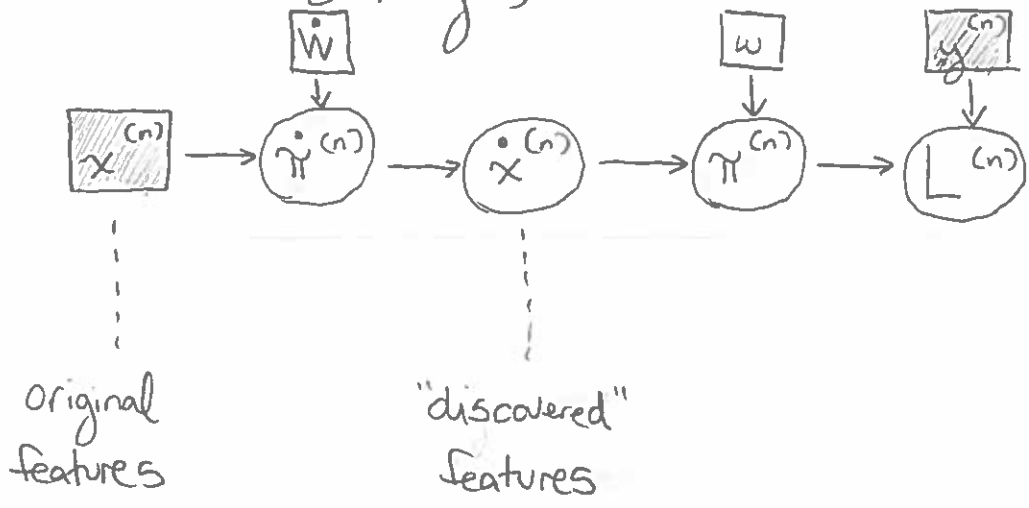
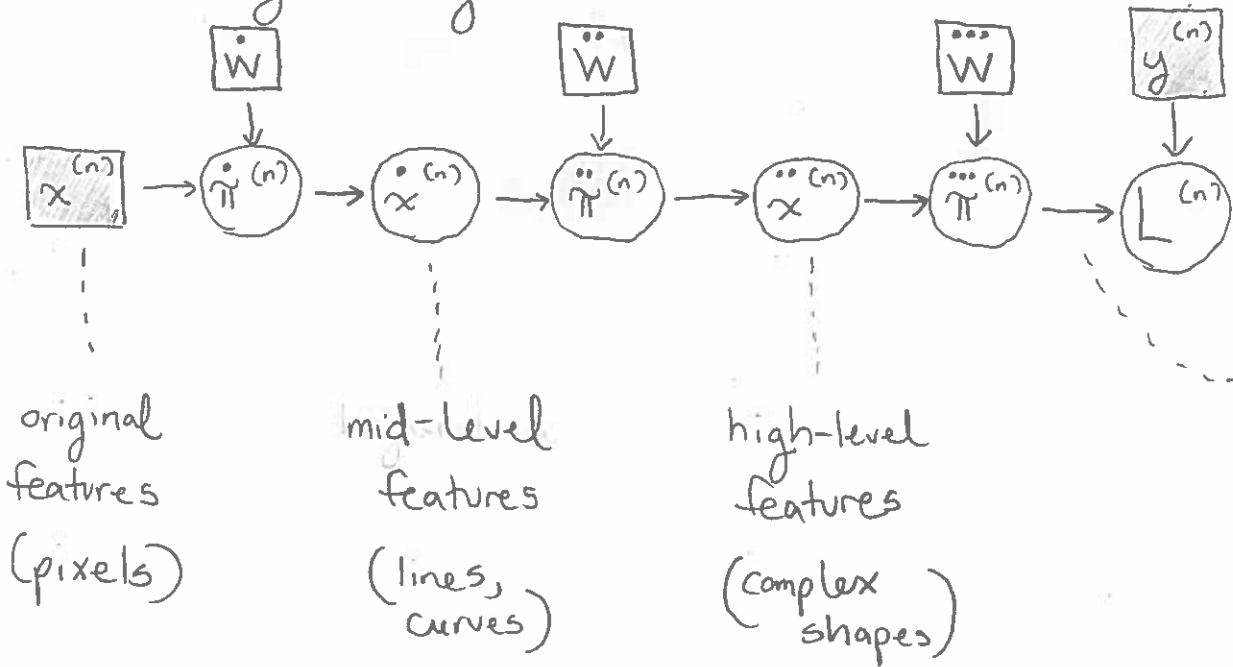


# NEURAL NETWORKS AND BACKPROPAGATION

① Consider again our feature discovery network (now drawn horizontally!):



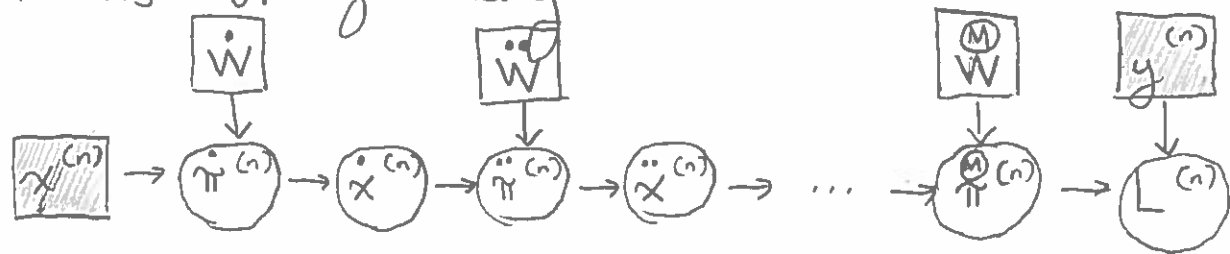
② We could consider generalizing this model to provide multiple layers of feature discovery, e.g. for image recognition:



and now you can predict a zebra versus a horse!

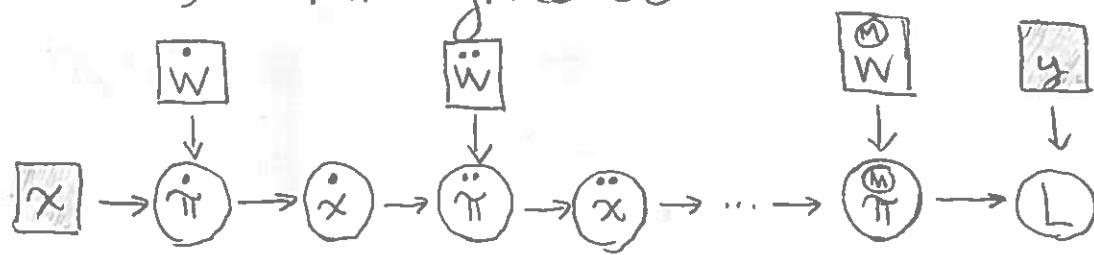
# NEURAL NETWORKS AND BACKPROPAGATION

③ In its full generality:

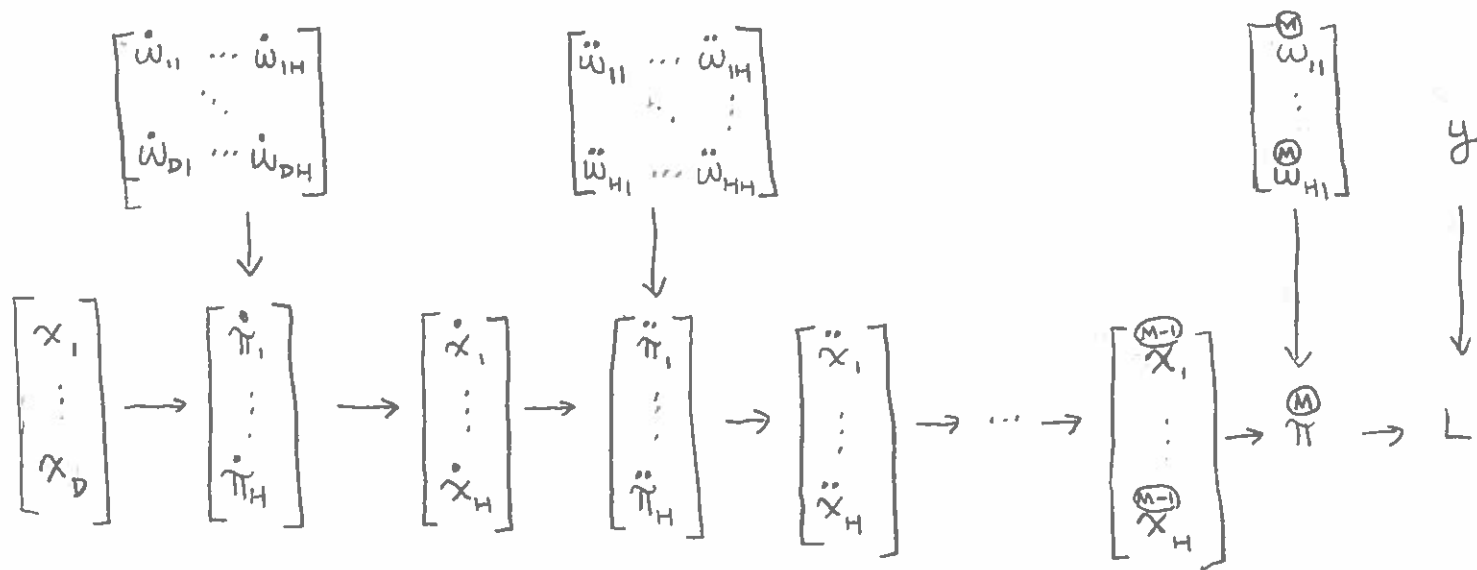


this is called an M-layer feedforward neural network.

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:

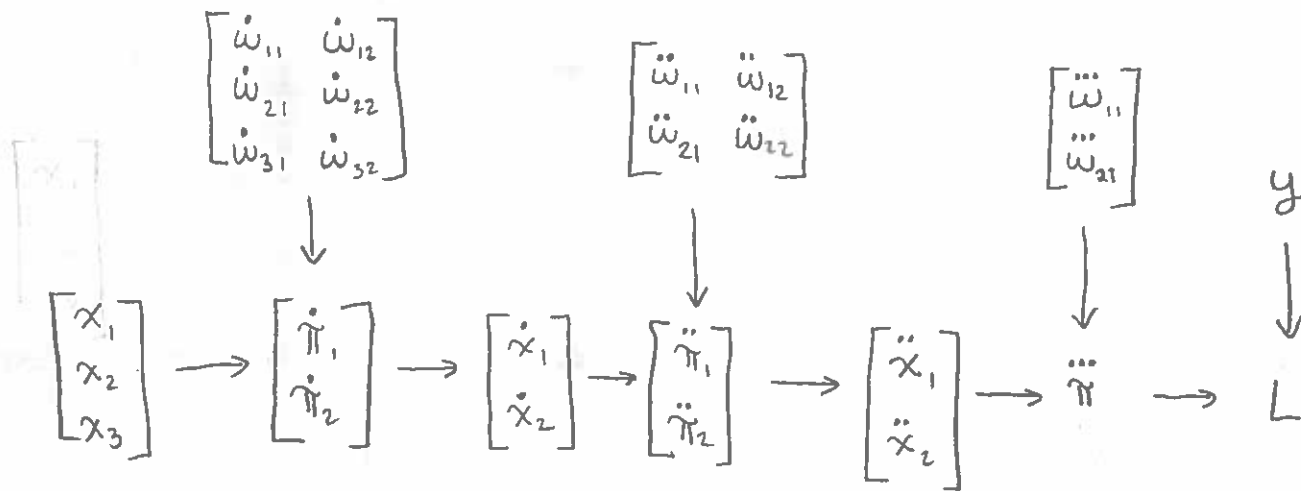


We assume each "feature discovery" layer discovers H features.

# NEURAL NETWORKS AND BACKPROPAGATION

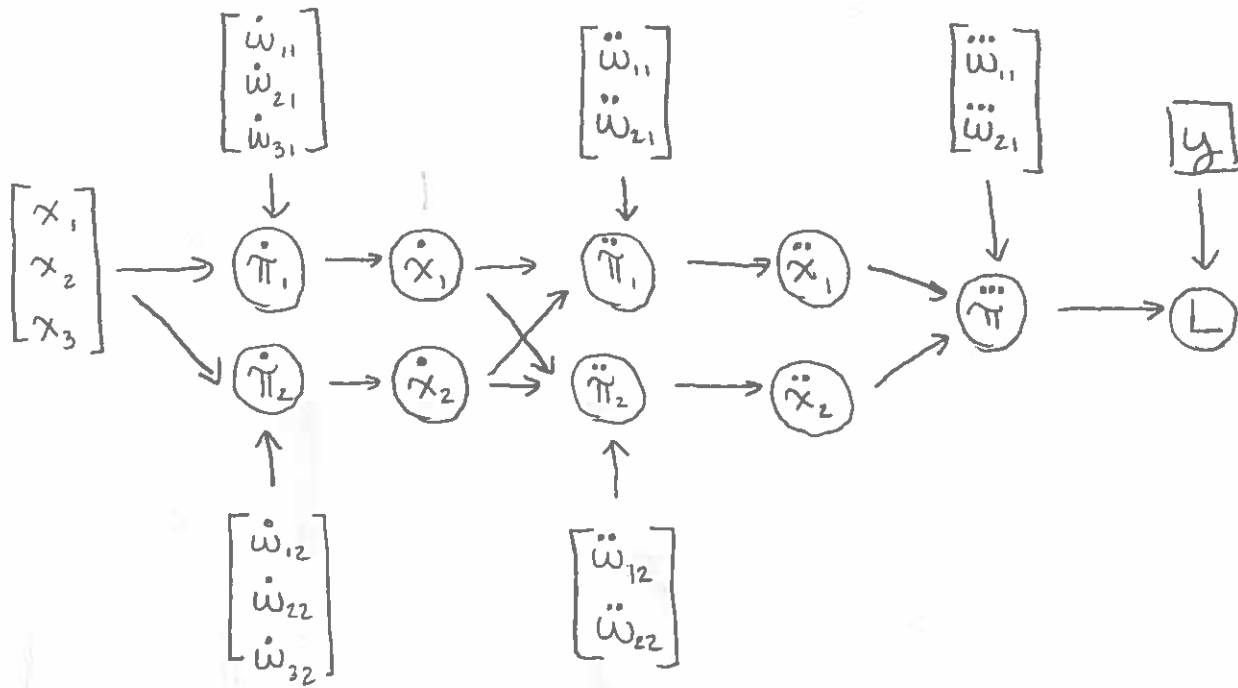
- ④ To train this model using gradient descent, we need to be able to compute  $\frac{\partial L}{\partial \tilde{w}_{ij}^{(m)}}$  for each weight  $\tilde{w}_{ij}^{(m)}$ .

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  $i, j$ ):

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

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

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

This is just the standard derivative of the loss function.

# NEURAL NETWORKS AND BACKPROPAGATION

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

It's straightforward for  $m=3$ :

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

⑧ What about  $m=2$ ?

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

⑨ What about  $m=1$ ?

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

# 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:

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

so how do we compute this term?

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

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

$\left[ \dot{x}_j \text{ separates } \ddot{\pi} \text{ from } \dot{\pi}_j, \right.$   
so Chain Rule applies

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

$\left[ \{ \ddot{\pi}_1, \ddot{\pi}_2 \} \text{ separates } \ddot{\pi} \text{ from } \dot{x}_j, \right.$   
so Chain Rule applies

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

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

but this can be computed recursively in the same way!



# NEURAL NETWORKS AND BACKPROPAGATION

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

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

and for the general case:

$$\frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{M}}{\pi}} = 1 \quad \leftarrow \text{base case}$$

$$\frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{m}}{\pi}_j} = a'(\overset{\textcircled{m}}{\pi}_j) \sum_{h=1}^H \overset{\textcircled{m+1}}{w}_{hj} \frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{m+1}}{\pi}_h} \quad \leftarrow \text{recursive step}$$

⑬ Putting it all together, we have cobbled together a strategy for computing every partial derivative  $\frac{\partial L}{\partial \overset{\textcircled{m}}{w}_{ij}}$ :

BACKPROPAGATION:

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

$$\text{compute } \frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{m}}{\pi}_j} = a'(\overset{\textcircled{m}}{\pi}_j) \sum_{h=1}^H \overset{\textcircled{m+1}}{w}_{hj} \frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{m+1}}{\pi}_h}$$

computed in an earlier iteration, since we loop backwards from  $M-1$

$$(b) \frac{\partial L}{\partial \overset{\textcircled{m}}{w}_{ij}} = \underbrace{\frac{\partial L}{\partial \overset{\textcircled{M}}{\pi}}}_{\text{just a simple derivative of the loss function}} \cdot \overset{\textcircled{m-1}}{x}_i \cdot \underbrace{\frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{m}}{\pi}_j}}_{\text{computed during (a)}}$$

## 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.