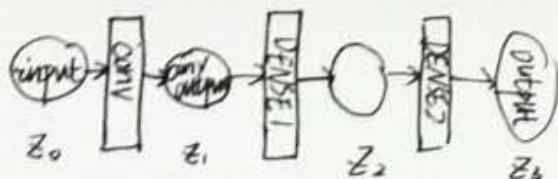


1.



令经过conv/dense运算, 没经过非线性的数据为 y_i , 经过非线性之后为 z_i ($i=1, \dots, 3$)

$$\text{loss } L = L(z_3) \quad z_i = f(y_i) \quad i=1, 2, 3$$

$$y_i = g_i(z_{i-1}) \quad i=1, \dots, 3 \quad g_1, \dots, g_3 \text{ 分别记为 } w_1, \dots, w_3$$

其中 g_1 为 conv 操作, g_2, g_3 为 dense 操作.

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_3} \quad \frac{\partial L}{\partial y_3} = \frac{\partial L}{\partial z_3} f'(y_3) \quad y_3 = w_3^T z_2$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial z_3} \frac{\partial z_3}{\partial y_3} \frac{\partial y_3}{\partial w_3} = \frac{\partial L}{\partial z_3} f'(y_3) z_2$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial z_3} f'(y_3) w_3 \quad y_2 = w_2^T z_1$$

$$\frac{\partial L}{\partial y_2} = \frac{\partial L}{\partial z_3} f'(y_3) w_3 f'(y_2) \quad \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y_2} \frac{\partial y_2}{\partial w_2} = \frac{\partial L}{\partial z_3} f'(y_3) w_3 f'(y_2) z_1$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial y_2} \frac{\partial y_2}{\partial z_1} = \frac{\partial L}{\partial z_3} f'(y_3) w_3 f'(y_2) w_2$$

$$\frac{\partial L}{\partial y_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial y_1} = \frac{\partial L}{\partial z_3} f'(y_3) w_3 f'(y_2) w_2 f'(y_1)$$

$$y_1 = \text{conv}(w_1, z_0) \quad w_1 = \begin{bmatrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{bmatrix} \{w_{11}, \dots, w_{19}, b_1\}$$

$$\frac{\partial L}{\partial w_{ij}} = \sum_{k=1}^4 \frac{\partial L}{\partial y_{ik}} \frac{\partial y_{ik}}{\partial w_{ij}}$$

$$= \sum_{k=1}^4 \frac{\partial L}{\partial y_{ik}} \frac{\partial y_{ik}}{\partial w_{ij}}$$

比如对于 $j=1$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y_{11}} z_{0,11}$$

$$+ \frac{\partial L}{\partial y_{12}} z_{0,12}$$

$$+ \frac{\partial L}{\partial y_{13}} z_{0,21} + \frac{\partial L}{\partial y_{14}} z_{0,22}$$

$$\therefore \frac{\partial L}{\partial w_{11}} = \text{conv}(\frac{\partial L}{\partial y_1}, z_0)$$

$$\text{而 } \frac{\partial L}{\partial b_1} = 1 \cdot \frac{\partial L}{\partial y_1}$$

$$\therefore K_{t+1} = K_t - \eta \frac{\partial L}{\partial w_{11}} = K_t - \eta \text{conv}(\frac{\partial L}{\partial y_1}, z_0)$$

$$B_{t+1} = B_t - \eta \frac{\partial L}{\partial b_1} = B_t - \eta \frac{\partial L}{\partial y_1}$$

2. 1) 序列条件分布 cross entropy 为

序列长度为 x

$$L(P_{\text{true}}, P_{\text{model}}) = - \sum_y P_{\text{true}}(y) \log P_{\text{model}}(y|x)$$

cross entropy loss 中 P_{true} 只在训练样本 y 处为 1 其它为 0. \therefore 设训练集为 $\{x^{(i)}, y^{(i)}\}_{i=1, \dots, N}$
 一共 N 组训练数据

$$L = - \sum_{i=1}^N \log P_{\text{model}}(y^{(i)} | x^{(i)})$$

每个 $y^{(i)}$ 长度为 $T=3$

包含 $y_1^{(i)}, y_2^{(i)}, y_3^{(i)}$

$$\text{其中 } P_{\text{model}}(y^{(i)} | x^{(i)}) = P(y_1^{(i)} | x_1^{(i)}) P(y_2^{(i)} | x_2^{(i)}, h_1^{(i)})$$

$$P(y_3^{(i)} | x_3^{(i)}, h_2^{(i)})$$

$$= \gamma_1^T (x_1^{(i)})^T y_1^{(i)} \gamma_2^T (x_2^{(i)})^T y_2^{(i)} \gamma_3^T (x_3^{(i)})^T y_3^{(i)}$$

其中 $\gamma_1(x), \gamma_2(x), \gamma_3(x)$ 为将 $x^{(i)}$ 输入 RNN 得到的每一个时间 $t=1, \dots, 3$ 得到的状态向量
 $y_1^{(i)}, y_2^{(i)}, y_3^{(i)}$ 讲道理应该是一个 $T=3$ 维指示向量。

2) 为表示方便, 样本仍用 $x^{(i)}, y^{(i)}$ 表示, 模型得到的用大写字 Y, Y, Y_2 表示。

① 对于任意参数 $w \in \{W, U, V\}$ 都有

$$\begin{aligned} \Delta w &= - \text{lr} \frac{\partial L}{\partial w} = \text{lr} \frac{\partial \log P_{\text{model}}(y^{(i)} | x^{(i)})}{\partial w} = \text{lr} \frac{\frac{\partial P_{\text{model}}(y^{(i)} | x^{(i)})}{\partial w}}{P_{\text{model}}(y^{(i)} | x^{(i)})} \\ &= \text{lr} \frac{1}{\prod_{j=1}^T \gamma_j^T (x^{(i)})^T y_j^{(i)}} \frac{\partial P_{\text{model}}(y^{(i)} | x^{(i)})}{\partial w} \end{aligned}$$

① 对于 V : $\frac{\partial P_{model}}{\partial V} = \sum_{t=1}^T \frac{\partial P_{model}}{\partial Y_t} \frac{\partial Y_t}{\partial V}$

$\frac{\partial Y_t}{\partial V} = \text{softmax}' \Big|_{res=Y_t} H_j^T$ $\text{softmax}' \Big|_{res=Y_t}$ 第 i, j 项为 $\begin{cases} i=j & Y_{ti}(1-Y_{ti}) \\ i \neq j & -Y_{ti}Y_{tj} \end{cases}$ (*)

$\frac{\partial P_{model}}{\partial Y_t} = \prod_{j=1}^T Y_j^T y_j^{(i)} \left(-\frac{y_t^{(i)}}{Y_t^T y_t^{(i)}} \right)$

$\frac{\partial P_{model}}{\partial V} = \sum_{t=1}^T \prod_{j=1}^T Y_j^T y_j^{(i)} \left(-\frac{y_t^{(i)}}{Y_t^T y_t^{(i)}} \right) \text{softmax}' \Big|_{res=Y_t} H_j^T$

$\Delta V = lr \cdot \frac{\frac{\partial P_{model}}{\partial V}}{\prod_{j=1}^T Y_j^T y_j^{(i)}} = \sum_{t=1}^T \frac{y_t^{(i)}}{Y_t^T y_t^{(i)}} \text{softmax}' \Big|_{res=Y_t} H_j^T$

其中 $\text{softmax}'$ 在上面 (*) 给出

② 对于 U : $\frac{\partial P_{model}}{\partial H_t}$

对于 $t=3$: $\frac{\partial P}{\partial H_3} = \frac{\partial P}{\partial Y_3} \frac{\partial Y_3}{\partial H_3} = \frac{\partial P}{\partial Y_3} \text{softmax}' \Big|_{Y_3} V$

其中 $\text{softmax}' \Big|_{Y_3}$ 与 $\frac{\partial P}{\partial Y_3}$ 都在上面给出

$t=2, t=1$: $\frac{\partial P}{\partial H_t} = \frac{\partial P}{\partial Y_t} \frac{\partial Y_t}{\partial H_t} + \frac{\partial P}{\partial H_{t+1}} \frac{\partial H_{t+1}}{\partial H_t}$

其中 $\frac{\partial H_{t+1}}{\partial H_t} = \text{sigmoid}' \Big|_{res=H_{t+1}} W$ 其中 $\text{sigmoid}' \Big|_{res=H_{t+1}} = \text{diag}[H_{t+1}(1-H_{t+1})]$

其中 $\frac{\partial P}{\partial Y_t}, \frac{\partial Y_t}{\partial H_t}$ 以及 $\frac{\partial P}{\partial H_{t+1}}$ 都在前面已经得到。

" \odot " 为 element wise 乘。

$\frac{\partial P}{\partial U} = \sum_{t=1}^T \frac{\partial P}{\partial H_t} \frac{\partial H_t}{\partial U}$ 其中 $\frac{\partial H_t}{\partial U} = \text{sigmoid}' \Big|_{res=H_t} \odot X_t^T$

$\frac{\partial P}{\partial H_t}$ 在上面给出

$\Delta U = lr \frac{\sum_{t=1}^T \frac{\partial P}{\partial H_t} \frac{\partial H_t}{\partial U}}{\prod_{j=1}^T Y_j^T y_j^{(i)}}$

③ 对于 W :

$\frac{\partial P}{\partial W} = \sum_{t=1}^T \frac{\partial P}{\partial H_t} \frac{\partial H_t}{\partial W} = \sum_{t=1}^T \frac{\partial P}{\partial H_t} \text{sigmoid}' \Big|_{res=H_t} H_{t+1}^T$ (设 $H_0 = 0$)

其中 $\text{sigmoid}' \Big|_{res=H_t}, \frac{\partial P}{\partial H_t}$ 都在上面给出。

有 $\Delta W = lr \frac{1}{\prod_{j=1}^T Y_j^T y_j^{(i)}} \cdot \sum_{t=1}^T \frac{\partial P}{\partial H_t} \frac{\partial H_t}{\partial W}$