

$$\text{输入: } a_t^{(\text{in})} = \sigma(W^{(\text{in})} x_t + b^{(\text{in})})$$

$$\text{隐层: } h_t = \sigma(U a_t^{(\text{in})} + V h_{t-1} + b_h)$$

$$o_t = \sigma(W h_t + b_o)$$

$$\text{输出: } h_t^{(\text{out})} = W^{(\text{out})} o_t + b^{(\text{out})}$$

$$\hat{r}_t = \text{Softmax}(h_t^{(\text{out})})$$

$$\text{Loss: } L = \sum_{t=1}^T L_t = \sum_{t=1}^T \text{Loss}(r_t, \hat{r}_t)$$

$$\text{使用交叉熵: } L = \sum_{t=1}^T -r_t \log \hat{r}_t$$

$$\frac{\partial \sigma(z)}{\partial z} = -\frac{-e^{-z}}{(1+e^{-z})^2} = \frac{e^{-z}}{(1+e^{-z})^2} = \sigma(z) \cdot (1 - \sigma(z))$$

$$\text{对于 SoftMax: 令 } \hat{z} = [x_1, \dots, x_c] \text{ 则 } \text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^c e^{x_j}}$$

$$\frac{\partial \text{Softmax}(x_i)}{\partial x_j} = -\frac{e^{x_j}}{(\sum_{j=1}^c e^{x_j})^2} \cdot e^{x_i} \quad (i \neq j)$$

$$\frac{\partial \text{Softmax}(x_i)}{\partial x_j} = \frac{e^{x_i} \cdot (\sum_{j=1}^c e^{x_j} - e^{x_j})}{(\sum_{j=1}^c e^{x_j})^2} \quad (i = j)$$

$$\therefore \frac{\partial L}{\partial h_{t,i}^{(\text{out})}} = \sum_{t=1}^T \sum_{a=1}^c \frac{\partial L_t}{\partial \hat{r}_{t,a}} \cdot \frac{\partial \hat{r}_{t,a}}{\partial h_{t,i}^{(\text{out})}} \quad \left| \quad \frac{\partial L_t}{\partial \hat{r}_{t,a}} = -\frac{r_t}{\hat{r}_t} \right.$$

$$\therefore L_t = -\sum_{i=1}^c \left[\frac{r_{t,i}}{\hat{r}_{t,i}} \hat{r}_{t,i} (1 - \hat{r}_{t,i}) - \sum_{j \neq i} \frac{r_{t,j}}{\hat{r}_{t,j}} \hat{r}_{t,i} \hat{r}_{t,j} \right]$$

$$= -\sum_{i=1}^c \left[r_{t,i} (1 - \hat{r}_{t,i}) - \sum_{j \neq i} r_{t,j} \hat{r}_{t,i} \right]$$

$$= -\sum_{i=1}^c [r_{t,i} - \hat{r}_{t,i} \sum_{j=1}^c r_{t,j}] = \sum_{i=1}^c (\hat{r}_{t,i} - r_{t,i})$$

根据 $h_t^{(out)} = W^{(out)} o_t + b^{(out)}$

$$\frac{\partial L}{\partial b_i^{(out)}} = \frac{\partial L}{\partial h_t^{(out)}} \cdot \frac{\partial h_t^{(out)}}{\partial b_i^{(out)}} = \sum_{t=1}^{Tr} (\hat{r}_t - r_t)$$

$$\frac{\partial L}{\partial W^{(out)}} = \frac{\partial L}{\partial h_t^{(out)}} \cdot \frac{\partial h_t^{(out)}}{\partial W^{(out)}} = \sum_{t=1}^{Tr} (\hat{r}_t - r_t) o_t^T$$

$$\frac{\partial L}{\partial o_t} = \frac{\partial L}{\partial h_t^{(out)}} \cdot \frac{\partial h_t^{(out)}}{\partial o_t} = \sum_{t=1}^{Tr} (\hat{r}_t - r_t) \cdot W^{(out)T}$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial o_t} \cdot \frac{\partial o_t}{\partial W} = \sum_{t=1}^{Tr} (\hat{r}_t - r_t) W^{(out)T} \cdot o_t (1 - o_t) \cdot h_t^T$$

$$\frac{\partial L}{\partial b_o} = \sum_{t=1}^{Tr} (\hat{r}_t - r_t) W^{(out)T} o_t (1 - o_t)$$

对 h_t 的梯度包含 h_{t+1} 和 o_t 而当 $t = Tr$ 时只有 o_t

$$\begin{aligned} \text{故 } t \neq Tr \text{ 时: } \frac{\partial L_t}{\partial h_t} &= \left(\frac{\partial o_t}{\partial h_t} \right)^T \frac{\partial L_t}{\partial o_t} + \left(\frac{\partial h_{t+1}}{\partial h_t} \right)^T \frac{\partial L_{t+1}}{\partial h_{t+1}} \\ &= o_t (1 - o_t) W^{(out)T} (\hat{r}_t - r_t) \cdot W^{(out)T} + h_{t+1} (1 - h_{t+1}) \cdot V^T \cdot \frac{\partial L_{t+1}}{\partial h_{t+1}} \end{aligned}$$

$t = Tr$: 上式成为前半部分

$$\therefore \frac{\partial L}{\partial U} = \sum_{t=1}^{Tr} h_t (1 - h_t) \frac{\partial L_t}{\partial h_t} \cdot a_t^{(in)T} \quad \left(\frac{\partial L_t}{\partial h_t} \text{ 见 D.B} \right)$$

$$\frac{\partial L}{\partial V} = \sum_{t=1}^{Tr} h_t (1 - h_t) \frac{\partial L_t}{\partial h_t} \cdot h_{t+1}^T \quad \frac{\partial L}{\partial W^{(in)}} = \frac{\partial L}{\partial a_t^{(in)}} \cdot \frac{\partial a_t^{(in)}}{\partial W^{(in)}}$$

$$\frac{\partial L}{\partial b_h} = \sum_{t=1}^{Tr} h_t (1 - h_t) \frac{\partial L_t}{\partial h_t} \quad = \sum_{t=1}^{Tr} h_t (1 - h_t) \frac{\partial L_t}{\partial h_t} U^T \cdot x_t^T$$

$$\frac{\partial L}{\partial a_t^{(in)}} = \sum_{t=1}^{Tr} h_t (1 - h_t) \frac{\partial L_t}{\partial h_t} U^T \quad \frac{\partial L}{\partial b^{(in)}} = \sum_{t=1}^{Tr} h_t (1 - h_t) \frac{\partial L_t}{\partial h_t} U^T$$