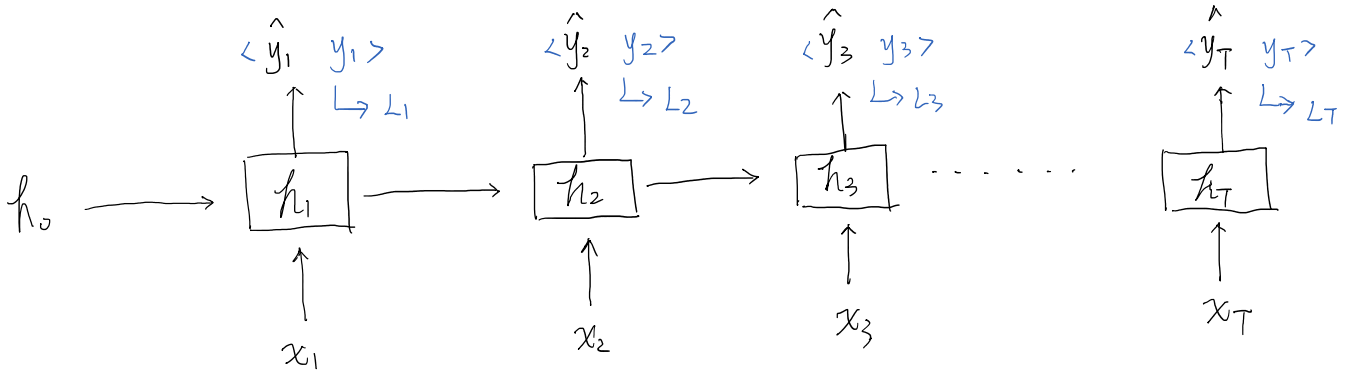


Backprop through time (BPTT)



$$L = \sum_{t=1}^T L_t \quad (\text{Summation of losses})$$

L_t can be:

- cross entropy: $-y \log \hat{y}$
- MSE: $(y - \hat{y})^2$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t + b)$$

\downarrow \downarrow
 W U

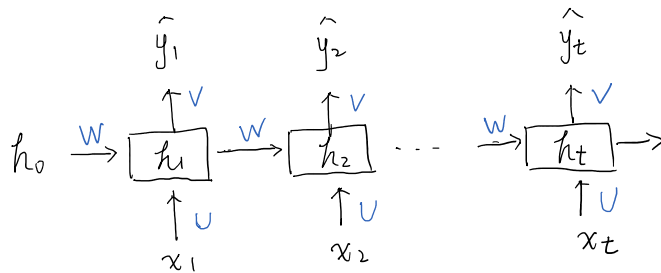
$$\hat{y}_t = g(W_{yh} h_t)$$

\downarrow
 V

non-linear function

$$h_t = \tanh(W h_{t-1} + U x_t + b)$$

$$\hat{y}_t = g(V h_t)$$



w, v, u are shared layers

To update weights:

compute $\frac{\partial L}{\partial W}$, $\frac{\partial L}{\partial V}$, $\frac{\partial L}{\partial U}$

Assumption:

- ① loss function is CE
- ② $g = \text{softmax}$ function

$$\textcircled{1} \frac{\partial Z}{\partial V}$$

$$\text{Let } Z_t = V h_t$$

Chain rule:

$$\begin{aligned} \frac{\partial Z}{\partial V} &= \sum_{t=1}^T \frac{\partial Z_t}{\partial V} \\ &= \sum_{t=1}^T \underbrace{\frac{\partial Z_t}{\partial \hat{y}_t}}_A \cdot \underbrace{\frac{\partial \hat{y}_t}{\partial Z_t}}_B \cdot \underbrace{\frac{\partial Z_t}{\partial V}}_C \end{aligned}$$

$$A: \frac{\partial Z_t}{\partial \hat{y}_t} = \frac{\partial (-y_t \log \hat{y}_t)}{\partial \hat{y}_t} = -y_t \frac{\partial \log \hat{y}_t}{\partial \hat{y}_t} = -\frac{y_t}{\hat{y}_t}$$

$$B: \frac{\partial \hat{y}_t}{\partial Z_t} = \frac{\partial (g(V h_t))}{\partial Z_t} = g' \cdot V$$

$$\begin{cases} Z_t = V h_t \\ g(Z_t) = \text{softmax}(Z_t) = \frac{e^{z_t}}{\sum_{k=1}^K e^{z_k}} \end{cases}$$

Compute g' :

$$\textcircled{1} \text{ case 1: } t=k \Rightarrow e^{z_t} = e^{z_k}$$

$$g' = \frac{\frac{\partial e^{z_t}}{\sum_{k=1}^K e^{z_k}}}{\partial Z_t} = \frac{e^{z_t}}{\sum_{k=1}^K e^{z_k}} - e^{z_t} \left[\frac{e^{z_t}}{\left(\sum_{k=1}^K e^{z_k}\right)^2} \right] = \hat{y}_t (1 - \hat{y}_t)$$

$$\textcircled{2} \text{ case 2: } t \neq k \Rightarrow e^{z_t} \neq e^{z_k} \quad (\text{treat } e^{z_k} \text{ as constant})$$

$$g' = \frac{-e^{z_t} \cdot e^{z_k}}{\left(\sum_{k=1}^K e^{z_k}\right)^2} = -\hat{y}_t \hat{y}_k$$

$$\textcircled{1} \text{ and } \textcircled{2}, \quad \frac{\partial \hat{y}_t}{\partial Z_t} = \begin{cases} \hat{y}_t (1 - \hat{y}_t), & t=k \\ -\hat{y}_t \hat{y}_k, & t \neq k \end{cases}$$

$$A \times B = \frac{\partial Z_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial Z_t} = -\frac{y_t}{\hat{y}_t} \begin{cases} \hat{y}_t (1 - \hat{y}_t), & t=k \\ -\hat{y}_t \hat{y}_k, & t \neq k \end{cases}$$

$$h_t = \tanh(W h_{t-1} + U x_t + b)$$

$$\hat{y}_t = g(V h_t)$$

$$Z_t = -y_t \log \hat{y}_t$$

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

$$A \times B \sim \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} = - \frac{\partial L}{\partial \hat{y}_t} \left\{ \begin{array}{l} -\hat{y}_t \hat{y}_k, \quad t \neq k \\ \hat{y}_t y_t - y_t, \quad t = k \\ y_t \hat{y}_k, \quad t \neq k \end{array} \right.$$

Summation of all k,

$$\rightarrow t=k \Rightarrow \hat{y}_k y_k - y_k$$

$$A \times B = \{ \hat{y}_t y_t - y_t \} + \sum_{t \neq k} y_t \hat{y}_k$$

$$= -y_k + \hat{y}_k \left[y_k + \sum_{t \neq k} y_t \right]$$

$$= -y_k + \hat{y}_k \sum_{j=1}^K y_j \rightarrow y_k = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \text{ one-hot vector}$$

$$= \hat{y}_k - y_k$$

$$C = \frac{\partial Z_t}{\partial v} = h_t$$

$$\frac{\partial L}{\partial v} = \sum_{t=1}^T (\hat{y}_t - y_t) \otimes h_t$$

outer product
vector \otimes vector \rightarrow matrix

$$\textcircled{2} \frac{\partial L}{\partial w}$$

Chain rule:

$$\frac{\partial L}{\partial w} = \sum_{t=1}^T \frac{\partial L_t}{\partial w} = \sum_{t=1}^T \underbrace{\frac{\partial L_t}{\partial \hat{y}_t}}_A \cdot \underbrace{\frac{\partial \hat{y}_t}{\partial h_t}}_B \cdot \underbrace{\frac{\partial h_t}{\partial w}}_C$$

$$\left. \begin{array}{l} h_t = \tanh(w h_{t-1} + U x_t + b) \\ \hat{y}_t = g(V h_t) \\ L_t = -y_t \log \hat{y}_t \\ \mathbf{1} = \sum_{t=1}^T \mathbf{1}_t \end{array} \right\}$$

$$A: \frac{\partial L_t}{\partial \hat{y}_t} = - \frac{y_t}{\hat{y}_t}$$

^

$$A \times B: \frac{\partial \hat{y}_t}{\partial h_t} = (\hat{y}_t - y_t) V$$

$$C: \frac{\partial h_t}{\partial w} = \frac{\partial \tanh(z_t)}{\partial w} = (1 - \tanh^2(z_t)) \cdot (h_{t-1} + w \cdot \frac{\partial h_{t-1}}{\partial w})$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\frac{\partial \tanh(x)}{\partial x} = 1 - \tanh^2(x)$$

recursion

$$h_{t-1} = \tanh(w h_{t-2} + U x_{t-1} + b)$$

(h_{t-1} is a function of w)

$$\frac{\partial h_{t-1}}{\partial w} = (1 - \tanh^2(z_{t-1})) \cdot (h_{t-2} + w \cdot \frac{\partial h_{t-2}}{\partial w})$$

$$\textcircled{3} \frac{\partial I}{\partial U}$$

$$h_t = \tanh(w h_{t-1} + U x_t + b)$$

Chain rule:

$$\frac{\partial I}{\partial U} = \sum_{t=1}^T \frac{\partial I_t}{\partial U} = \sum_{t=1}^T \frac{\partial I_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \frac{\partial h_t}{\partial U}$$

$$= \sum_{t=1}^T (\hat{y}_t - y_t) V \cdot \frac{\partial h_t}{\partial U}$$

$$\frac{\partial h_t}{\partial U} = (1 - \tanh^2(z_t)) \left(x_t + \frac{\partial (w h_{t-1})}{\partial U} \right)$$

$$\frac{\partial (w h_{t-1})}{\partial U} = w \cdot \frac{\partial h_{t-1}}{\partial U} + h_{t-1} \cdot \frac{\partial w}{\partial U} = 0$$

$$= w \cdot \frac{\partial (\tanh(w h_{t-2} + U x_{t-1} + b))}{\partial U}$$

(h_{t-1} is a function of U)

$$\frac{\partial h_t}{\partial U} = (1 - \tanh^2(z_t)) \left(x_t + w \cdot \frac{\partial h_{t-1}}{\partial U} \right)$$

recursion

