# Neural Networks and Backpropagation

① Consider again our feature discovery network (now drawn horizontally!):
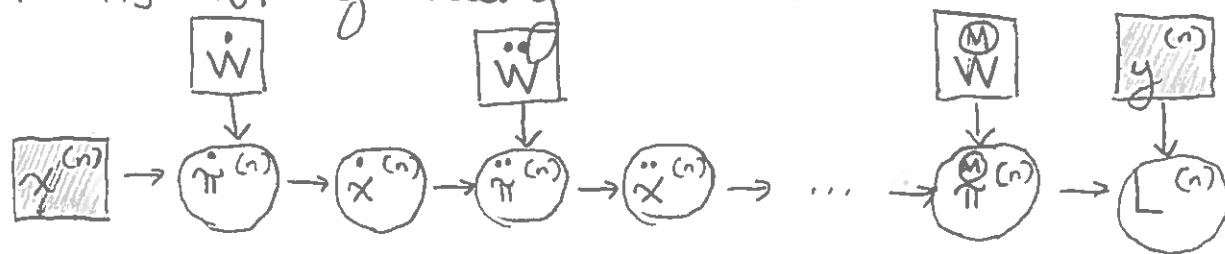
$$\boxed{\dot{W}} \qquad \boxed{W} \qquad \boxed{\overset{\text{\tiny ///}}{\cancel{W}}{}^{(n)}}$$

$$\boxed{x}^{(n)} \rightarrow \dot{\gamma}^{(n)} \rightarrow \dot{x}^{(n)} \rightarrow \gamma^{(n)} \rightarrow L^{(n)}$$

original features $\qquad\qquad$ "discovered" features

---

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

$$\boxed{\dot{W}} \qquad \boxed{\ddot{W}} \qquad \boxed{\dddot{W}} \qquad \boxed{y}^{(n)}$$

$$\boxed{x}^{(n)} \rightarrow \dot{\gamma}^{(n)} \rightarrow \dot{x}^{(n)} \rightarrow \ddot{\gamma}^{(n)} \rightarrow \ddot{x}^{(n)} \rightarrow \dddot{\gamma}^{(n)} \rightarrow L^{(n)}$$

original features (pixels) $\qquad$ mid-level features (lines, curves) $\qquad$ high-level features (complex shapes) $\qquad$ and now you can predict a zebra versus a horse!
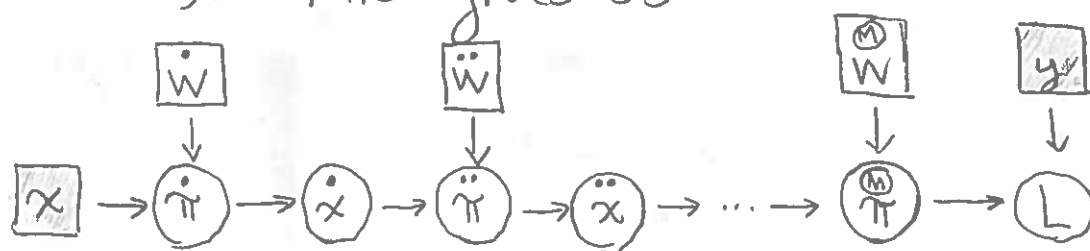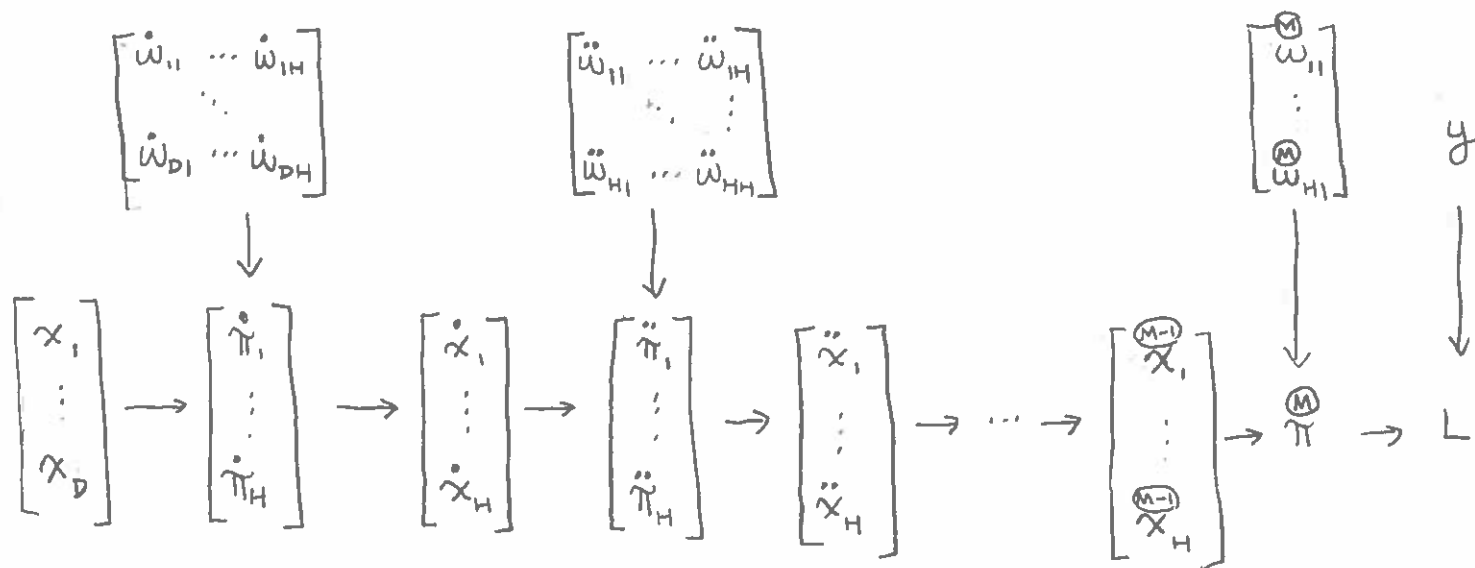
③ In its full generality:



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:



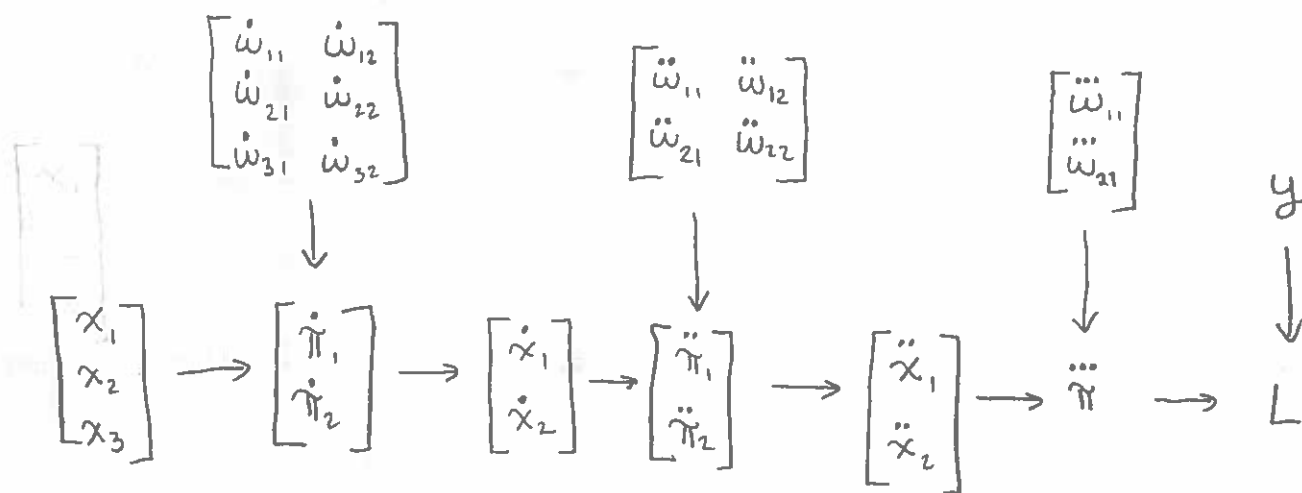Just in case we've forgotten which of these are vectors and which are matrices, here it is explicitly:

$$
\begin{bmatrix} \dot{w}_{11} & \cdots & \dot{w}_{1H} \\ \vdots & \ddots & \vdots \\ \dot{w}_{D1} & \cdots & \dot{w}_{DH} \end{bmatrix}
\qquad
\begin{bmatrix} \ddot{w}_{11} & \cdots & \ddot{w}_{1H} \\ \vdots & \ddots & \vdots \\ \ddot{w}_{H1} & \cdots & \ddot{w}_{HH} \end{bmatrix}
\qquad
\begin{bmatrix} \overset{M}{w}_{11} \\ \vdots \\ \overset{M}{w}_{H1} \end{bmatrix}
\qquad y
$$

$$
\begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} \rightarrow
\begin{bmatrix} \dot{\pi}_1 \\ \vdots \\ \dot{\pi}_H \end{bmatrix} \rightarrow
\begin{bmatrix} \dot{x}_1 \\ \vdots \\ \dot{x}_H \end{bmatrix} \rightarrow
\begin{bmatrix} \ddot{\pi}_1 \\ \vdots \\ \ddot{\pi}_H \end{bmatrix} \rightarrow
\begin{bmatrix} \ddot{x}_1 \\ \vdots \\ \ddot{x}_H \end{bmatrix} \rightarrow \cdots \rightarrow
\begin{bmatrix} \overset{M-1}{x}_1 \\ \vdots \\ \overset{M-1}{x}_H \end{bmatrix} \rightarrow \overset{M}{\pi} \rightarrow L
$$

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

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

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 $D=3$.

$$
\begin{bmatrix} \dot{w}_{11} & \dot{w}_{12} \\ \dot{w}_{21} & \dot{w}_{22} \\ \dot{w}_{31} & \dot{w}_{32} \end{bmatrix}
\qquad
\begin{bmatrix} \ddot{w}_{11} & \ddot{w}_{12} \\ \ddot{w}_{21} & \ddot{w}_{22} \end{bmatrix}
\qquad
\begin{bmatrix} \dddot{w}_{11} \\ \dddot{w}_{21} \end{bmatrix}
\qquad y
$$

$$
\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}
\longrightarrow
\begin{bmatrix} \dot{\pi}_1 \\ \dot{\pi}_2 \end{bmatrix}
\longrightarrow
\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix}
\longrightarrow
\begin{bmatrix} \ddot{\pi}_1 \\ \ddot{\pi}_2 \end{bmatrix}
\longrightarrow
\begin{bmatrix} \ddot{x}_1 \\ \ddot{x}_2 \end{bmatrix}
\longrightarrow
\dddot{\pi}
\longrightarrow
L
$$

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 $i, j$):

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

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

$$\frac{\partial L}{\partial \overset{m}{w}_{ij}} = \underbrace{\frac{\partial L}{\partial \dddot{\pi}}}_{} \cdot \frac{\partial \dddot{\pi}}{\partial \overset{m}{w}_{ij}}$$

This is the just the standard derivative of the loss function.

# Neural Networks and Backpropagation

⑦ So the challenge is to compute $\dfrac{\partial \dddot{\pi}}{\partial \overset{\tiny\textcircled{m}}{w}_{ij}}$ for any layer $m$.

It's straightforward for $m=3$:

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

⑧ What about $m=2$?

$$\frac{\partial \dddot{\pi}}{\partial \ddot{w}_{ij}} = \frac{\partial \dddot{\pi}}{\partial \ddot{\pi}_j} \cdot \frac{\partial \ddot{\pi}_j}{\partial \ddot{w}_{ij}}$$

$\left[ \ddot{\pi}_j \text{ separates } \dddot{\pi} \text{ from } \ddot{w}_{ij}, \text{ so Chain Rule applies} \right]$

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

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

⑨ What about $m=1$?

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

$\left[ \dot{\pi}_j \text{ separates } \dddot{\pi} \text{ from } \dot{w}_{ij}, \text{ so Chain Rule applies} \right]$

$$= \frac{\partial \dddot{\pi}}{\partial \dot{\pi}_j} \cdot x_i$$

$\left[ \dfrac{\partial \dot{\pi}_j}{\partial \dot{w}_{ij}} = x_i \right]$

NEURAL NETWORKS AND BACKPROPAGATION

⑩ In summary:

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

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

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

for the general case:

$$\frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{m}}{w}_{ij}} = \overset{\textcircled{m-1}}{x}_i \cdot \frac{\partial \overset{\textcircled{M}}{\pi}}{\partial \overset{\textcircled{m}}{\pi}_j}$$

so how do we compute this term?

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

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

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

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

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

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

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

but this can be computed recursively in the same way!

# NEURAL NETWORKS AND BACKPROPAGATION

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

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

and for the general case:

$$\frac{\partial \overset{(M)}{\pi}}{\partial \overset{(M)}{\pi}} = \boxed{\phantom{xxxxx}} \qquad \leftarrow \text{base case}$$

$$\frac{\partial \overset{(M)}{\pi}}{\partial \overset{(m)}{\pi}_j} = \boxed{\phantom{xxxxxxxxxxxx}} \qquad \leftarrow \text{recursive step}$$

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

BACKPROPAGATION:

   (a) for m in range$\left( \boxed{\phantom{xxx}} \right)$ and j in range$\left( \boxed{\phantom{xxx}} \right)$ :

      Compute $\dfrac{\partial \overset{(M)}{\pi}}{\partial \overset{(m)}{\pi}_j} = \boxed{\phantom{xxxxxx}}$

   (b) $\dfrac{\partial L}{\partial \overset{(m)}{w}_{ij}} = \boxed{\phantom{xxxx}} \cdot \dfrac{\partial \overset{(M)}{\pi}}{\partial \overset{(m)}{\pi}_j}$

Because we compute the partial derivatives $\dfrac{\partial \overset{(M)}{\pi}}{\partial \overset{(m)}{\pi}_j}$ starting from the final layer M and moving back, we call it backpropagation.