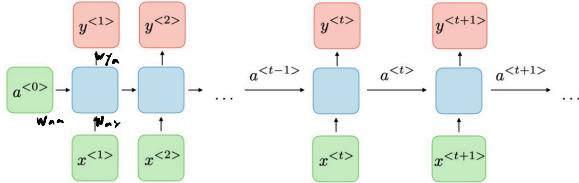


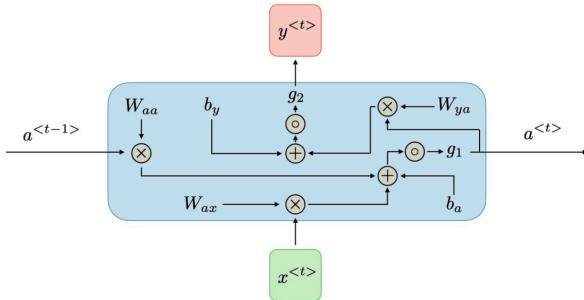
□ **Architecture of a traditional RNN** — Recurrent neural networks, also known as RNNs, are a class of neural networks that allow previous outputs to be used as inputs while having hidden states. They are typically as follows:



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}, b_a, b_y are coefficients that are shared temporally and g_1, g_2 activation functions.

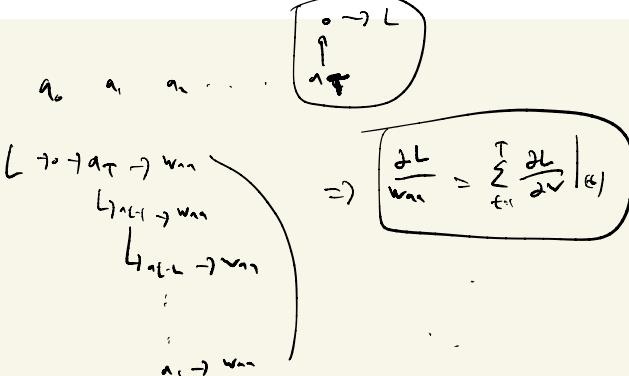


□ **Loss function** — In the case of a recurrent neural network, the loss function \mathcal{L} of all time steps is defined based on the loss at every time step as follows:

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

□ **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:

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



$$C \in \text{loss}, (K_{\text{class}}), \text{label}_j = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} (\delta_{\text{class}})$$

$$L = - \sum_{i=1}^k y_i \log(g_i) = - \log(g_m), (\text{label}_m)$$

$$\boxed{\frac{\partial L}{\partial \theta_j} = -\frac{y_j}{g_j}}$$

$$\text{logit } o \in \mathbb{R}^k$$

$$\vec{g} = \text{softmax}(o)$$

$$S = \sum_k \exp(o_k)$$

$$g_i = \frac{e^{o_i}}{S}$$

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

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

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

Kronecker delta

$$\delta_{ij} \in \text{setS}$$

$$\delta_{ij} = \begin{cases} 1 & i=j \\ 0 & \text{if } j \end{cases} \Rightarrow$$

$$\boxed{\frac{\partial \vec{g}_i}{\partial o_j} = \vec{g}_i (\delta_{ij} - \vec{g}_j)}$$

$$f(x+o^2) - f(x) = f'(x)o^2 = o +$$

$$L(\vec{g}) = L(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n)$$

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

$$L = -\log(\vec{g}_n)$$

$$\delta L \frac{\partial L}{\partial \vec{g}_n} \delta \vec{g}_n = -\frac{1}{\vec{g}_n} \delta \vec{g}_n$$

$$\vec{g}_n = \frac{e^{0_n}}{S}$$

$$\delta \vec{g}_n = \frac{\partial \vec{g}_n}{\partial o_j} \delta o_j$$

$$\frac{\partial \vec{g}_i}{\partial o_j} = \vec{g}_i (\delta_{ij} - \vec{g}_j)$$

$$\frac{\partial \vec{g}_m}{\partial o_j} = -\vec{g}_m \cdot \vec{g}_j$$

$$\delta \vec{g}_n = \vec{g}_n \cdot \vec{g}_j \cdot \delta o_j$$

$$\delta L = \vec{g}_j \cdot \delta o_j$$

$$\delta L = \sum_{i=1}^k \frac{\partial L}{\partial \vec{g}_i} \cdot \delta \vec{g}_i$$

$$= \sum_{i=1}^k \left(\frac{g_i}{\vec{g}_i} \right) \cdot \delta \vec{g}_i$$

$$\boxed{\delta L = -\frac{1}{\vec{g}_n} \delta \vec{g}_n}$$

$$\delta L = \sum \left(\frac{\partial L}{\partial \vec{g}_i} \frac{\partial \vec{g}_i}{\partial o_j} \right) \delta o_j$$

$$\frac{\partial L}{\partial o_j} = \sum \left(\frac{\partial L}{\partial \vec{g}_i} \frac{\partial \vec{g}_i}{\partial o_j} \right)$$

$$\frac{\partial L}{\partial o_j} = \sum \frac{\partial L}{\partial \vec{g}_i} \frac{\partial \vec{g}_i}{\partial o_j}$$

$$L_t = - \sum_{i,j} g_{t,i} \log(g_{t,i})$$

$$\frac{\partial L_t}{\partial \hat{g}_{t,i}} = - \frac{g_{t,i}}{\hat{g}_{t,i}}$$

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

$$\frac{\partial L_t}{\partial \theta_{t,j}} = \sum_i \frac{\partial L_t}{\partial \hat{g}_{t,i}} \frac{\partial \hat{g}_{t,i}}{\partial \theta_{t,j}}$$

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

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

$$= \hat{g}_{t,j} - g_{t,j}$$

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

$$\frac{\partial L_t}{\partial (w_{t,j})_i} = \left\{ \frac{\partial L_t}{\partial \theta_{t,h}} \right\} \frac{\partial \theta_{t,h}}{\partial (w_{t,j})_i}$$

$$\frac{\partial \theta_{t,h}}{\partial (w_{t,j})_i} = \delta_{kj} (\mathbf{1}_n)_j , \quad \frac{\partial L_t}{\partial \theta_{t,h}} = (\delta_t^y)_k$$

$$\frac{\partial L_c}{\partial (w_k)_{ij}} = \sum_k \frac{\partial L_c}{\partial w_{ki}} \delta_{ki} (w_k)_j$$

$$= (w_k)_j \left\{ \frac{\partial L_c}{\partial w_{ki}} \delta_{ki} \right.$$

$$= (w_k)_j \frac{\partial L_c}{\partial (w_k)_i} = (w_k)_j \left(\frac{\partial^2 L_c}{\partial w_{ki}^2} \right)_i$$

$$\frac{\partial L_c}{\partial w_{ki}} = \delta_{ki}^{-1} \mathbf{h}_k^\top$$

$$\frac{\partial L}{\partial w_{ki}} = \left\{ \begin{array}{l} \delta_{ki}^{-1} \mathbf{h}_k^\top \\ \vdots \end{array} \right.$$

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

$$\frac{1}{1+e^{-\text{tanh}(z)}}=\text{tanh}(z)$$

$$\begin{aligned} \mathcal{L} &= \sum_{i=1}^N \left[\frac{1}{2} \|y_i - \hat{y}_i\|^2 + \frac{\lambda}{2} \left(\sum_{j=1}^J w_j^2 \right) \right] \\ &= \frac{1}{2} \|Y - \hat{Y}\|^2 + \frac{\lambda}{2} \|W\|^2 \end{aligned}$$

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

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

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

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

$(A, B \text{ is matrix})$

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

$$T = 3$$

$$\delta_i = \dots + \left(p_i w_m^T D_2 w_m^T D_3 \right) \left(w_m^T \delta_{j+2} \right)$$

$$J_j = w_m^T D_j \quad (\text{jacobian})$$

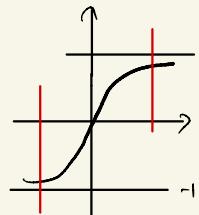
$$p_i w_m^T D_2 w_m^T D_3 = J_1 J_2 J_3$$

$$\|JT\| \leq \|J_1\| \|J_2\| \|J_3\|$$

$$D_j = \frac{\partial \log(\tau h_j)}{\partial}$$

$$0 \leq \tau h_j \leq 1, \quad h = \tanh(u)$$

$$(\tau h_j) \approx 0$$



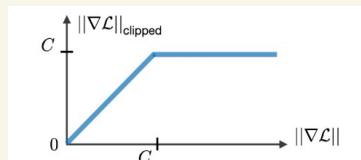
$$(h \approx \pm 1)$$

$$\left\| \prod_{t=1}^T J_t \right\| \leq 0 \Rightarrow \text{gradient vanishes}$$

Explizit

- Gradient of \mathcal{L}_{CE}

$$\delta_{CE} = g \cdot \ln \left(1, \frac{C}{\|g\|} \right)$$



$$\frac{\partial h_{T+1}}{\partial x_i} \approx 0$$

$$\frac{\partial L_T}{\partial x_1} = \frac{\partial L_T}{\partial h_{T-1}} \cdot \frac{\partial h_{T-1}}{\partial x_1} = 0$$

start error

$$\frac{\partial L}{\partial w_{in}} = \sum_{k=1}^T \sum_{t=k}^T \frac{\partial L_t}{\partial w_{in}(k)}$$

$$\frac{\partial L}{\partial w_{m,n}(k)} = \delta_k^{u_i(T)} \cdot h_m^T$$

$$\delta_K^{w_i(T)} = 0 \quad (\Rightarrow) \quad \frac{\partial h_T}{\partial w_m(x)} = 0$$

If the gradient from L_T doesn't reach the early W_{aa} , then the early parameters can't be trained to reduce L_T . As a result, the model doesn't learn to maintain (store) information over time, and it fails to handle long-term dependencies.