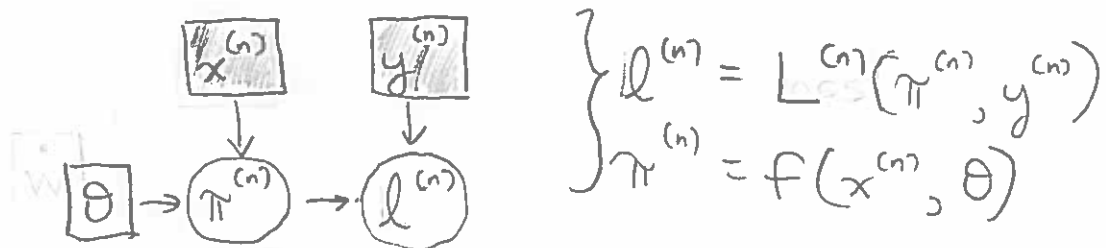


TRAINING FEATURE DISCOVERY NETWORKS

① So, given the causal model



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

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

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

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

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

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

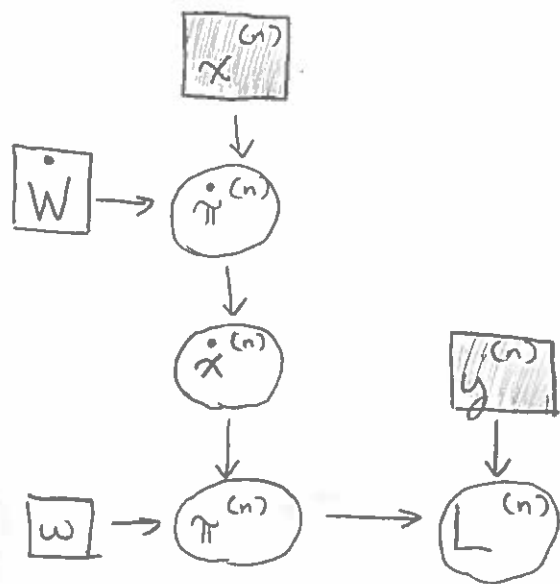
$$\frac{\partial}{\partial \pi^{(n)}} l^{(n)} = \frac{\partial}{\partial \pi^{(n)}} (1 - y^{(n)})\pi^{(n)} + \log(1 + e^{-\pi^{(n)}})$$

$$= (1 - y^{(n)}) - \frac{e^{\pi^{(n)}}}{1 + e^{\pi^{(n)}}} \quad \left(\text{for } l^{(n)} = L_{\text{logistic}}^{(n)}(\pi^{(n)}, y^{(n)}) \right)$$

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 $\theta = \{w_i \in w\} \cup \{\dot{w}_{ij} \in \dot{W}\}$,

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

(a) for each $w_i \in w$:

$$\frac{\partial}{\partial w_i} \pi^{(n)} = \frac{\partial}{\partial w_i} w^T \dot{x}^{(n)} = \dot{x}_i^{(n)}$$



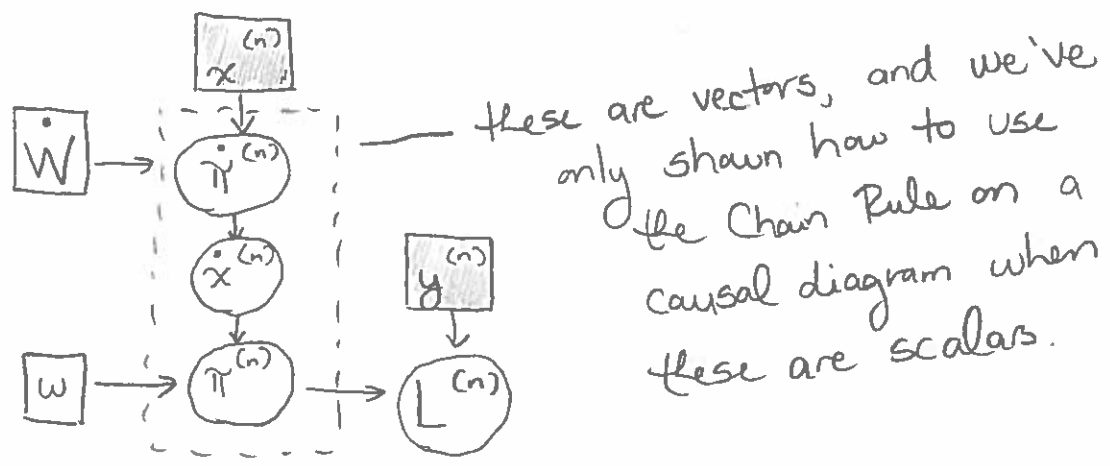
(b) for each $\dot{w}_{ij} \in \dot{W}$:

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

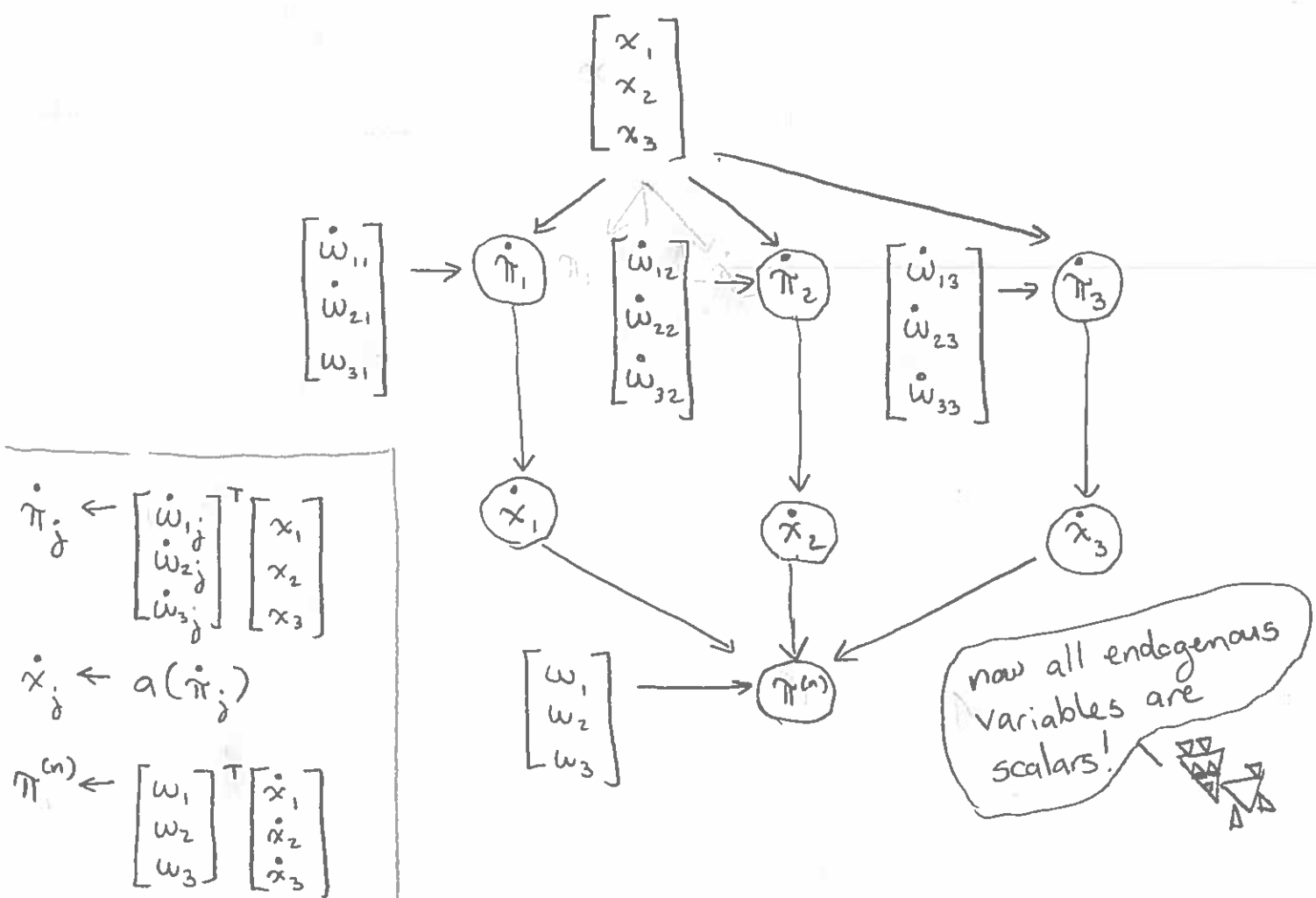


TRAINING FEATURE DISCOVERY NETWORKS

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 $\frac{\partial \pi^{(n)}}{\partial \dot{w}_{ij}}$ by repeated application

of the Chain Rule:

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

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

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

$$= w_j \cdot a'(\dot{\pi}_j) \cdot x_i$$

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

TRAINING FEATURE DISCOVERY NETWORKS

③ Putting this together with the result from ①:

$$\frac{\partial L^{(n)}}{\partial w_{ij}} = \left(\frac{\partial l^{(n)}}{\partial \pi^{(n)}} \right) w_j a'(\pi_j) x_i$$

and

$$\frac{\partial L^{(n)}}{\partial w_i} = \left(\frac{\partial l^{(n)}}{\partial \pi^{(n)}} \right) \dot{x}_i^{(n)}$$

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.

④ Let's generalize about our goal & 4.