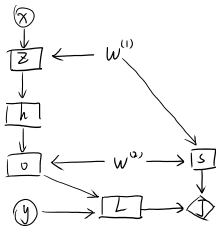


2018年7月31日 20:21

• 设 $Y=f(X)$, $Z=g(Y)$

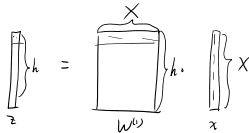
$$\frac{\partial Z}{\partial X} = \text{prod} \left(\frac{\partial Z}{\partial Y}, \frac{\partial Y}{\partial X} \right)$$

• 例: 正则化的多层感知机



FP:

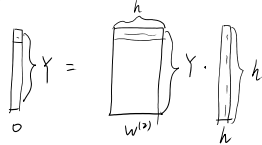
$$(1) \quad z = w^{(1)} x$$



(2) $h = f(z)$

$y(\cdot)$ element wise

$$(3) \quad 0 = W^{(2)} \cdot h$$

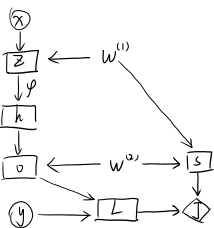


$$(4) \mathcal{L} = \ell(o, y)$$

scalar

$$(5) \quad S = \frac{\lambda}{2} (\|w^{(1)}\|_2^2 + \|w^{(2)}\|_2^2)$$

$$(b) J = L + S$$



BP:

$$c) \frac{\partial J}{\partial \lambda} = 1, \quad \frac{\partial J}{\partial \xi} = 1$$

$$\begin{aligned} (2) \quad \frac{\partial J}{\partial \theta} &= \text{prod} \left(\frac{\partial J}{\partial L} \cdot \frac{\partial L}{\partial \theta} \right) \\ &= 1 \cdot \frac{\partial \ell(\theta, y)}{\partial \theta} \\ &= \frac{\partial \ell}{\partial \theta} \end{aligned}$$

$$(3) \frac{\partial S}{\partial W^{(1)}} = \lambda W^{(1)}, \quad \frac{\partial S}{\partial W^{(2)}} = \lambda W^{(2)}$$

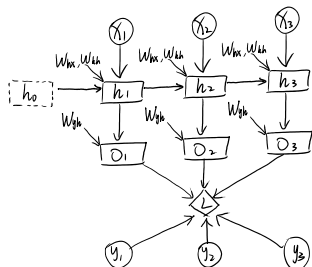
$$(4) \frac{\partial J}{\partial W^{(2)}} = \text{prod} \left(\frac{\partial J}{\partial o}, \frac{\partial o}{\partial W^{(2)}} \right) + \text{prod} \left(\frac{\partial J}{\partial s}, \frac{\partial s}{\partial W^{(2)}} \right)$$

$$\begin{aligned} \frac{\partial J}{\partial h} &= \text{prod} \left(\frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial h} \right) \\ &= W^{(2)T} \cdot \frac{\partial l}{\partial o} \end{aligned}$$

$$(b) \frac{\partial J}{\partial z} = \text{prod} \left(\frac{\partial J}{\partial h}, \frac{\partial h}{\partial z} \right) \\ = (W^{(2)T} \frac{\partial L}{\partial O}) \odot \varphi'(z)$$

$$(7) \frac{\partial J}{\partial W^{(1)}} = \text{prod} \left(\frac{\partial J}{\partial Z}, \frac{\partial Z}{\partial W^{(1)}} \right) + \text{prod} \left(\frac{\partial J}{\partial S}, \frac{\partial S}{\partial W^{(1)}} \right)$$

- RNN



FP:

FP:

$$(1) \quad h_t^{h \times 1} = \rho (W_{hx}^{h \times x} \cdot x_t^{x \times 1} + W_{hh}^{h \times h} \cdot h_{t-1}^{h \times 1}) \quad h_{t+1} = W_{hh} \cdot h_t$$

$$(2) \quad O_t = W_{y_h}^{y \times h} \cdot h_t^{h \times 1}$$

$$(3) \quad L = \frac{1}{T} \sum_{t=1}^T \ell(a_t, y_t)$$

BP:

$$(1) \frac{\partial L}{\partial \theta_1} = \frac{1}{T} \cdot \frac{\partial \ell(\theta_1, y_1)}{\partial \theta_1}$$

$$(2) \frac{\partial L^{y \times h}}{\partial W_{y_h}} = \sum_{t=1}^T \text{prod} \left(\frac{\partial L^{y \times h}}{\partial O_t}, \frac{\partial O_t}{\partial W_{y_h}} \right) \quad h_t: h \times 1$$

$$= \sum_{t=1}^T \frac{\partial L}{\partial \theta_t} h_t^T$$

$$(3) \frac{\partial L}{\partial h_T} = \text{Prod} \left(\frac{\partial L}{\partial O_T} \cdot \frac{\partial O_T}{\partial h_T} \right)$$

$$= W_{yh}^T \cdot \frac{\partial L}{\partial O_T}$$

$$(3) \frac{\partial L}{\partial h_T} = \text{Prod} \left(\frac{\partial O_T}{\partial T}, \frac{\partial h_T}{\partial T} \right)$$

$$= W_{yh}^T \cdot \frac{\partial L}{\partial O_T}$$

T 是时序的最大值, 如图中 $T=3$

$$\frac{\partial L}{\partial h_t} = \underbrace{\text{Prod} \left(\frac{\partial L}{\partial h_{t+1}}, \frac{\partial h_{t+1}}{\partial h_t} \right)}_{\substack{\text{error from } t+1 \\ \text{BP through time}}} + \underbrace{\text{Prod} \left(\frac{\partial L}{\partial O_t}, \frac{\partial O_t}{\partial h_t} \right)}_{\text{error from } t}$$

$$= W_{hh}^T \cdot \frac{\partial L}{\partial h_{t+1}} + W_{yh}^T \cdot \frac{\partial L}{\partial O_t} \quad (t < T)$$

calculate recursively from $\frac{\partial L}{\partial h_T}$ * gradient explosion

$$(4) \frac{\partial L}{\partial W_{hx}} = \sum_{t=1}^T \text{Prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hx}} \right) = \sum_{t=1}^T \left[\left(\frac{\partial L}{\partial h_t} \odot \psi'(\dots) \right) \cdot X_t^T \right]$$

$$\frac{\partial L}{\partial W_{hh}} = \sum_{t=1}^T \text{Prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hh}} \right) = \sum_{t=1}^T \left[\left(\frac{\partial L}{\partial h_t} \odot \psi'(\dots) \right) \cdot h_{t-1}^T \right]$$

$$\text{假设: } \frac{\partial L}{\partial h_t} = \sum_{i=t}^T (W_{hh}^T)^{T-i} \cdot W_{yh}^T \frac{\partial L}{\partial O_{T+i}}$$

• when $i \rightarrow 0$, $\lim_{i \rightarrow 0} (W_{hh}^T)^{T-i} = (W_{hh}^T)^T$ is big \Rightarrow gradient explosion