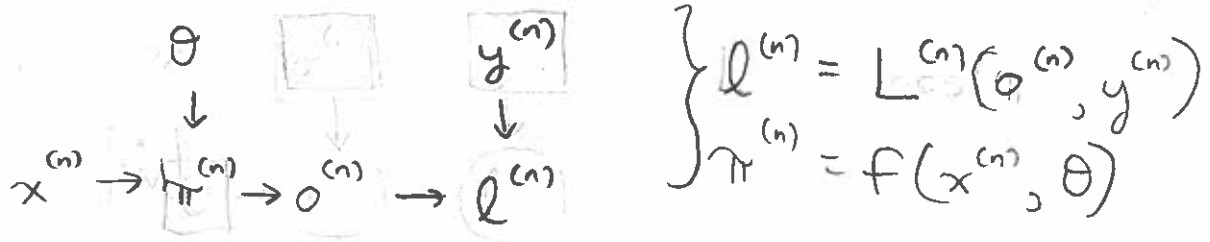


TRAINING FEATURE DISCOVERY NETWORKS

① So, given the causal model



we want to compute $\frac{\partial l^{(n)}}{\partial \theta}$.

By the Chain Rule: (since $\pi^{(n)}$ depends on θ from $L^{(n)}$):

$$\frac{\partial l^{(n)}}{\partial \theta} = \frac{\partial l^{(n)}}{\partial o^{(n)}} \frac{\partial o^{(n)}}{\partial \theta} = \frac{\partial l^{(n)}}{\partial o^{(n)}} \frac{\partial o^{(n)}}{\partial \pi^{(n)}} \frac{\partial \pi^{(n)}}{\partial \theta}$$

② This part should be straightforward to compute (assuming we chose some reasonable loss function).

e.g. $\frac{\partial}{\partial o^{(n)}} l^{(n)} = \frac{\partial}{\partial \pi^{(n)}} (y^{(n)} - o^{(n)})^2$

$$= -2(y^{(n)} - \pi^{(n)}) \quad \left(\text{for } l^{(n)} = L_{\text{lin}}^{(n)}(\pi^{(n)}, y^{(n)}) \right)$$

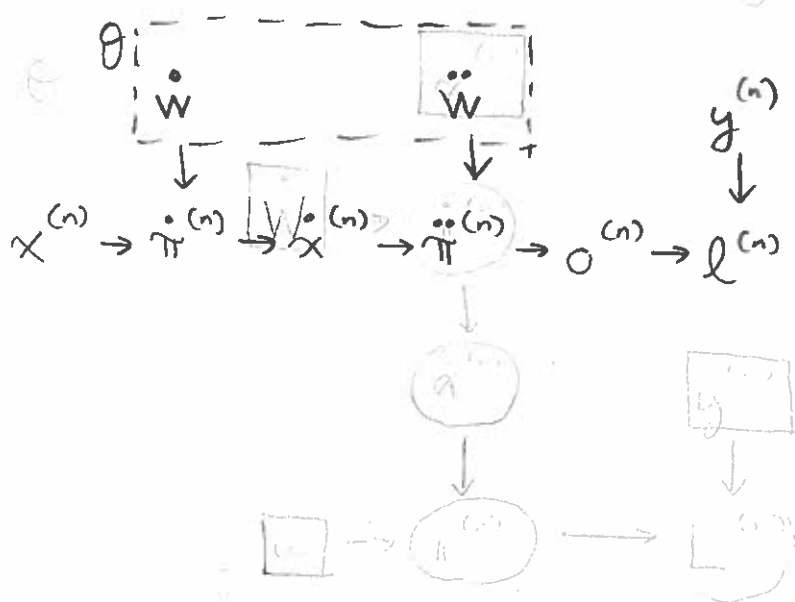
②b This part is also straightforward. For instance, with logistic regression, $o^{(n)} = \sigma(\pi^{(n)})$, thus:

$$\frac{\partial o^{(n)}}{\partial \pi^{(n)}} = \sigma'(\pi^{(n)}) = \sigma(\pi^{(n)}) (1 - \sigma(\pi^{(n)}))$$

TRAINING FEATURE DISCOVERY NETWORKS

③ The challenge lies in computing $\frac{\partial}{\partial \theta} \pi^{(n)} = \frac{\partial}{\partial \theta} f(x^{(n)}, \theta)$

Let's suppose it's our "feature discovery" network:

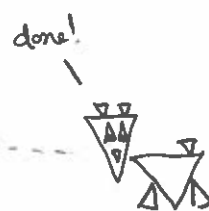


④ Recall that $\theta = \{\ddot{w}_{ij} \in \ddot{W}\} \cup \{\dot{w}_{ij} \in \dot{W}\}$,

so to compute $\frac{\partial}{\partial \theta} \ddot{\pi}^{(n)}$, we need to compute:

(a) for each $\ddot{w}_{ij} \in \ddot{W}$:

$$\frac{\partial}{\partial \ddot{w}_{ij}} \ddot{\pi}^{(n)} = \frac{\partial}{\partial \ddot{w}_{ij}} \dot{x}^{(n)} \ddot{W} = \dot{x}_i^{(n)}$$



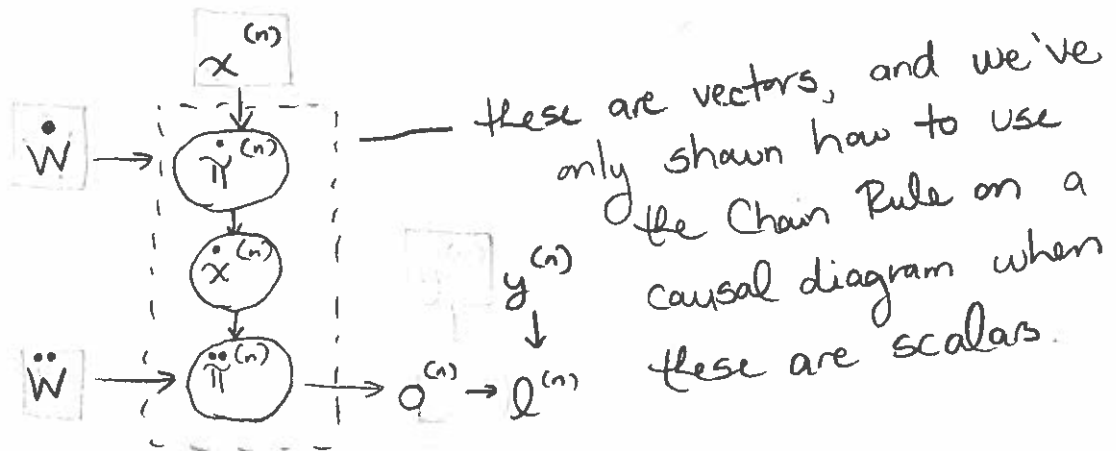
(b) for each $\dot{w}_{ij} \in \dot{W}$:

$$\frac{\partial}{\partial \dot{w}_{ij}} \ddot{\pi}^{(n)}$$

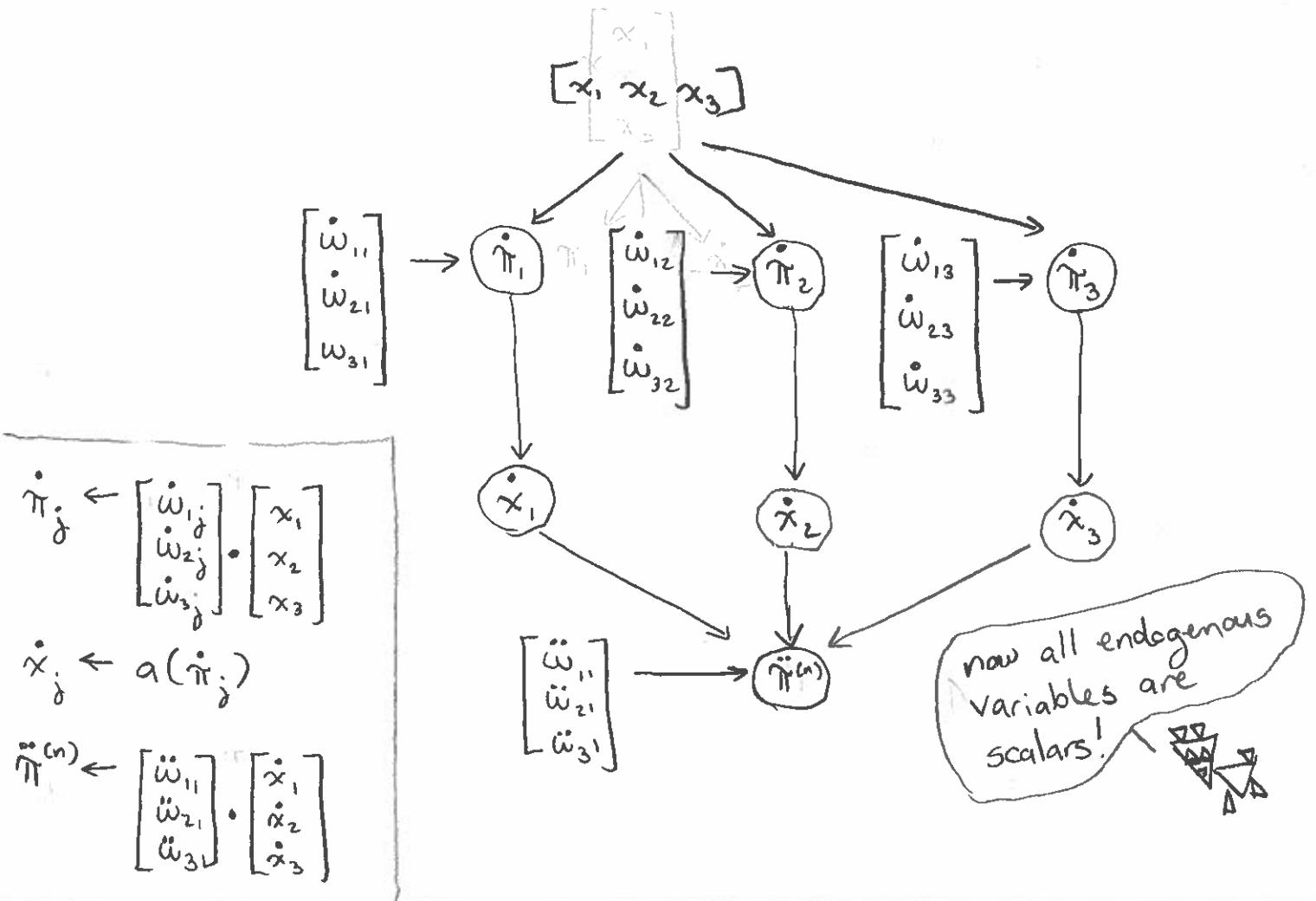


TRAINING FEATURE DISCOVERY NETWORKS

⑤ Sounds like a job for the Chain Rule of Partial Derivatives, but there's one issue:



⑥ No problem though, let's just explicitly represent the components of these vectors and matrices:



TRAINING FEATURE DISCOVERY NETWORKS

⑦ Now we can compute $\frac{\partial \ddot{\pi}^{(n)}}{\partial \dot{w}_{ij}}$ by repeated application

of the Chain Rule:

$$\frac{\partial \ddot{\pi}^{(n)}}{\partial \dot{w}_{ij}} = \frac{\partial \ddot{\pi}^{(n)}}{\partial \dot{x}_j} \cdot \frac{\partial \dot{x}_j}{\partial \dot{w}_{ij}} \quad \left[\dot{x}_j \text{ separates } \ddot{\pi}^{(n)} \text{ from } \dot{w}_{ij}, \text{ so Chain Rule applies} \right]$$

$$= \frac{\partial \ddot{\pi}^{(n)}}{\partial \dot{x}_j} \cdot \frac{\partial \dot{x}_j}{\partial \dot{\pi}_j} \cdot \frac{\partial \dot{\pi}_j}{\partial \dot{w}_{ij}} \quad \left[\dot{\pi}_j \text{ separates } \dot{x}_j \text{ from } \dot{w}_{ij}, \text{ so Chain Rule applies} \right]$$

$$= \left(\frac{\partial}{\partial \dot{x}_j} \ddot{x}^{(n)} \begin{bmatrix} \ddot{w}_{1j} \\ \ddot{w}_{2j} \\ \ddot{w}_{3j} \end{bmatrix} \right) \left(\frac{\partial}{\partial \dot{\pi}_j} a(\dot{\pi}_j) \right) \left(\frac{\partial}{\partial \dot{w}_{ij}} \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} \dot{w}_{1j} \\ \dot{w}_{2j} \\ \dot{w}_{3j} \end{bmatrix} \right)$$

$$= \ddot{w}_{ji} \cdot a'(\dot{\pi}_j) \cdot x_i$$

↑
this is just
the standard
derivative of
the ReLU
function.

TRAINING FEATURE DISCOVERY NETWORKS

⑧ Putting this together with the result from ①:

$$\frac{\partial \mathcal{L}^{(n)}}{\partial \dot{w}_{ij}} = \frac{\partial \mathcal{L}^{(n)}}{\partial o^{(n)}} \frac{\partial o^{(n)}}{\partial \ddot{\pi}^{(n)}} \cdot w_j \cdot a'(\ddot{\pi}_j^{(n)}) x_i^{(n)}$$

$$\text{and } \frac{\partial \mathcal{L}^{(n)}}{\partial \ddot{w}_{ij}} = \frac{\partial \mathcal{L}^{(n)}}{\partial o^{(n)}} \frac{\partial o^{(n)}}{\partial \ddot{\pi}^{(n)}} \dot{x}_i^{(n)}$$

So we can compute the partial derivative of the loss with respect to any of our weights. Thus we can use gradient descent to compute the weights that minimize our loss.

⑨ Let's generalize what we just did