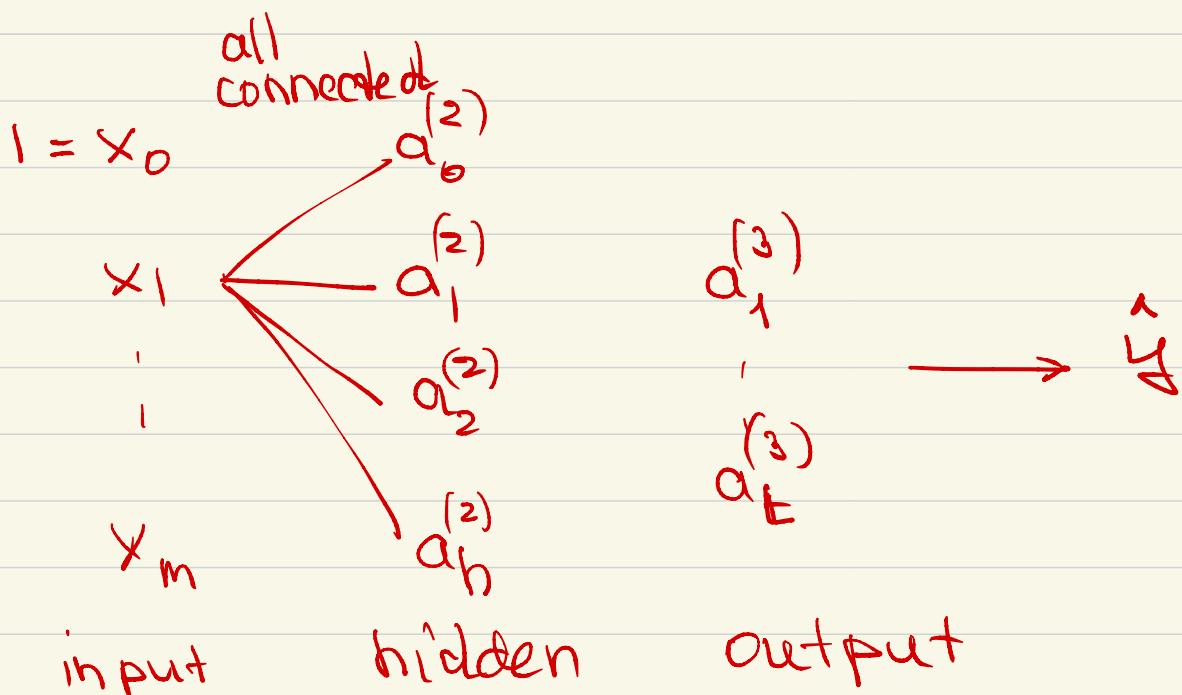


E. Multi-layer perceptron (or ANN)

This 'machine' consists of an input layer of m units, a number n_h of hidden layers of h_k units, $k = 1, \dots, n_h$ and an output layer of t units.

For simplicity we will consider

$$n_h = 1 \text{ and } h_1 = h.$$



Weights : $w_{j,k}^{(e)}$
 connecting units
 k in layer e with
 j in layer $e+1$

Feed-forward network :

$$\tilde{z}_1^{(2)} = a_0^{(1)} w_{1,0}^{(1)} + \dots + a_m^{(1)} w_{1,m}^{(1)}$$

with

$$\begin{pmatrix} a_0^{(1)} \\ a_1^{(1)} \\ \vdots \\ a_m^{(1)} \end{pmatrix} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_m \end{pmatrix}$$

$$a_1^{(2)} = \phi(\tilde{z}_1^{(2)}) \quad \text{with}$$

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

More general :

$$\underline{\underline{z}}_j^{(2)} = \sum_{k=0}^m \underline{a}_k^{(1)} \underline{w}_{j,k}^{(1)} \quad j = 1, \dots, h$$

$$\text{or } \underline{\underline{z}}^{(2)} = \underline{W}^{(1)} \underline{a}^{(1)}$$

$$\underline{a}^{(2)} = \phi(\underline{\underline{z}}^{(2)}) \quad \begin{matrix} \underline{W}^{(1)}: h \times (m+1) \\ \underline{a}^{(1)}: (m+1) \times 1 \\ \underline{\underline{z}}^{(2)}, \underline{a}^{(2)}: h \times 1 \end{matrix}$$

Suppose, we have now N samples of training data, $i = 1, \dots, N$, then

$$\underline{\underline{z}}^{(2)} = \underline{W}^{(1)} (\underline{A}^{(1)})^T$$

where $\underline{A}^{(1)} : N \times m+1$

$$\underline{\underline{z}}^{(2)} : h \times N$$

and $\underline{A}^{(2)} = \phi(\underline{\underline{z}}^{(2)}) \quad \underline{A}^{(2)} : h \times N$

↑ componentwise

For the output layer (layer(3))
we can also write :

$$\underline{\underline{z}}^{(3)} = W^{(2)} A^{(2)}$$

$$W^{(2)} : t \times h$$

$$A^{(2)} : h \times N$$

$$A^{(3)}, \underline{\underline{z}}^{(3)} : t \times N$$

with $A^{(3)} = \phi(\underline{\underline{z}}^{(3)})$

Output : $\begin{pmatrix} y_i \\ \vdots \\ y_j \\ \vdots \\ y_t \end{pmatrix}_j = A_{j,i}^{(3)}$
 $j = 1, \dots, t$
 $i = 1, \dots, N$

F. Training an ANN

Cost function (logistic)

To be compatible with notation of Raschka, let $a_j^i = h_j^{(i)}$, then the cost function is given by

$$\begin{aligned} \mathcal{J}(w) = & - \sum_{i=1}^N \sum_{j=1}^t \left\{ y_j^i \ln \phi(z_j^i) + \right. \\ & \left. + (1 - y_j^i) \ln (1 - \phi(z_j^i)) \right\} \\ & + \frac{\lambda}{2} \sum_{l=1}^{L-1} \sum_{i=1}^{u_l} \sum_{j=1}^{u_{l+1}} w_{j,i}^{(l)} \end{aligned}$$

weight over all layers.

where also $z_j^i = z_{j,i}^{(3)}$

Again the weights will be updated by the gradient of J , and hence we need,

$$\frac{\partial}{\partial w_{j,k}^{(e)}} J(w^{(1)}, w^{(2)})$$

where $w^{(1)}$ connects the hidden layer to the input layer and $w^{(2)}$ the hidden layer to the output layer. The weights will be updated through

$$w_{k+1}^{(e)} = w_k^{(e)} - \eta \Delta_{k+1}^{(e)}$$

where $\Delta^{(e)}$ is determined from the so-called backpropagation algorithm.