

RNN

1. 递归模型

变量 y_t 的历史信息来预测自己

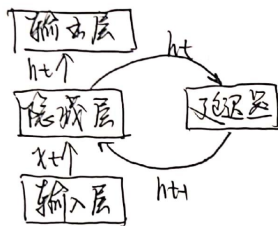
$$y_t = w_0 + \sum_{i=1}^p w_i y_{t-i} + \epsilon_t$$

有外部输入的非线性递归

$h_t =$

2. 循环 RNN

$$h_t = f(h_{t-1}, x_t) \quad h_0 = 0$$



可逆系统 (非线性动力系统)

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

3. BPTT

$$L = \sum_{t=1}^T L_t$$

$$\frac{\partial L}{\partial U} = \sum_{t=1}^T \frac{\partial L_t}{\partial U}$$

$$\frac{\partial L_t}{\partial U_{ij}} = \sum_{k=1}^t \frac{\partial z_k}{\partial U_{ij}} \frac{\partial L_t}{\partial z_k}$$

$$z_k = Uh_{k-1} + Wx_k + b$$

$$\frac{\partial z_k}{\partial U_{ij}} = I_i([h_{k-1}]_j)$$

$$\delta_{t,k} = \frac{\partial L_t}{\partial z_k} = \frac{\partial L_t}{\partial z_k} \frac{\partial z_k}{\partial h_{k-1}} \frac{\partial h_{k-1}}{\partial U} = \text{diag}(f'(z_k)) U^T \delta_{t,k+1}$$

第t时刻的损失

对第k时刻隐藏

神经层输入 z_k

的导数

$$\frac{\partial L_t}{\partial U_{ij}} = \sum_{k=1}^t [\delta_{t,k}]_i [h_{k-1}]_j$$

$$\frac{\partial L_t}{\partial U} = \sum_{k=1}^t \delta_{t,k} h_{k-1}^T$$

$$\text{网络} \frac{\partial L}{\partial U} = \sum_{t=1}^T \sum_{k=1}^t \delta_{t,k} h_{k-1}^T$$

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \sum_{k=1}^t \delta_{t,k} x_k^T$$

$$\frac{\partial L}{\partial b} = \sum_{t=1}^T \sum_{k=1}^t \delta_{t,k}$$

4. 长期依赖问题

难以建模长时间间隔之间的关系

$$\delta_{t,k} = \frac{\partial L_t}{\partial z_k} = \frac{\partial L_t}{\partial z_k} \prod_{\tau=k}^{t-1} (\text{diag}(f'(z_\tau)) U^T) \delta_{t,t}$$

$$\text{令 } \gamma \approx \|\text{diag}(f'(z_k)) U^T\|, \text{ 则}$$

$$\delta_{t,k} \approx \gamma^{t-k} \delta_{t,t}$$

$$\text{若 } t \rightarrow \infty, \begin{cases} \gamma > 1, & \delta_{t,k} \rightarrow \infty \\ \gamma < 1, & \delta_{t,k} \rightarrow 0 \end{cases} \quad \checkmark$$

注意 RNN 中说的梯度消失并非指

$$\frac{\partial L_t}{\partial U} \text{ 消失了, } \because \frac{\partial L_t}{\partial U} = \sum_{k=1}^t \delta_{t,k} h_{k-1}^T$$

$$\text{而 } \delta_{t,k} \approx \gamma^{t-k} \delta_{t,t}$$

\therefore 指的是那些离 t 很远的 k , 也就是

更新主要依靠当前时刻 t 的 U 相邻

状态 h_k , 长距离的状态对 U 没有影响

$$\text{基础改进: } h_t = h_{t-1} + g(x_t, h_{t-1}, 0)$$

1) 梯度爆炸

2) 记忆容量

5. 基于门控的 RNN (LSTM)

3 门: 遗忘门: $f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$

输入门: $i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$

输出门: $o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

$$\tilde{c}_t = i_t \odot \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

梯度问题

$$h_t = o_t \odot \tanh(c_t)$$

注: 一般 RNN 网络参数初始值比较小, 但 LSTM 参数初值会使梯度消失, 导致大量信息丢失, 所以一般初始值比较大

6. 门控循环神经网络 (GRU)

$$h_t = z_t \odot h_{t+1} + (1 - z_t) \odot g(x_t, h_{t+1}; \theta)$$

$$z_t = \sigma(W_z x_t + U_z h_{t+1} + b_z)$$

$$\tilde{h}_t = g(x_t, h_{t+1}; \theta)$$

$$= \tanh(W_h x_t + U_h (r_t \odot h_{t+1}) + b_h)$$

r_t 为重置门，是否忽略 h_{t+1}

$$r_t = \sigma(W_r x_t + U_r h_{t+1} + b_r)$$