

$$(x_1, x_2, \dots, x_T) \quad x_t \in \mathbb{R}^d$$

sequence classification / regression

$$f(x) = R \sigma(Wx_T + U \sigma(Wx_{T-1} + U \sigma(\dots \sigma(Wx_1 + U h_0 + b) \dots) + b) + b) + c$$

$$f(x) = Rh_T + c, \text{ where}$$

$$h_t = \sigma(Wx_t + Uh_{t-1} + b)$$

$$l(f(x), y) = \|f(x) - y\|^2$$

$$\frac{\partial l}{\partial R} = \frac{\partial l}{\partial f} \cdot \frac{\partial f}{\partial R}, \quad \frac{\partial l}{\partial c} = \frac{\partial l}{\partial f} \cdot \frac{\partial f}{\partial c}$$

$$f(x) = R \sigma([W]^T x_T + [U]^T \sigma([W]^T x_{T-1} + [U]^T \sigma(\dots \sigma([W]^T x_1 + [U]^T h_0 + [b]) \dots) + [b]) + [b]) + c$$

$$\frac{\partial l}{\partial u} = \frac{\partial l}{\partial f} \left(\sum_{t=1}^T \frac{\partial f}{\partial h_t} \cdot \frac{\partial h_t}{\partial [u]^T} \right)$$

$$\frac{\partial l}{\partial W} = \frac{\partial l}{\partial f} \left(\sum_{t=1}^T \frac{\partial f}{\partial h_t} \cdot \frac{\partial h_t}{\partial [W]^T} \right)$$

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial f} \left(\sum_{t=1}^T \frac{\partial f}{\partial h_t} \cdot \frac{\partial h_t}{\partial [b]^T} \right)$$

$$\frac{\partial h_T}{\partial h_{T-1}}(h_{T-1})$$

$$\frac{\partial f}{\partial h_t} = \frac{\partial f}{\partial h_T} \cdot \frac{\partial h_T}{\partial h_{t-1}} \cdot \dots \cdot \frac{\partial h_{t+1}}{\partial h_t}$$

temporal derivative

$$\frac{\partial h_t}{\partial h_{t-1}} = \text{diag}(\sigma'(a_t)) U$$

$$h_t = \sigma(Wx_t + Uh_{t-1} + b)$$

$$\sigma'(a) = \begin{cases} 0 & \text{if } a \leq 0 \\ 1 & \text{o.w.} \end{cases}, \text{ where } \sigma(a) = \max(0, a)$$

$$\sigma: \mathbb{R}^{d'} \rightarrow (-1, 1)^{d'} \quad \underline{\sigma(a) = \tanh(a)} \quad \checkmark$$

$$0 < \underline{\tanh'(a)} \leq \underline{1}.$$

$$\left\| \frac{\partial h_t}{\partial h_{t-1}} \right\| = \left\| \text{diag}(\sigma'(a_t)) U \right\| \leq \left\| \text{diag}(\sigma'(a_t)) \right\| \|U\|$$

$$\left\| \text{diag}(\sigma'(a_t)) \right\| \|U\| = \max \sigma'(a_t) \|U\| \leq \|U\|$$

$$\left\| \frac{\partial f}{\partial h_t} \right\| \leq \left\| \frac{\partial f}{\partial h_T} \right\| \left\| \frac{\partial h_T}{\partial h_{T-1}} \right\| \cdots \left\| \frac{\partial h_{t+1}}{\partial h_t} \right\| \leq \left\| \frac{\partial f}{\partial h_T} \right\| \cdot \|U\|^{T-t}$$

• exploding gradient may happen when $\|U\| \gg 1$.

• Vanishing gradient when $\|U\| < 1$.

$$\left\| \frac{\partial \mathcal{L}}{\partial [u]^t} \right\| \rightarrow 0 \quad \dots \quad \left\| \frac{\partial \mathcal{L}}{\partial u} \right\| \rightarrow 0$$

Linear shortcuts.

$$h_t = \sigma(U h_{t-1} + W x_t + b) + \sum_{t'=1}^{t-1} g_t(h_{t'})$$

$$h_t = \sigma(U h_{t-1} + W x_t + b) + \sum_{t'=1}^{t-1} g_t([h_{t'}]^t)$$

$$\frac{\partial h_t}{\partial h_{t'}} = \frac{1}{1} \frac{\partial h_{t-t'+1}}{\partial h_{t-t'}} + \frac{\partial h_t}{\partial g_t([h_{t'}]^t)} \cdot \frac{\partial g_t([h_{t'}]^t)}{\partial [h_{t'}]^t}$$

$$t' \ll t$$

① Residual connection



$$g_{t+1}(a) = I_a$$

$$g_{t'}(a) = 0 \quad \text{for } t' < t-1$$

$$h_t = \sigma(Uh^{t-1} + Wx^t + b) + \underbrace{h_{t-1}}_{= \sigma(Uh_{t-2} + Wx_{t-1} + b) + \underbrace{h_{t-2}}_{\vdots}}$$

$$h_t = \sigma(Uh^{t-1} + Wx^t + b) + \sigma(Uh_{t-2} + Wx_{t-1} + b) + h_{t-2}$$

$$= \sigma(Uh_{t-1} + Wx_t + b) + \sigma(Uh_{t-2} + Wx_{t-1} + b) + \sigma(Uh_{t-3} + Wx_{t-2} + b) + h_{t-3}$$

$$= \sum_{t' \leq t} \sigma(Uh_{t'-1} + Wx_{t'} + b)$$

$$\|h_t\| \leq \sum_{t' \leq t} \|\sigma(Uh_{t'-1} + Wx_{t'} + b)\| \quad \text{formal explosion}$$

② GRU

element-wise
multiplication

$$h_t = u_t \odot \sigma_h(Uh^{t-1} + Wx^t + b) + (1 - u_t) \odot h_{t-1}$$

$$u_t = \sigma_u(U_u h_{t-1} + W_u x_t + b_u) \in [0, 1]^d$$

update gate

$$h_t = u_t \odot \sigma_h(U(r_t \odot h^{t-1}) + Wx^t + b) + (1 - u_t) \odot h_{t-1}$$

$$r_t = \sigma_r(U_r h_{t-1} + W_r x_t + b_r) \in [0, 1]^d$$