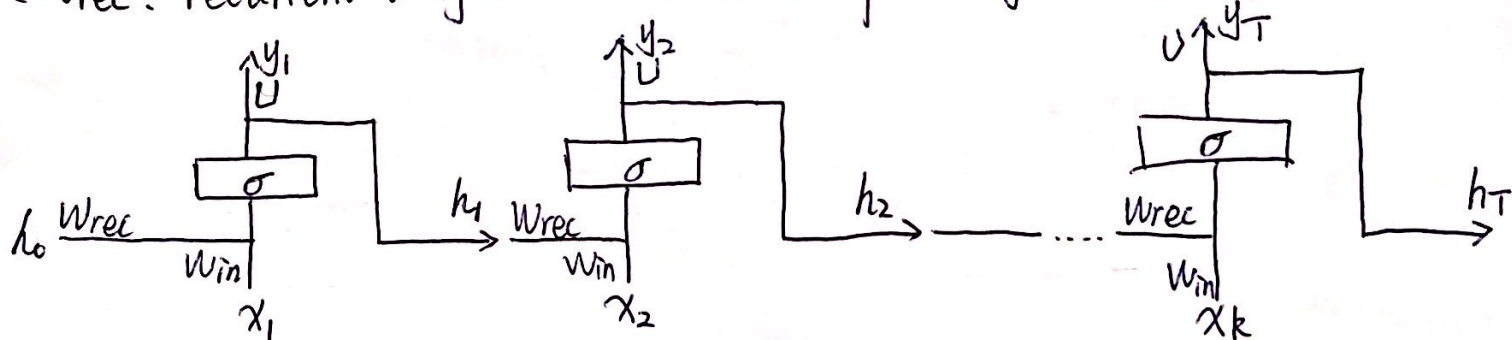


Recurrent Neural Network (RNN)

RNN基本结构

(W_{rec} : recurrent weight matrix, W_{in} : input weight matrix)



$$h_t = \sigma(W_{rec}h_{t-1} + W_{in}x_t + b_t)$$

Back-prop ~~in~~ through time (BPTT). let $W = [W_{rec}, W_{in}]$

$$\frac{\partial E}{\partial W} = \bar{z}_{t-1}^T \frac{\partial E_t}{\partial W} \quad (t: \text{第 } t \text{ 个 timestep})$$

$$\frac{\partial E_k}{\partial W} = \frac{\partial E_k}{\partial y_k} \cdot \frac{\partial y_k}{\partial h_k} \left(\prod_{t=2}^k \frac{\partial h_t}{\partial h_{t-1}} \right) \cdot \frac{\partial h_1}{\partial W}$$

其中: $\frac{\partial h_t}{\partial h_{t-1}} = \text{diag}(\sigma'(W_{rec}h_{t-1} + W_{in}x_t + b_t)) \cdot W_{rec}$

$$\text{so: } \prod_{t=2}^k \frac{\partial h_t}{\partial h_{t-1}} = \prod_{t=2}^k (\text{diag}(\sigma'(W_{rec}h_{t-1} + W_{in}x_t + b_t)) \cdot W_{rec}) \quad \textcircled{1}$$

σ 为 sigmoid 时, σ' 最大为 $\frac{1}{4}$; σ 为 tanh 或 ReLU, 最大为 1. 当网络较深时, 即 k 较大时, ①式在连乘下会变得极小导致梯度消失; 而当 W_{rec} 过大时, ①式又会

导致梯度爆炸.