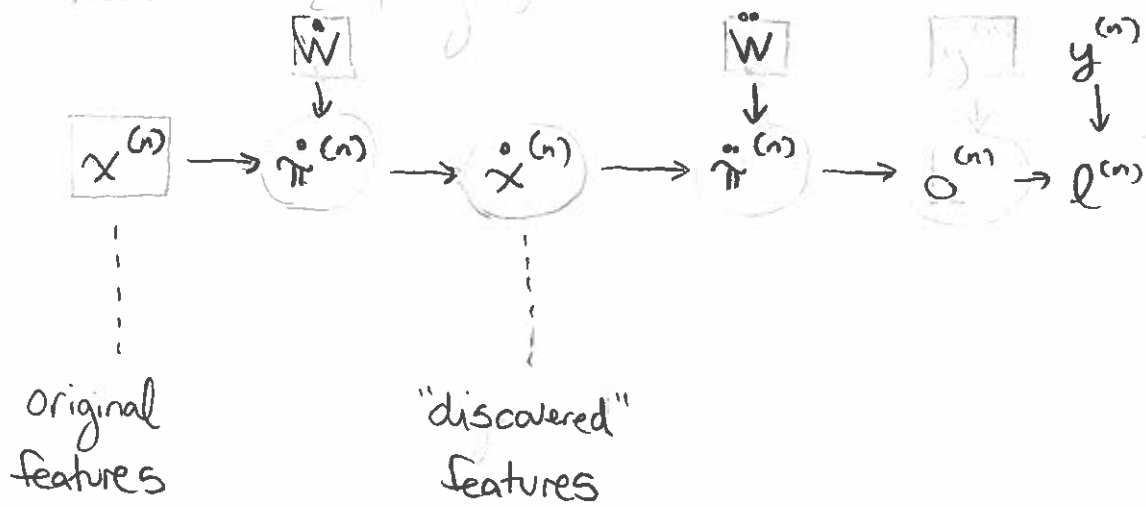
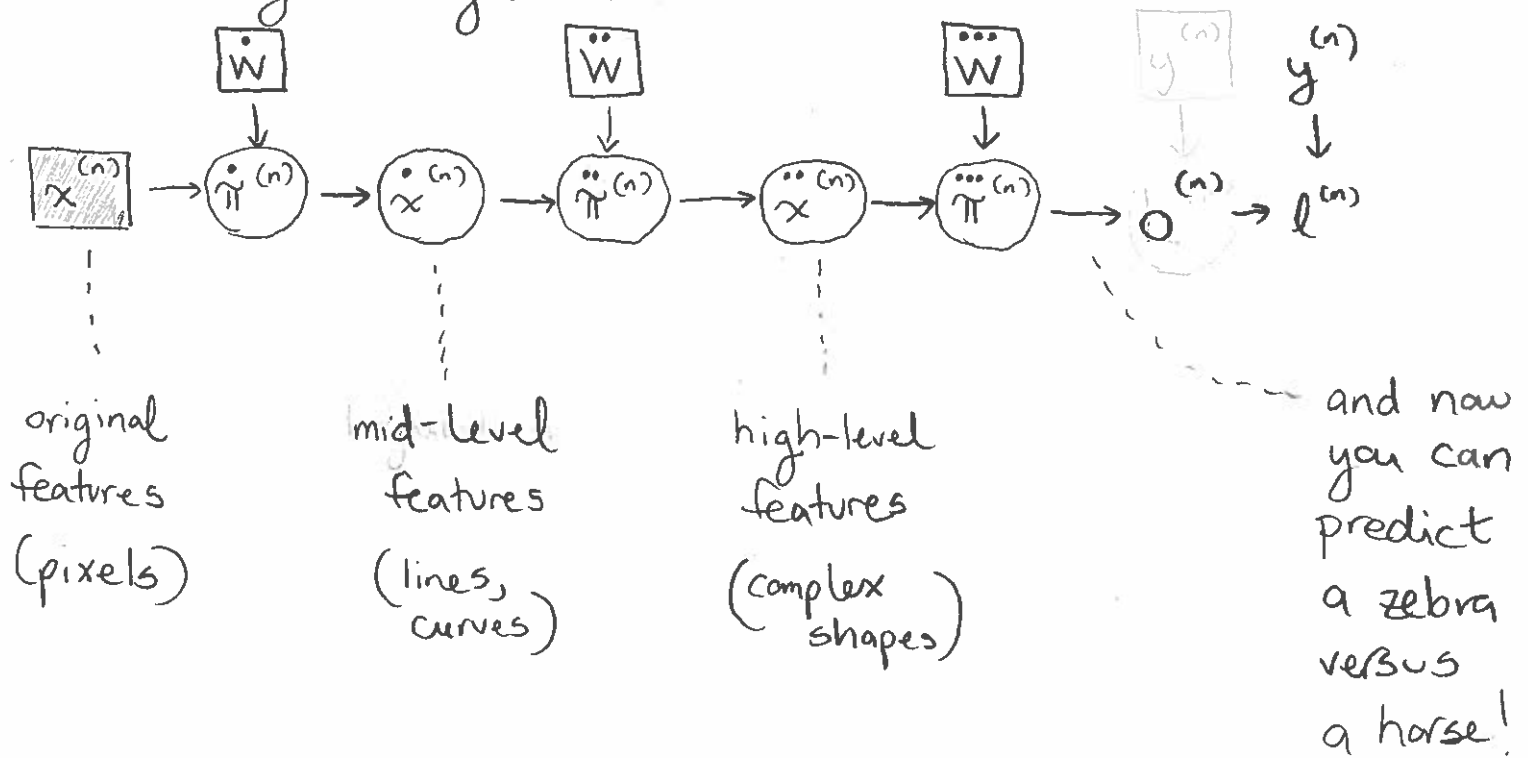


NEURAL NETWORKS AND BACKPROPAGATION

① Consider again our feature discovery network: (revisited)

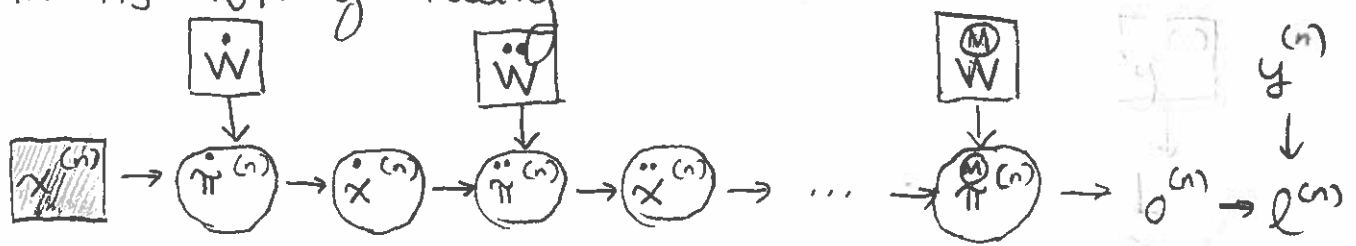


② We could consider generalizing this model to provide multiple layers of feature discovery, e.g. for image recognition:



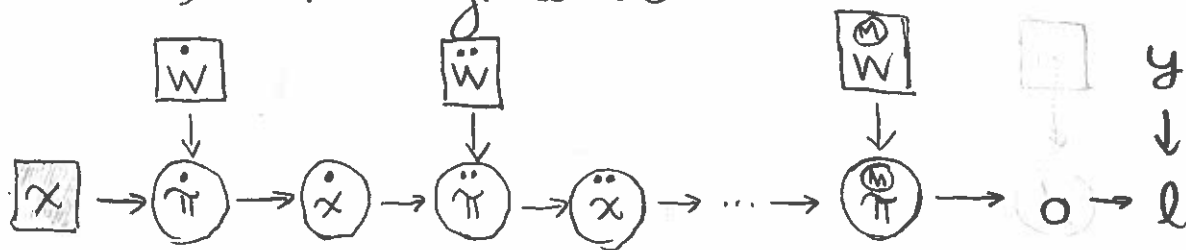
NEURAL NETWORKS AND BACKPROPAGATION

③ In its full generality:

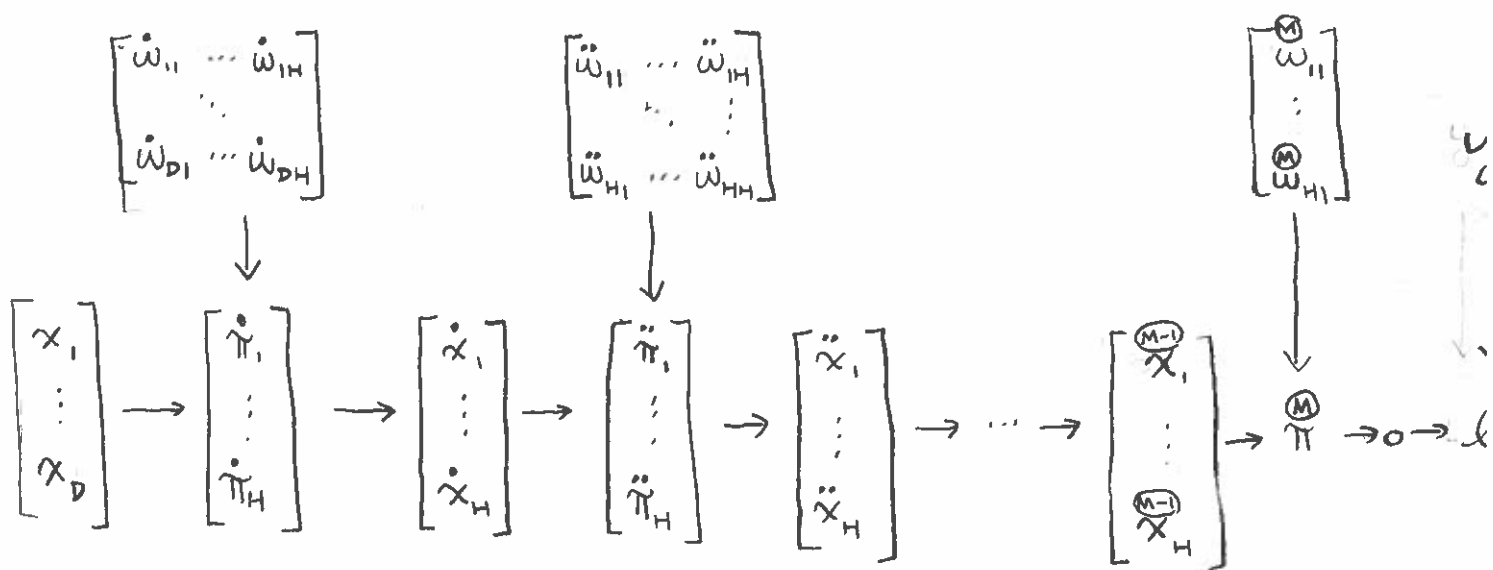


this is called an M-layer feedforward neural network.

Let's drop all those (n) superscripts for convenience (we'll bring them back when needed to avoid confusion). This gives us:



Just in case we've forgotten which of these are vectors and which are matrices, here it is explicitly:

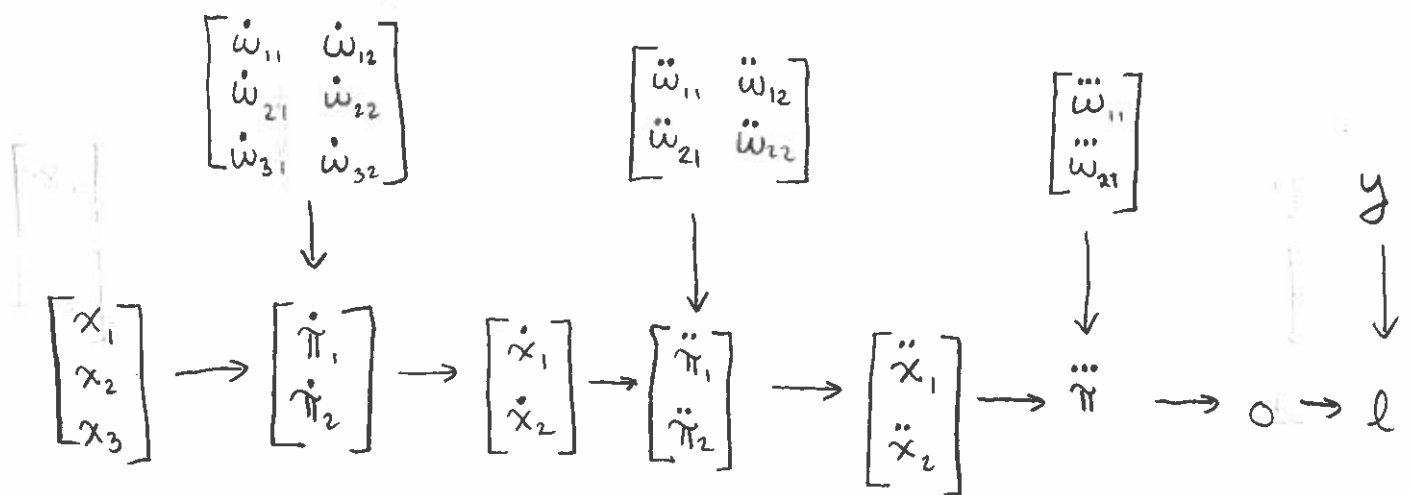


We assume each "feature discovery" layer discovers H features.

NEURAL NETWORKS AND BACKPROPAGATION

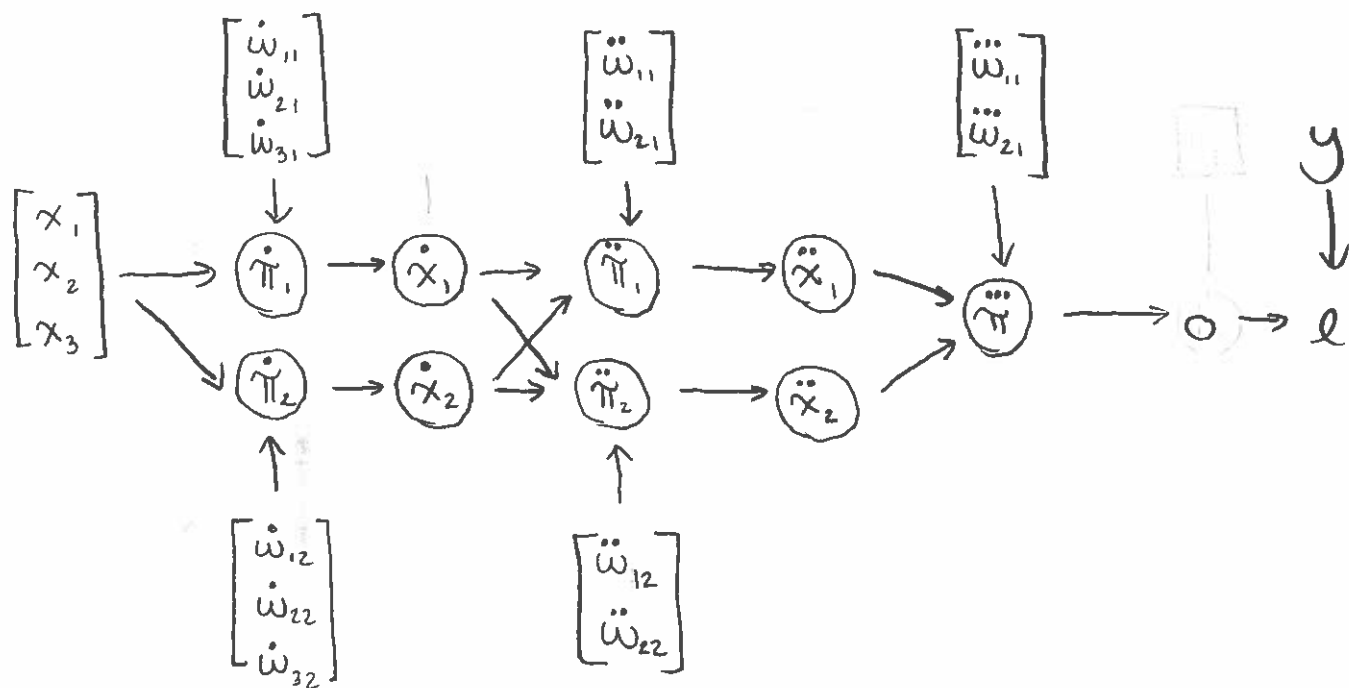
- ④ To train this model using gradient descent, we need to be able to compute $\frac{\partial \mathcal{L}}{\partial \overset{(m)}{w}_{ij}}$ for each weight $\overset{(m)}{w}_{ij}$.

Before doing this in its full generality, let's see how we can compute these derivatives for a 3-layer network where $H=2$ and $D=3$.



NEURAL NETWORKS AND BACKPROPAGATION

- ⑤ As we did before for the feature discovery network, let's break down the endogenous variables into scalars to make it easier to apply the Chain Rule of Partial Derivatives:



- ⑥ Our goal is to compute (for all relevant i, j):

$$\frac{\partial \mathcal{L}}{\partial \dot{w}_{ij}} \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \ddot{w}_{ij}} \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \ddot{w}_{ij}}$$

First, we can observe that $\ddot{\pi}$ separates \mathcal{L} from all \dot{w}_{ij} , so:

$$\frac{\partial \mathcal{L}}{\partial \dot{w}_{ij}} = \frac{\partial \mathcal{L}}{\partial \ddot{\pi}} \cdot \frac{\partial \ddot{\pi}}{\partial \dot{w}_{ij}}$$

This is just the standard derivative of the loss function.

NEURAL NETWORKS AND BACKPROPAGATION

⑦ So the challenge is to compute $\frac{\partial \ddot{\pi}}{\partial \ddot{w}_{ij}}$ for any layer m .

It's straightforward for $m=3$:

$$\frac{\partial \ddot{\pi}}{\partial \ddot{w}_{ij}} = \frac{\partial}{\partial \ddot{w}_{ij}} \left(\begin{bmatrix} \ddot{x}_1 & \ddot{x}_2 \end{bmatrix} \begin{bmatrix} \ddot{w}_{11} \\ \ddot{w}_{21} \end{bmatrix} \right) = \ddot{x}_i$$

⑧ What about $m=2$?

$$\frac{\partial \ddot{\pi}}{\partial \ddot{w}_{ij}} = \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}_j} \cdot \frac{\partial \ddot{\pi}_j}{\partial \ddot{w}_{ij}} \quad \left[\begin{array}{l} \ddot{\pi}_j \text{ separates } \ddot{\pi} \text{ from } \ddot{w}_{ij}, \\ \text{so Chain Rule applies} \end{array} \right]$$

$$= \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}_j} \cdot \dot{x}_i \quad \left[\frac{\partial \ddot{\pi}_j}{\partial \ddot{w}_{ij}} = \dot{x}_i \right]$$

⑨ What about $m=1$?

$$\frac{\partial \ddot{\pi}}{\partial \dot{w}_{ij}} = \frac{\partial \ddot{\pi}}{\partial \dot{\pi}_j} \frac{\partial \dot{\pi}_j}{\partial \dot{w}_{ij}} \quad \left[\begin{array}{l} \dot{\pi}_j \text{ separates } \ddot{\pi} \text{ from } \dot{w}_{ij}, \\ \text{so Chain Rule applies} \end{array} \right]$$

$$= \frac{\partial \ddot{\pi}}{\partial \dot{\pi}_j} \cdot x_i \quad \left[\frac{\partial \dot{\pi}_j}{\partial \dot{w}_{ij}} = x_i \right]$$

NEURAL NETWORKS AND BACKPROPAGATION

⑩ In summary:

$$\frac{\partial \ddot{\pi}}{\partial \ddot{w}_{ij}} = \ddot{x}_i$$

$$\frac{\partial \ddot{\pi}}{\partial \ddot{w}_{ij}} = \dot{x}_i \frac{\partial \ddot{\pi}}{\partial \dot{\pi}_j}$$

$$\frac{\partial \ddot{\pi}}{\partial \dot{w}_{ij}} = x_i \frac{\partial \ddot{\pi}}{\partial \dot{\pi}_j}$$

for the general case:

$$\frac{\partial \pi^{(M)}}{\partial \dot{w}_{ij}^{(m)}} = \overset{(m-1)}{x_i} \cdot \frac{\partial \pi^{(M)}}{\partial \dot{\pi}_j^{(m)}}$$

so how do we compute this term?

⑪ Consider $\frac{\partial \ddot{\pi}}{\partial \dot{\pi}_j}$ for our 3-layer network.

$$\frac{\partial \ddot{\pi}}{\partial \dot{\pi}_j} = \frac{\partial \ddot{\pi}}{\partial \dot{x}_j} \frac{\partial \dot{x}_j}{\partial \dot{\pi}_j}$$

$\left[\begin{array}{l} \dot{x}_j \text{ separates } \ddot{\pi} \text{ from } \dot{\pi}_j, \\ \text{so Chain Rule applies} \end{array} \right]$

$$= \left(\sum_{h=1}^2 \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}_h} \frac{\partial \ddot{\pi}_h}{\partial \dot{x}_j} \right) \frac{\partial \dot{x}_j}{\partial \dot{\pi}_j}$$

$\left[\begin{array}{l} \{ \ddot{\pi}_1, \ddot{\pi}_2 \} \text{ separates } \ddot{\pi} \text{ from } \dot{x}_j \\ \text{so Chain Rule applies} \end{array} \right]$

$$= \frac{\partial \dot{x}_j}{\partial \dot{\pi}_j} \sum_{h=1}^2 \frac{\partial \ddot{\pi}_h}{\partial \dot{x}_j} \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}_h}$$

$$= a'(\dot{\pi}_j) \sum_{h=1}^2 \ddot{w}_{hj} \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}_h}$$

but this can be computed recursively in the same way!



NEURAL NETWORKS AND BACKPROPAGATION

⑫ In summary (for our 3-layer example):

$$\frac{\partial \ddot{\pi}}{\partial \ddot{\pi}_j} = a'(\ddot{\pi}_j) \sum_{h=1}^2 \ddot{w}_{hj} \frac{\partial \ddot{\pi}}{\partial \ddot{\pi}_h}$$

and for the general case:

$$\frac{\partial \pi^{(M)}}{\partial \pi^{(M)}} = 1 \quad \leftarrow \text{base case}$$

$$\frac{\partial \pi^{(M)}}{\partial \pi_j^{(m)}} = a'(\pi_j^{(m)}) \sum_{h=1}^H w_{hj}^{(m+1)} \frac{\partial \pi^{(M)}}{\partial \pi_h^{(m+1)}} \quad \leftarrow \text{recursive step}$$

⑬ Putting it all together, we have cobbled together a strategy for computing every partial derivative $\frac{\partial \mathcal{L}}{\partial w_{ij}^{(m)}}$:

BACKPROPAGATION:

(a) for m in $\{M-1, \dots, 1\}$ and j in $\{1, \dots, H\}$:

compute $\frac{\partial \pi^{(M)}}{\partial \pi_j^{(m)}} = a'(\pi_j^{(m)}) \sum_{h=1}^H w_{hj}^{(m+1)} \frac{\partial \pi^{(M)}}{\partial \pi_h^{(m+1)}}$

computed in an earlier iteration, so we loop backwards for $M-1$

(b) $\frac{\partial \mathcal{L}}{\partial w_{ij}^{(m)}} = \frac{\partial \mathcal{L}}{\partial \pi^{(m)}} \cdot x_i^{(m-1)} \cdot \frac{\partial \pi^{(M)}}{\partial \pi_j^{(m)}}$

just a simple derivative of the loss function

computed during (a)

NEURAL NETWORKS AND BACKPROPAGATION

⑭ Because we compute the partial derivatives $\frac{\partial \pi^{(M)}}{\partial \pi^{(m)}}$

starting from the final layer M and moving backwards, this algorithm is known as backpropagation.