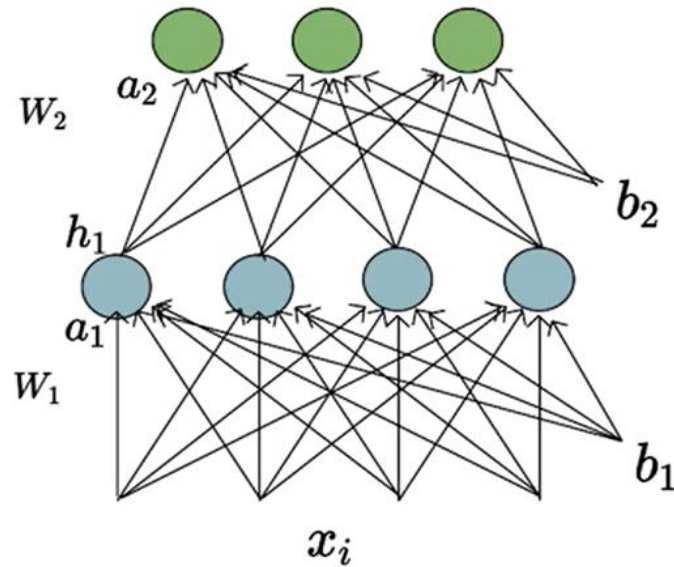
$$\hat{y} = \text{sigmoid}(a_2) = 0.7152$$

Cross Entropy Loss:

$$L(\Theta) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

$$\begin{aligned} L(\Theta) &= -1 * \log(1 - 0.7152) \\ &= 1.2560 \end{aligned}$$

Loss function for multi class classification



$$b = [0 \ 0]$$

$$W_1 = \begin{bmatrix} 0.1 & 0.3 & 0.8 & -0.4 \\ -0.3 & -0.2 & 0.5 & 0.5 \\ -0.3 & 0.1 & 0.5 & 0.4 \\ 0.2 & 0.5 & -0.9 & 0.7 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} 0.3 & 0.8 & -0.2 & -0.4 \\ 0.5 & -0.2 & -0.3 & 0.5 \\ 0.2 & 0.1 & 0.6 & 0.6 \end{bmatrix}$$

$$x = [0.6 \ 0.4 \ 0.6 \ 0.1] \quad y = [0 \ 0 \ 1]$$

Output :

$$a_1 = W_1 * x + b_1 = [0.62 \ 0.09 \ 0.2 \ -0.15]$$

$$h_1 = \text{sigmoid}(a_1) = [0.65 \ 0.52 \ 0.55 \ 0.46]$$

$$a_2 = W_2 * h_1 + b_2 = [0.32 \ 0.29 \ 0.85]$$

$$\hat{y} = \text{softmax}(a_2) = [0.2718 \ 0.2634 \ 0.4648]$$

Cross Entropy Loss:

$$L(\Theta) = - \sum_{i=1}^k y_i \log(\hat{y}_i)$$

$$\begin{aligned} L(\Theta) &= -1 * \log(0.4648) \\ &= 0.7661 \end{aligned}$$

Hyper parameter tuning

Algorithms

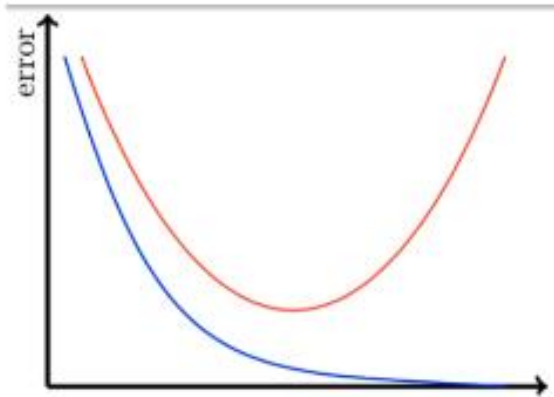
- Vanilla/Momentum /Nesterov GD
- AdaGrad
- RMSProp
- Adam

Strategies

- Batch
- Mini-Batch (32, 64, 128)
- Stochastic
- Learning rate schedule

Network Architectures

- Number of layers
- Number of neurons



Initialization Methods

- Xavier
- He

Activation Functions

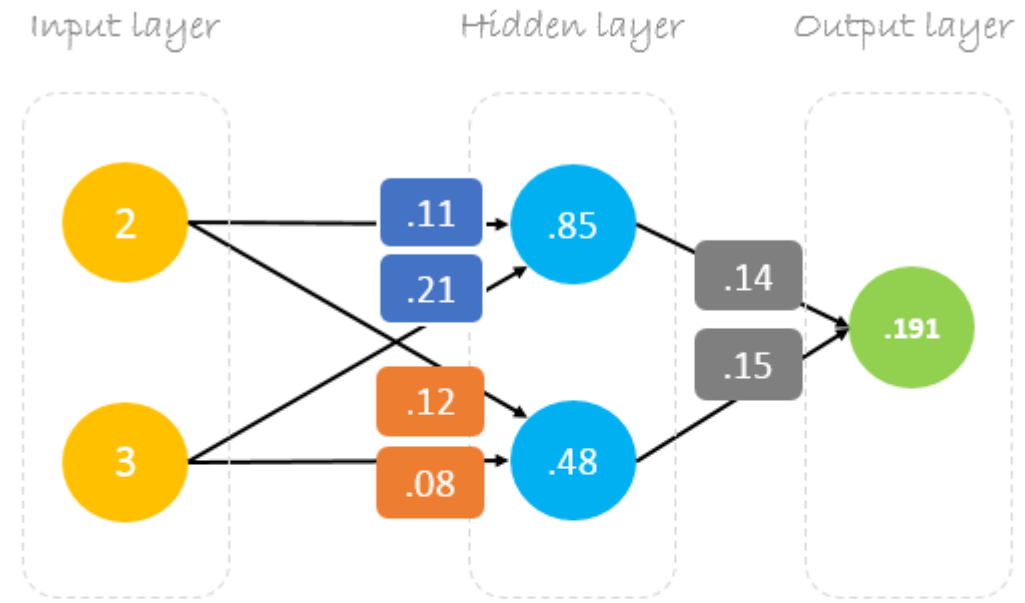
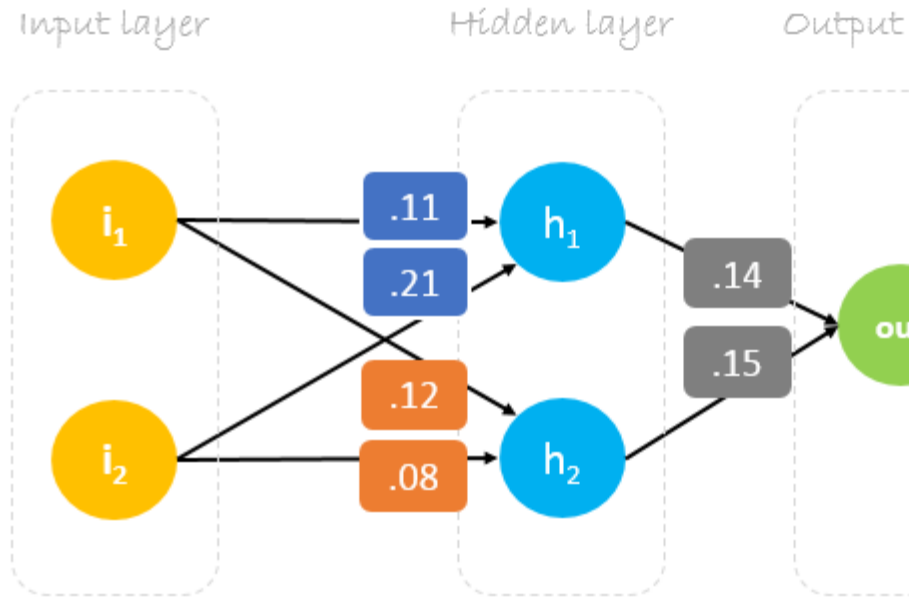
- tanh (RNNs)
- relu (CNNs, DNNs)
- leaky relu (CNNs)

Regularization

- L2
- Early stopping
- Dataset augmentation
- Drop-out
- Batch Normalizat

Namah Shivaya

<https://hmkcode.com/a>



$$\begin{bmatrix} 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 0.11 & 0.12 \\ 0.21 & 0.08 \end{bmatrix} = \begin{bmatrix} 0.85 & 0.48 \end{bmatrix} \cdot \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} = \begin{bmatrix} 0.191 \end{bmatrix}$$

$$2 \times .11 + 3 \times .21 = .85$$

$$.85 \times .14 + .48 \times .15 = .191$$

$$2 \times .12 + 3 \times .08 = .48$$

Matrix multiplication

Details

<https://visualstudiomagazine.com/articles/2014/04/01/neural-network-cross-entropy-error.aspx>