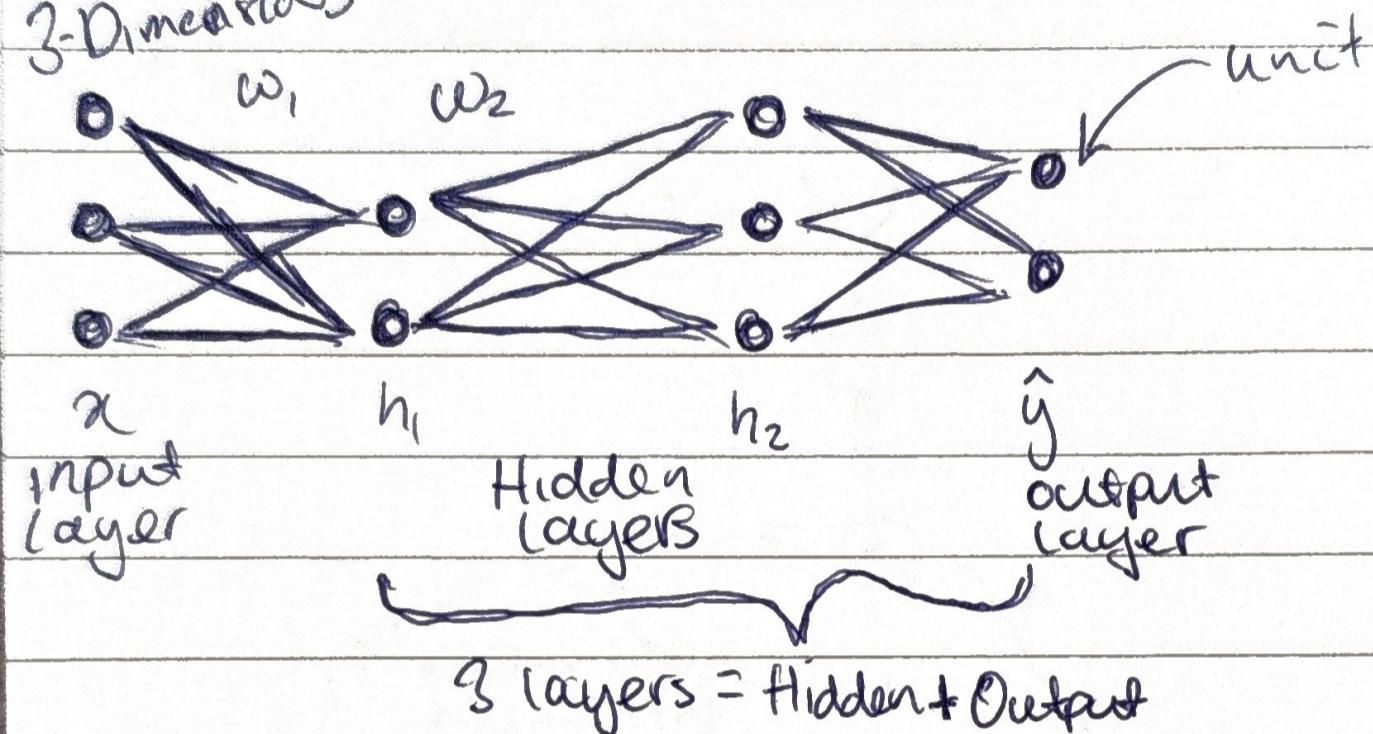


(MCP)

Multi-layer Perceptron

Smallest MCP = 2 layers

3 dimensions



- each unit connects to each unit in the next
- No intra layer connections

feed forward example:

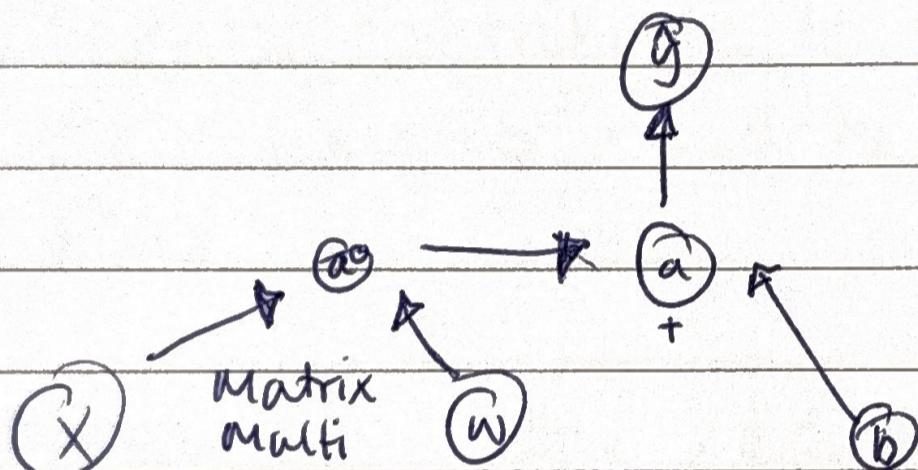
Input $\xrightarrow{\text{linear map}}$ Hidden \rightarrow Output

$$h(x) = w_1 x$$

Computational Graphs

Nodes = variables (inputs, output, parameters)
 Edges = mapping

$$\hat{y} = \text{relu}(x \cdot w + b)$$



(3)

Backpropagation

backpropagating the loss gradient from the output layer to all of the layers

All layers contribute to the loss so their weights need to be updated in SGD

Backprop is the process by which the loss gradient computed @ the output is propagated backward to the inner layers

Output layer is a function of a function
so we need chain rule for diff

because we have multiple dimensions (i.e. units) we need to compute partial derivatives $\frac{\partial L}{\partial x}$ instead of $\frac{dL}{dx}$

Backpropagation Equation Derivation

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_{ji}} \quad \text{(chain rule)}$$

where \rightarrow loss gradient for the weight w_{ji} connecting unit i in layer $T-1$ to unit j in layer T

$$z_j = \sum_i a_i w_{ji} \quad w_{ji} = \text{linear output model for unit } j \\ g(z_j)$$

a_i = the input for unit j = activation output for unit i

\rightarrow Denote $\frac{\partial L}{\partial z_j}$ as ∇_j

$$= \frac{\partial L}{\partial w_{ji}} = \nabla_j \cdot \frac{\partial z_j}{\partial w_{ij}}$$

Sub $z_j = \sum_i a_i w_{ji}$ & perform DIFF

$$\frac{\partial L}{\partial w_{ji}} = \nabla_j \cdot a_i$$

$$\text{Sub } a_i = g(z_i) = \nabla_j \cdot g(z_i)$$

(5)

Solving ∇_j

j represents a layer & a layer has n number of units

u will represent a unit in the layer

$$\Rightarrow \nabla_u = \frac{\partial L}{\partial z_u}$$

apply the chain rule

$$\nabla_u = \sum \frac{\partial L}{\partial z_k} \cdot \frac{\partial z_k}{\partial z_u}$$

$$\text{denoting } \frac{\partial L}{\partial z_k} \text{ as } \nabla_k : \nabla_u = \sum \nabla_k \cdot \frac{\partial z_k}{\partial z_u}$$

$$\text{recall } z_k = \sum a_j \cdot w_{kj} \text{ & then } a_j = g(z_j)$$

More steps in slides

final form

$$\frac{\partial L}{\partial w_{ji}} = \left(g(z_j) \sum v_k \cdot w_{kj} \right) \cdot g(z_i)$$

↓

Derivative of activation
for unit j in layer T

activation for
unit i in layer
 $T-1$

loss gradient for unit
 k in layer $T+1$

Output layer has $t+1$ layer

$$= \frac{\partial L}{\partial y} \cdot g(z_i)$$

Input layer has no $t-1$ activation layer

$$= (\underline{m}) \cdot \underline{n_i}$$