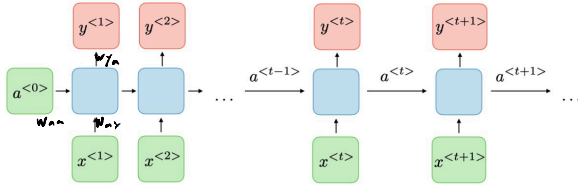




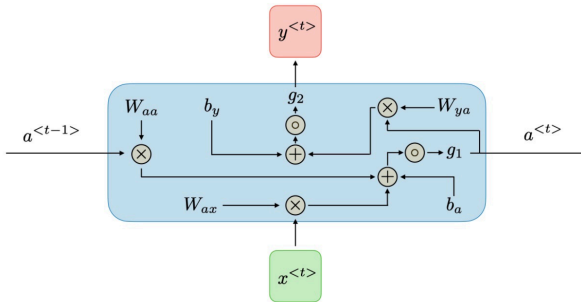
□ **Architecture of a traditional RNN** — Recurrent neural networks, also known as RNNs, are a class of neural networks that allow previous outputs to be used as inputs while having hidden states. They are typically as follows:



For each timestep t , the activation $a^{<t>}$ and the output $y^{<t>}$ are expressed as follows:

$$a^{<t>} = g_1(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a) \quad \text{and} \quad y^{<t>} = g_2(W_{ya}a^{<t>} + b_y)$$

where W_{ax} , W_{aa} , W_{ya} , b_a , b_y are coefficients that are shared temporally and g_1, g_2 activation functions.

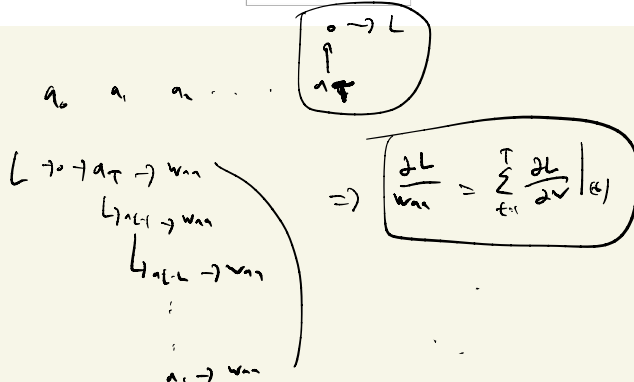


□ **Loss function** — In the case of a recurrent neural network, the loss function \mathcal{L} of all time steps is defined based on the loss at every time step as follows:

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}(\hat{y}^{<t>}, y^{<t>})$$

□ **Backpropagation through time** — Backpropagation is done at each point in time. At timestep T , the derivative of the loss \mathcal{L} with respect to weight matrix W is expressed as follows:

$$\frac{\partial \mathcal{L}^{(T)}}{\partial W} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(T)}}{\partial W} \Big|_{(t)}$$



$$C \in \mathbb{R}^{(n \times 1)}, \quad (K \text{ class}) \quad , \quad \text{label } y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad (n \times 1)$$

$$L = - \sum_{i=1}^n y_i \log(\hat{y}_i) = - \log(\hat{y}_n) \quad , \quad (n=1 \rightarrow n)$$

$$\frac{\partial L}{\partial \hat{y}_i} = - \frac{y_i}{\hat{y}_i}$$

$$\log \pi \quad 0 \in \mathbb{R}^k$$

$$\vec{y} = \text{softmax}(\mathbf{0})$$

$$S = \sum_k \exp(\mathbf{0}_k)$$

$$\hat{y}_i = \frac{e^{o_i}}{S}$$

$$\left(\left(\frac{f}{g} \right)' \right)' = \frac{f'g - fg'}{g^2}$$

$$\begin{aligned} \frac{\partial \hat{y}_i}{\partial o_j} &= \frac{\partial}{\partial o_j} \left(\frac{e^{o_i}}{S} \right) = \frac{e^{o_i} S - e^{o_i} \cdot e^{o_j}}{S^2} = \frac{e^{o_i} S}{S^2} - \frac{e^{o_i} e^{o_j}}{S^2} \\ &= \frac{e^{o_i}}{S} - \frac{e^{o_i} e^{o_j}}{S^2} \\ &= \frac{e^{o_i}}{S} \left(1 - \frac{e^{o_j}}{S} \right) = \hat{y}_i (1 - \hat{y}_j) \end{aligned}$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \frac{\partial}{\partial o_j} \left(\frac{e^{o_i}}{S} \right) = -e^{o_i} \frac{1}{S^2} \cdot \frac{\partial S}{\partial o_j} = - \frac{e^{o_i} \cdot e^{o_j}}{S^2} = -\hat{y}_i \cdot \hat{y}_j$$

Kronecker delta

$$\delta_{ij} \in \{0, 1\}$$

$$\delta_{ij} = \begin{pmatrix} 1 & i=j \\ 0 & \text{if } j \end{pmatrix} \quad \therefore$$

$$\frac{\partial \hat{y}_i}{\partial o_j} = \hat{y}_i (\delta_{ij} - \hat{y}_j)$$

$$f(x+\delta x) - f(x) = f'(x) \delta x = 0 +$$

$$\delta \hat{y}_i = \frac{\partial \hat{y}_i}{\partial \omega_j} \delta \omega_j$$

$$L = -\log(\hat{y}_n)$$

$$\frac{\partial L}{\partial \hat{y}_n} \delta \hat{y}_n = -\frac{1}{\hat{y}_n} \delta \hat{y}_n$$

$$\hat{y}_n = \frac{e^{\omega_n}}{S}$$

$$\delta \hat{y}_n = \frac{\partial \hat{y}_n}{\partial \omega_j} \delta \omega_j$$

$$\frac{\partial \hat{y}_i}{\partial \omega_j} = \hat{y}_i (\delta_{ij} - \hat{y}_j)$$

$$\frac{\partial \hat{y}_m}{\partial \omega_j} = -\hat{y}_m \hat{y}_j$$

$$\delta \hat{y}_m = -\hat{y}_m \hat{y}_j \cdot \delta \omega_j$$

$$\delta L = \hat{y}_j \cdot \delta \omega_j$$

$$L(\hat{y}) = L(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n)$$

$$\delta L = \sum_{i=1}^n \frac{\partial L}{\partial \hat{y}_i} \delta \hat{y}_i$$

$$= \sum_{i=1}^n \left(-\frac{\hat{y}_i}{\hat{y}_i} \right) \cdot \delta \hat{y}_i$$

$$\delta L = -\frac{1}{\hat{y}_n} \delta \hat{y}_n$$

$$\delta L = \sum \left(\frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \omega_j} \right) \delta \omega_j$$

$$\frac{\partial L}{\partial \omega_j} = \sum \left(\frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \omega_j} \right)$$

$$\frac{\partial L}{\partial \omega_j} = \sum \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \omega_j}$$

$$L_t = - \sum_{i=1}^n y_{t,i} \ln(\hat{y}_{t,i})$$

$$\frac{\partial L_t}{\partial \hat{y}_{t,i}} = - \frac{y_{t,i}}{\hat{y}_{t,i}}$$

$$\frac{\partial \hat{y}_{t,i}}{\partial \theta_{t,j}} = \hat{y}_{t,i} (\delta_{ij} - \hat{y}_{t,j})$$

$$\frac{\partial L_t}{\partial \theta_{t,j}} = \sum \frac{\partial L_t}{\partial \hat{y}_{t,i}} \frac{\partial \hat{y}_{t,i}}{\partial \theta_{t,j}}$$

$$= \sum_i -y_{t,i} (\delta_{ij} - \hat{y}_{t,j})$$

$$= \sum_i -y_{t,i} \delta_{ij} + \sum_i y_{t,i} \hat{y}_{t,j}$$

$$= \hat{y}_{t,j} - y_{t,j}$$

$$\frac{\partial L_t}{\partial \theta_t} = \hat{y}_t - y_t = \delta_t^y$$

$$\frac{\partial L_t}{\partial (w_{t+1})_j} = \sum \frac{\partial L_t}{\partial \phi_{t+1,k}} \frac{\partial \phi_{t+1,k}}{\partial (w_{t+1})_j}$$

$$\frac{\partial \phi_{t+1,k}}{\partial (w_{t+1})_j} = \delta_{k,j} \phi_{t+1,j}, \quad \frac{\partial L_t}{\partial \phi_{t+1,k}} = (\delta_t^y)_k$$

$$\frac{\partial L_c}{\partial (w_1)_j} = \sum_k \frac{\partial L_c}{\partial b_k} \delta x_k (h_c)_j$$

$$= (h_c)_j \sum_k \frac{\partial L_c}{\partial b_k} \delta x_k$$

$$= (h_c)_j \frac{\partial L_c}{\partial (b_c)_1} = (h_c)_j \left(\delta^2 \right)_1$$

$$\frac{\partial L_c}{\partial w_1} = \delta_c^2 h_c^T$$

$$\frac{\partial L}{\partial w_1} = \sum_c \delta_c^2 h_c^T$$

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

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

$$u_t = W_{aa}h_{t-1} + W_{ax}x_t + b_a, h_t = \tanh(u_t), o_t = W_{ya}h_t + b_y, \hat{y}_t = \text{softmax}(o_t), L_t = -\sum_{i=1}^K y_{t,i} \log(\hat{y}_{t,i}), \mathcal{L} = \sum_{t=1}^T L_t$$

$$\delta_t^y \equiv \frac{\partial L_t}{\partial o_t} = \hat{y}_t - y_t, \delta_{T+1}^u \equiv 0, \delta_t^h \equiv \frac{\partial \mathcal{L}}{\partial h_t} = W_{ya}^\top \delta_t^y + W_{aa}^\top \delta_{t+1}^u, \delta_t^u \equiv \frac{\partial \mathcal{L}}{\partial u_t} = (1 - h_t^2) \odot \delta_t^h.$$

$$\frac{\partial \mathcal{L}}{\partial W_{ya}} = \sum_{t=1}^T \delta_t^y h_t^\top, \frac{\partial \mathcal{L}}{\partial b_y} = \sum_{t=1}^T \delta_t^y, \frac{\partial \mathcal{L}}{\partial W_{aa}} = \sum_{t=1}^T \delta_t^u h_{t-1}^\top, \frac{\partial \mathcal{L}}{\partial W_{ax}} = \sum_{t=1}^T \delta_t^u x_t^\top, \frac{\partial \mathcal{L}}{\partial b_a} = \sum_{t=1}^T \delta_t^u.$$

$$W_{ya} \leftarrow W_{ya} - \eta \frac{\partial \mathcal{L}}{\partial W_{ya}}, b_y \leftarrow b_y - \eta \frac{\partial \mathcal{L}}{\partial b_y}, W_{aa} \leftarrow W_{aa} - \eta \frac{\partial \mathcal{L}}{\partial W_{aa}}, W_{ax} \leftarrow W_{ax} - \eta \frac{\partial \mathcal{L}}{\partial W_{ax}}, b_a \leftarrow b_a - \eta \frac{\partial \mathcal{L}}{\partial b_a}.$$

(A, B is matrix)

$$\|AD\| \leq \|A\| \|D\|$$

$$T \approx 3$$

$$\delta_j^n = \dots + \dots + (p_r w_{rn}^T D_r w_{rn}^T D_s) (w_{rn}^T \delta_s^n)$$

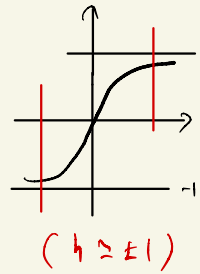
$$J_j = w_{jn}^T D_j \quad (\text{Jacobian})$$

$$p_r w_{rn}^T D_r w_{rn}^T D_s = p_r J_r J_s$$

$$\|J_r J_s\| \leq \|J_r\| \|J_s\|$$

$$D_j = \text{diag}(1 - h_{ji}^2)$$

$$0 \leq 1 - h_{ji}^2 \leq 1, \quad h = \tanh(u)$$



$$(1 - h_{ji}^2) \geq 0$$

$$\left\| \prod_{t \in \mathcal{T}} J_t \right\| \geq 0 \Rightarrow \underline{\text{gradient vanishing}}$$

exploding

- Gradient clipping

$$g_{\text{clipped}} = g \cdot \min\left(1, \frac{C}{\|g\|}\right)$$

