

$$\underline{a}^{(t)} = b + \underline{W} \underline{h}^{(t-1)} + \underline{U} \underline{x}^{(t)}$$

$$\underline{b}^{(t)} = \tanh(\underline{a}^{(t)})$$

$$\underline{O}^{(t)} = c + \underline{V} \underline{b}^{(t)}$$

$$\underline{L} = \sum_t \underline{L}^{(t)}$$

$$\nabla_{\underline{O}^{(t)}} \underline{L} = \underline{\hat{y}}^{(t)} - \underline{y}^{(t)}$$

$$* J_{ij} = \partial f_i / \partial x_j, f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f = A x, \partial f / \partial x = A^T$$

bias:

$$\nabla_{\underline{c}} \underline{L} = \sum_t \left(\frac{\partial \underline{O}^{(t)}}{\partial \underline{c}} \right)^T (\nabla_{\underline{O}^{(t)}} \underline{L})$$

$$= \sum_t \underline{\nabla_{\underline{c}} O}^{(t)} \cdot \underline{\nabla_{\underline{O}^{(t)}} L}^{(t)}$$

$$= \sum_t \underline{\nabla_{\underline{O}^{(t)}} L}^{(t)} = \sum_t \underline{\nabla_{\underline{O}^{(t)}} L} *$$

$$\nabla_{\underline{h}^{(t)}} \underline{L} = \left(\frac{\partial \underline{h}^{(t+1)}}{\partial \underline{h}^{(t)}} \right)^T (\nabla_{\underline{h}^{(t+1)}} \underline{L}) + \left(\frac{\partial \underline{O}^{(t)}}{\partial \underline{h}^{(t)}} \right)^T (\nabla_{\underline{O}^{(t)}} \underline{L})$$

$$= \underline{W}^T \underline{H}^{(t+1)} (\nabla_{\underline{h}^{(t+1)}} \underline{L}) + \underline{V}^T (\nabla_{\underline{O}^{(t)}} \underline{L})$$

* $\underline{h}^{(t+1)}$ 依賴 $\underline{h}^{(t)}$ 及 $\underline{O}^{(t)}$

where

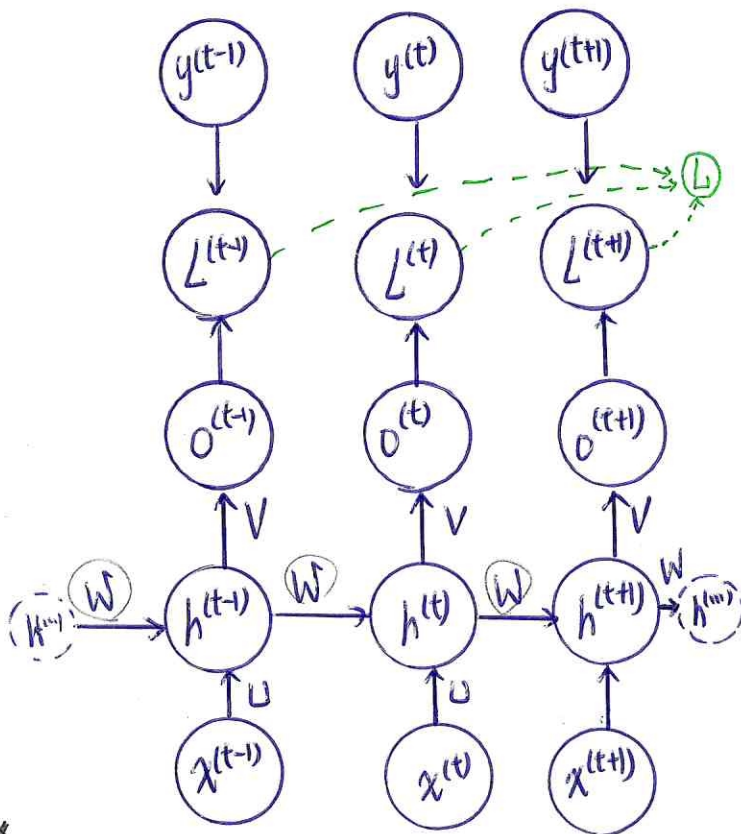
$$\underline{H}^{(t+1)} = \left(\frac{\partial \underline{h}^{(t+1)}}{\partial \underline{a}^{(t+1)}} \right)^T$$

$$= \begin{bmatrix} 1 - (\underline{h}_1^{(t+1)})^2 & 0 & \dots & 0 \\ 0 & 1 - (\underline{h}_2^{(t+1)})^2 & & \vdots \\ \vdots & & \ddots & \\ 0 & 0 & \dots & 1 - (\underline{h}_n^{(t+1)})^2 \end{bmatrix}$$

$$\frac{d \tanh(x)}{dx} = 1 - \tanh^2(x)$$

$$\nabla_{\underline{b}} \underline{L} = \sum_t \left(\frac{\partial \underline{h}^{(t)}}{\partial \underline{b}} \right)^T \nabla_{\underline{h}^{(t)}} \underline{L} = \sum_t \text{diag}(1 - (\underline{h}^{(t)})^2) (\nabla_{\underline{h}^{(t)}} \underline{L}) = \sum_t \underline{H}^{(t)} \cdot \nabla_{\underline{h}^{(t)}} \underline{L} *$$

$$\sum_t \sum_i \left(\frac{\partial \underline{L}}{\partial \underline{h}_i^{(t)}} \right) (\nabla_{\underline{b}} \underline{h}_i^{(t)})$$



$$\nabla_{\underline{v}} L = \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) (\nabla_{\underline{v}} h_i^{(t)}) = \sum_t \left(\frac{\partial \sigma^{(t)}}{\partial v} \right) (\nabla_{\underline{0}^{(t)}} L) = \sum_t (\nabla_{\underline{0}^{(t)}} L) (h^{(t)})^T$$

$$\begin{aligned} \nabla_{\underline{w}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) (\nabla_{\underline{w}} h_i^{(t)}) = \sum_t \left(\frac{\partial L}{\partial h^{(t)}} \right) \left(\frac{\partial h^{(t)}}{\partial \underline{w}} \right) = \sum_t \text{diag}(1 - (h^{(t)})^2) (\nabla_{\underline{h}^{(t)}} L) (h^{(t-1)})^T \\ &= \sum_t H^{(t)} (\nabla_{\underline{h}^{(t)}} L) (h^{(t-1)})^T \end{aligned}$$

$$\begin{aligned} \nabla_{\underline{u}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) (\nabla_{\underline{u}} h_i^{(t)}) = \sum_t \left(\frac{\partial L}{\partial h^{(t)}} \right) \left(\frac{\partial h^{(t)}}{\partial \underline{u}} \right) \\ &= \sum_t \text{diag}(1 - (h^{(t)})^2) (\nabla_{\underline{h}^{(t)}} L) (x^{(t)})^T = \sum_t H^{(t)} (\nabla_{\underline{h}^{(t)}} L) (x^{(t)})^T \end{aligned}$$