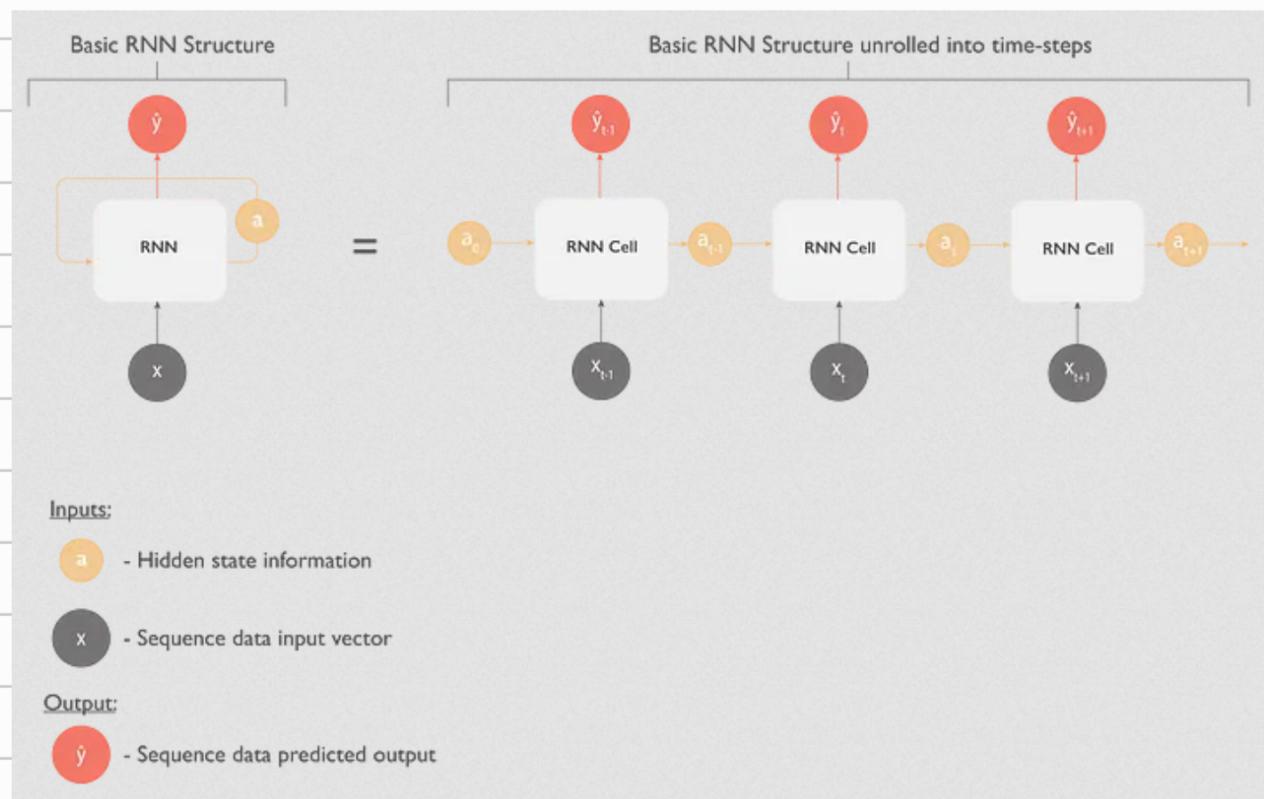


Basic Architecture



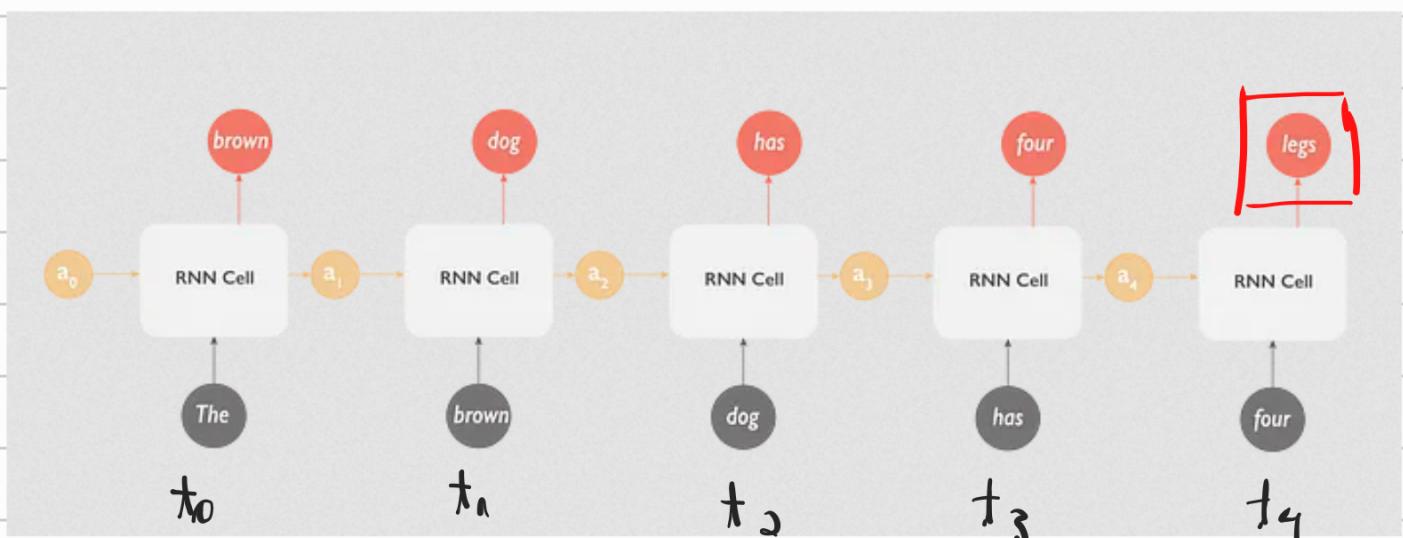
We unroll RNN basic structure (remove cycles) into repetitive chain of cells to correspond to the length of data

The input and output of RNN structure can vary

RNN types

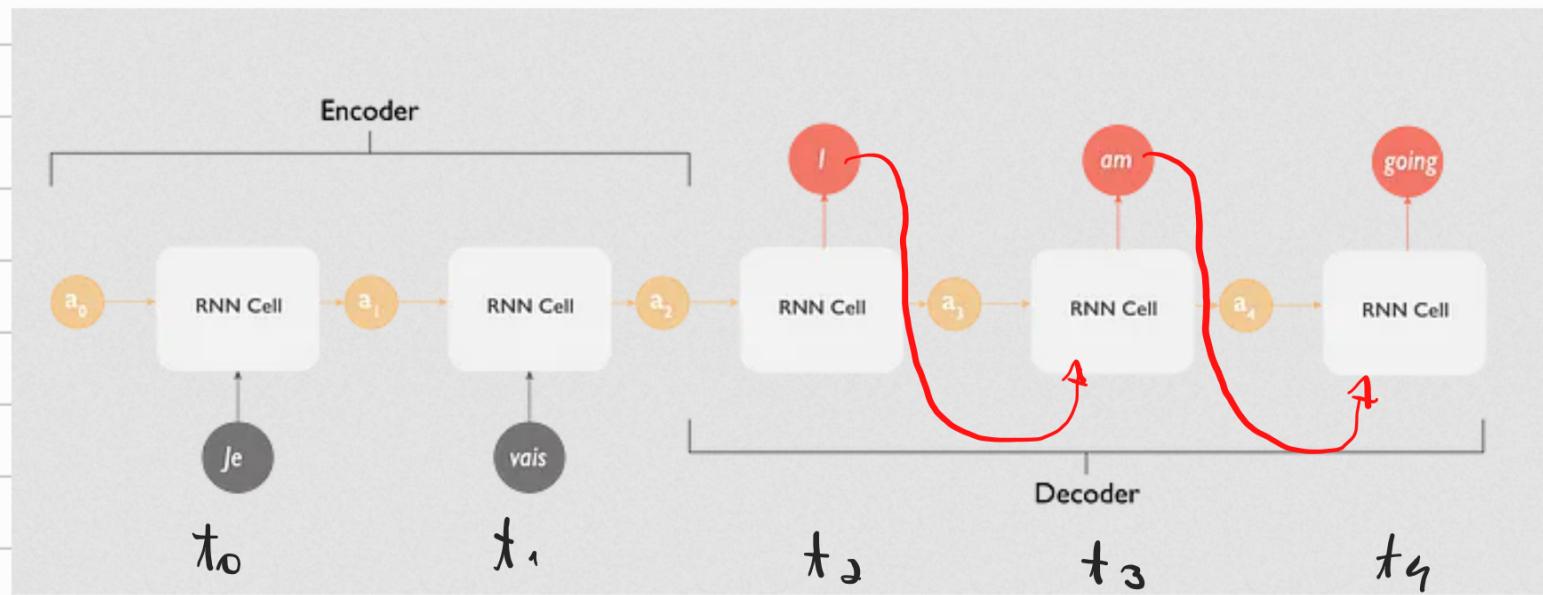
Many-to-Many (Same length)

Input number = Output number every time-step



Many-to-Many (Different length)

Input number ≠ Output number

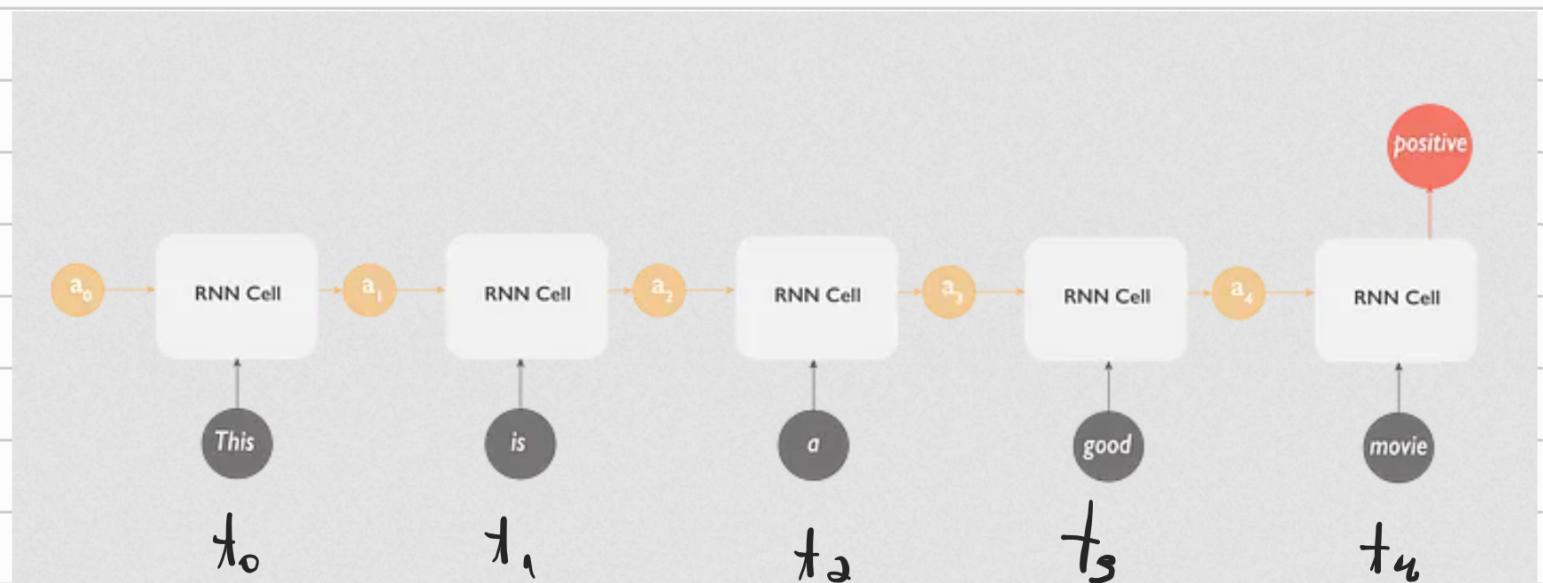


Encoder → Takes an input seq and maps to some internal state

Decoder → Takes internal state and generate some output

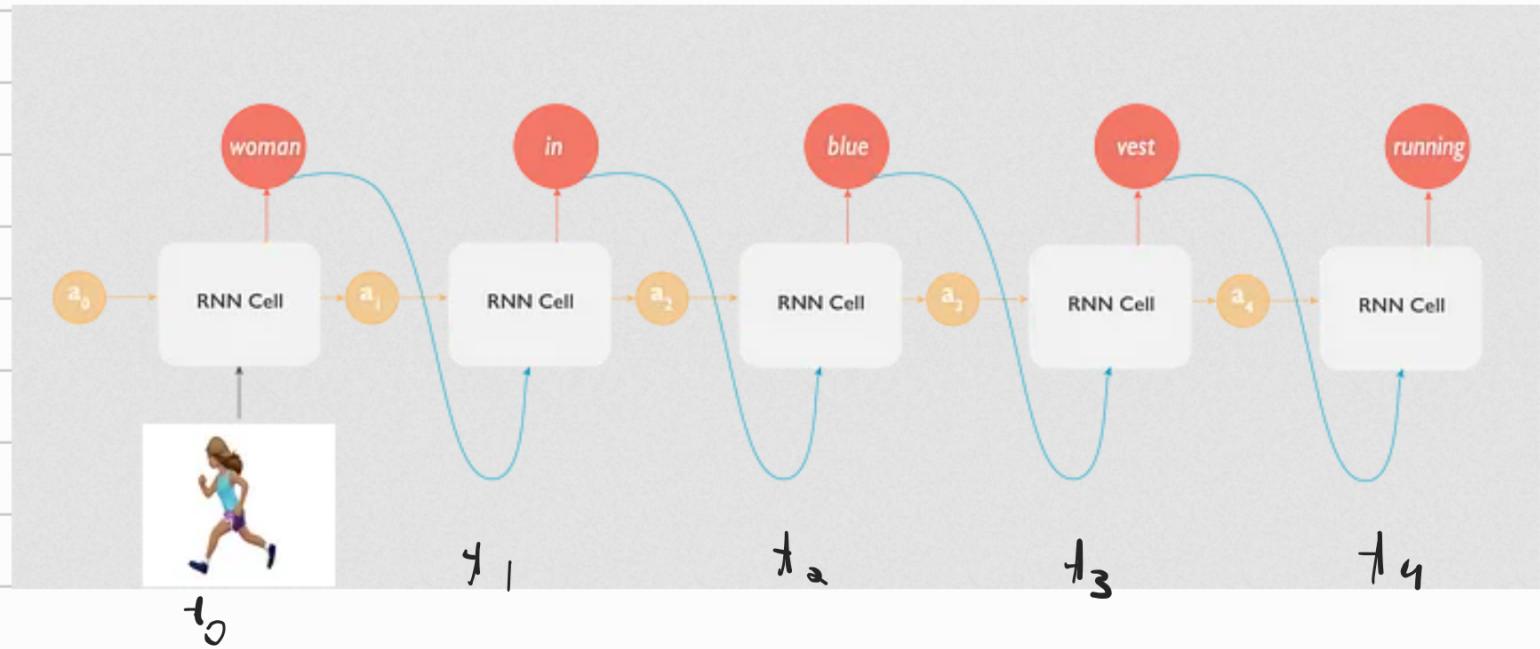
Many-to-One

Has many inputs at each time step, but outputs just a single value in last time step

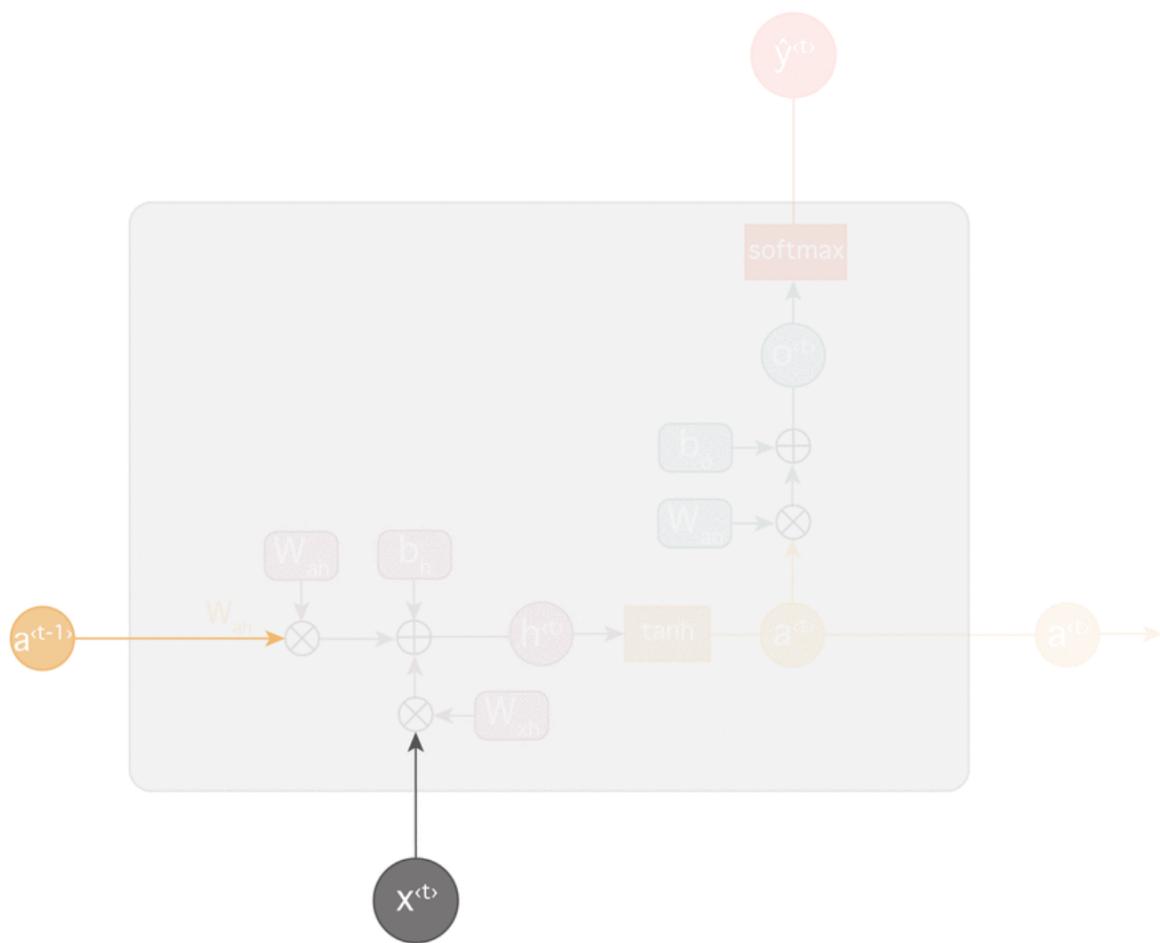


One-to-Many

Takes a single input at the first time step and outputs a sequence of values in remaining time steps



Forward RNN



Step 1) The cell receives 2 inputs $x^{<t>}$ and $a^{<t-1>}$ from Input cell t

Forward Propagation Equations

$$1. h^{<t>} = (W_{xh} * x^{<t>}) + (W_{ah} * a^{<t-1>}) + b_h$$

$$2. a^{<t>} = \tanh(h^{<t>})$$

$$3. o^{<t>} = (W_{ao} * a^{<t>}) + b_o$$

$$5. \hat{y}^{<t>} = \text{softmax}(o^{<t>})$$

Inputs:

$a^{<t-1>}$ - hidden state activation at previous time-step $t-1$

$x^{<t>}$ - input at time-step t

Vector Operators:

\otimes - matrix multiplication

\oplus - matrix addition

Activation Functions:

\tanh - hyperbolic tangent function

softmax - softmax function

Weight Matrices & Bias Vectors:

W_{xh} - hidden-to-hidden weight matrix

W_{ah} - input-to-hidden weight matrix

b_h - hidden state bias vector

W_{ao} - hidden-to-output weight matrix

b_o - output bias vector

Outputs:

$h^{<t>}$ - hidden state at time-step t

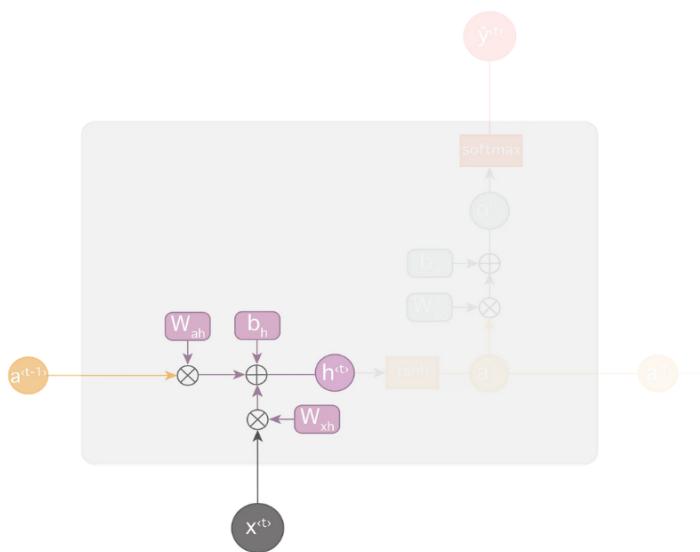
$a^{<t>}$ - hidden state activation at time-step t

$\hat{o}^{<t>}$ - Predicted output at time-step t

$\hat{y}^{<t>}$ - Predicted output activation at time-step t

Step 2) Do matrix multiplication between W_{xh} & $x^{<t>}$ and between W_{ah} & $a^{<t-1>}$. Next, it computes $h^{<t>}$ by adding the above together with bias term

$$(i) \quad h^{<t>} = W_{xh} * x^{<t>} + W_{ah} * a^{<t-1>} + b_h$$



Forward Propagation Equations

$$1. h^{<t>} = (W_{xh} * x^{<t>}) + (W_{ah} * a^{<t-1>}) + b_h$$

$$2. a^{<t>} = \tanh(h^{<t>})$$

$$3. o^{<t>} = (W_{ao} * a^{<t>}) + b_o$$

$$5. \hat{y}^{<t>} = \text{softmax}(o^{<t>})$$

Inputs:

$a^{<t-1>}$ - hidden state activation at previous time-step $t-1$

$x^{<t>}$ - input at time-step t

Vector Operators:

\otimes - matrix multiplication

\oplus - matrix addition

Activation Functions:

\tanh - hyperbolic tangent function

softmax - softmax function

Weight Matrices & Bias Vectors:

W_{ah} - hidden-to-hidden weight matrix

W_{xh} - input-to-hidden weight matrix

b_h - hidden state bias vector

W_{ao} - hidden-to-output weight matrix

b_o - output bias vector

Outputs:

$h^{<t>}$ - hidden state at time-step t

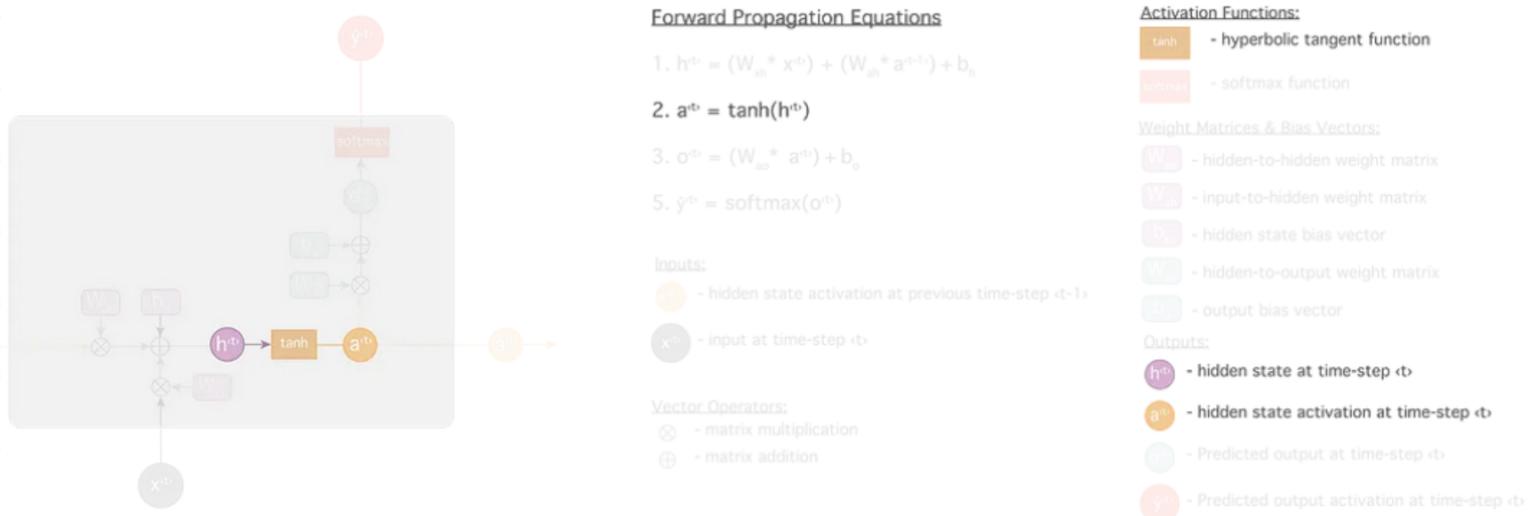
$a^{<t>}$ - hidden state activation at time-step t

$\hat{o}^{<t>}$ - Predicted output at time-step t

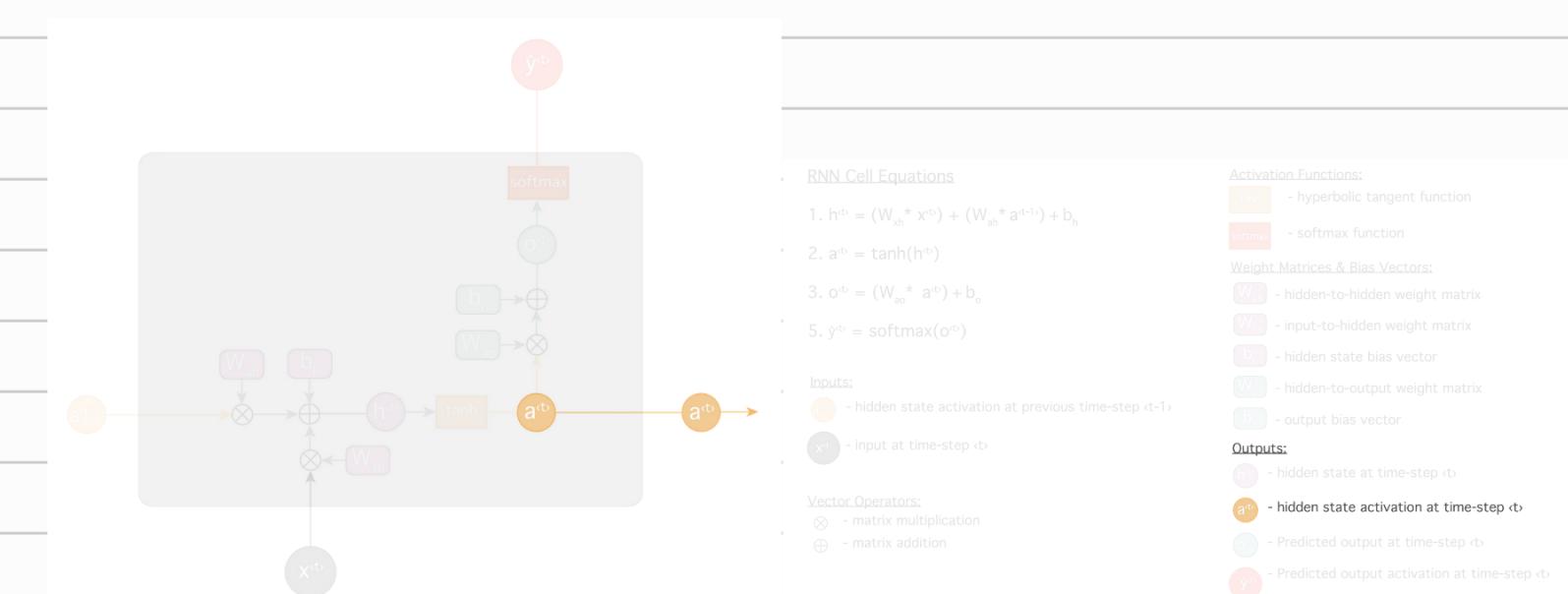
$\hat{y}^{<t>}$ - Predicted output activation at time-step t

Step 3) Calculate $a^{<+>}$ by passing $h^{<+>}$ through an activation function f (tanh or ReLU). I will use tanh as example

$$(ii) \quad a^{<+>} = f(h^{<+>})$$

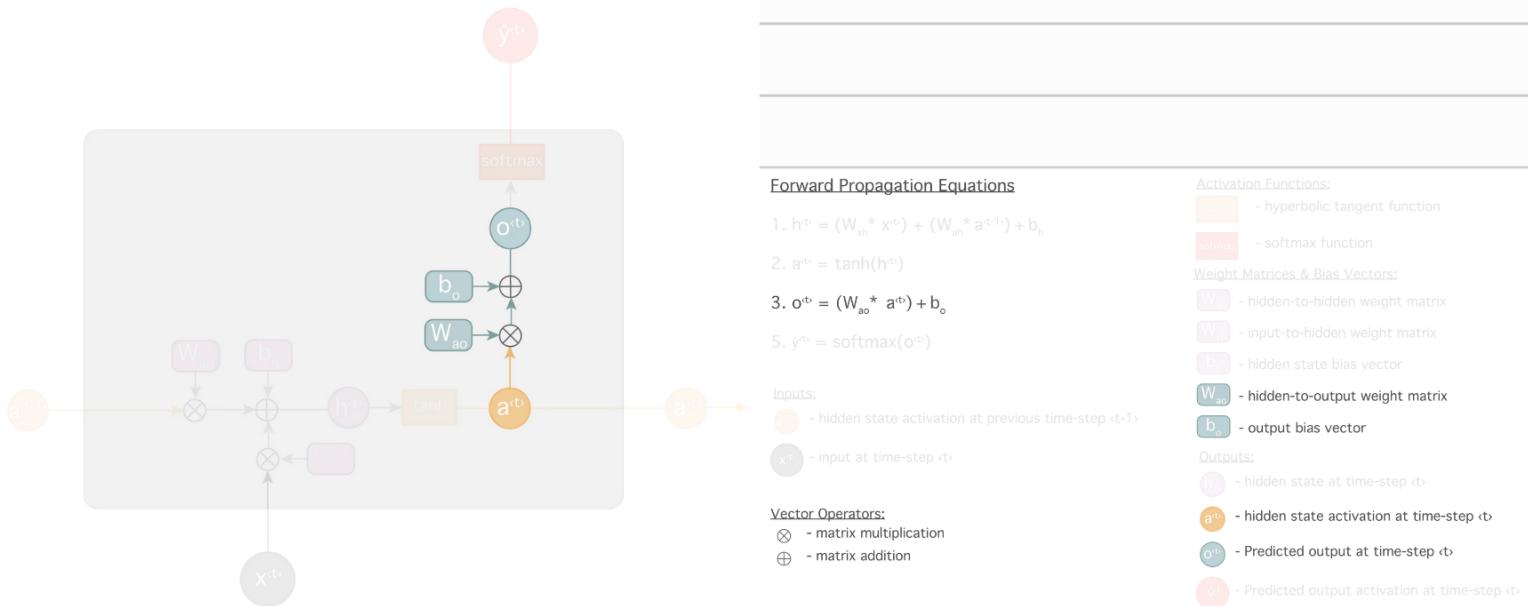


Step 4) The value $a^{<+>}$ is passed to the next cell



Step 5) Calculates the unnormalized log probability of each possible value of true output $\hat{o}^{<t>}$

$$(iii) \hat{o}^{<t>} = W_{ao} * a^{<t>} + b_o$$

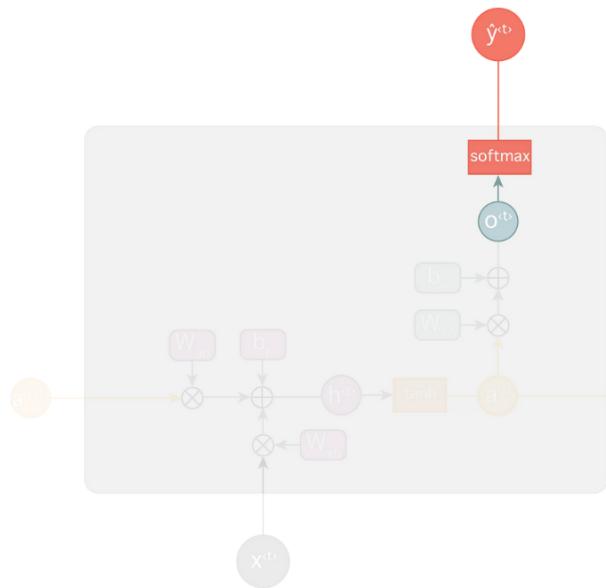


Step 6) Calculates the vector of normalized probabilities $\gamma^{<t>}$ by passing $\hat{o}^{<t>}$ through an activation function g (sigmoid or softmax)

The choice of g depends on the expected output, for example, for binary class output use sigmoid (remember logistic regression), but for multi-class output use softmax

$$(iv) \gamma^{<t>} = g(\hat{o}^{<t>})$$

Note: In examples I'll use softmax



Forward Propagation Equations

$$1. h^t = (W_{xh} * x^{t+1}) + (W_{ah} * a^{t-1}) + b_h$$

$$2. a^t = \tanh(h^t)$$

$$3. o^t = (W_{ao} * a^t) + b_o$$

$$5. \hat{y}^{t+1} = \text{softmax}(o^t)$$

Inputs:

- hidden state activation at previous time-step a^{t-1}

- input at time-step x^{t+1}

Vector Operators:

\otimes - matrix multiplication

\oplus - matrix addition

Activation Functions:

hyperbolic tangent function

softmax function

Weight Matrices & Bias Vectors:

W_{xh} - hidden-to-hidden weight matrix

W_{ah} - input-to-hidden weight matrix

b_h - hidden state bias vector

W_{ao} - hidden-to-output weight matrix

b_o - output bias vector

Outputs:

a^t - hidden state at time-step t

h^t - hidden state activation at time-step t

o^t - Predicted output at time-step t

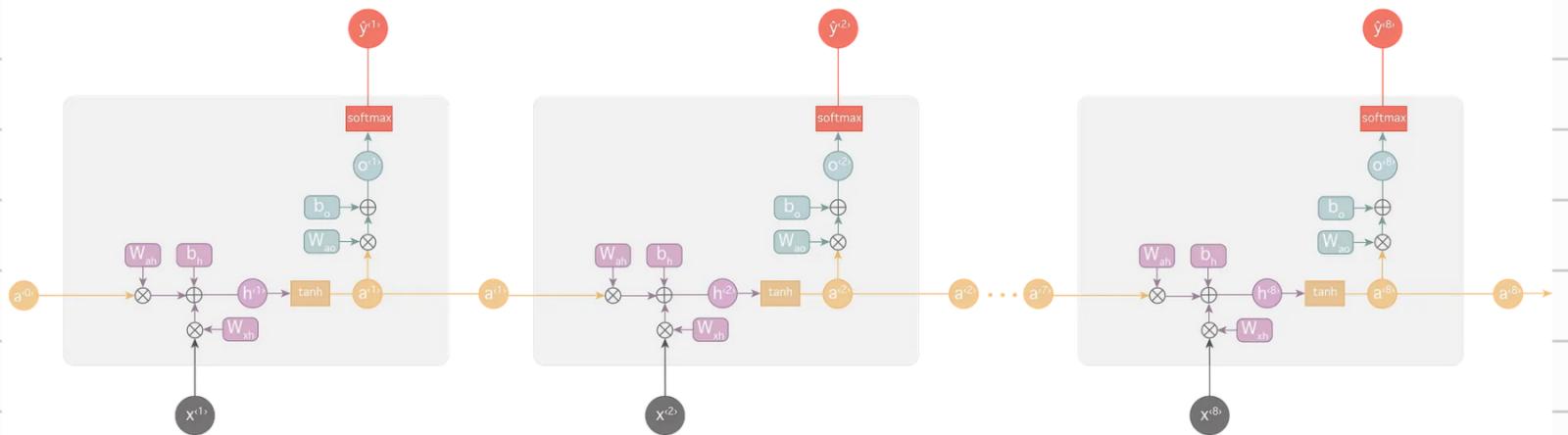
\hat{y}^{t+1} - Predicted output activation at time-step t

Forward propagation through time

1) Initializes a hidden state $a^{<0>}$ and $W_{xa}, W_{ah}, W_{ao}, b_h, b_o$

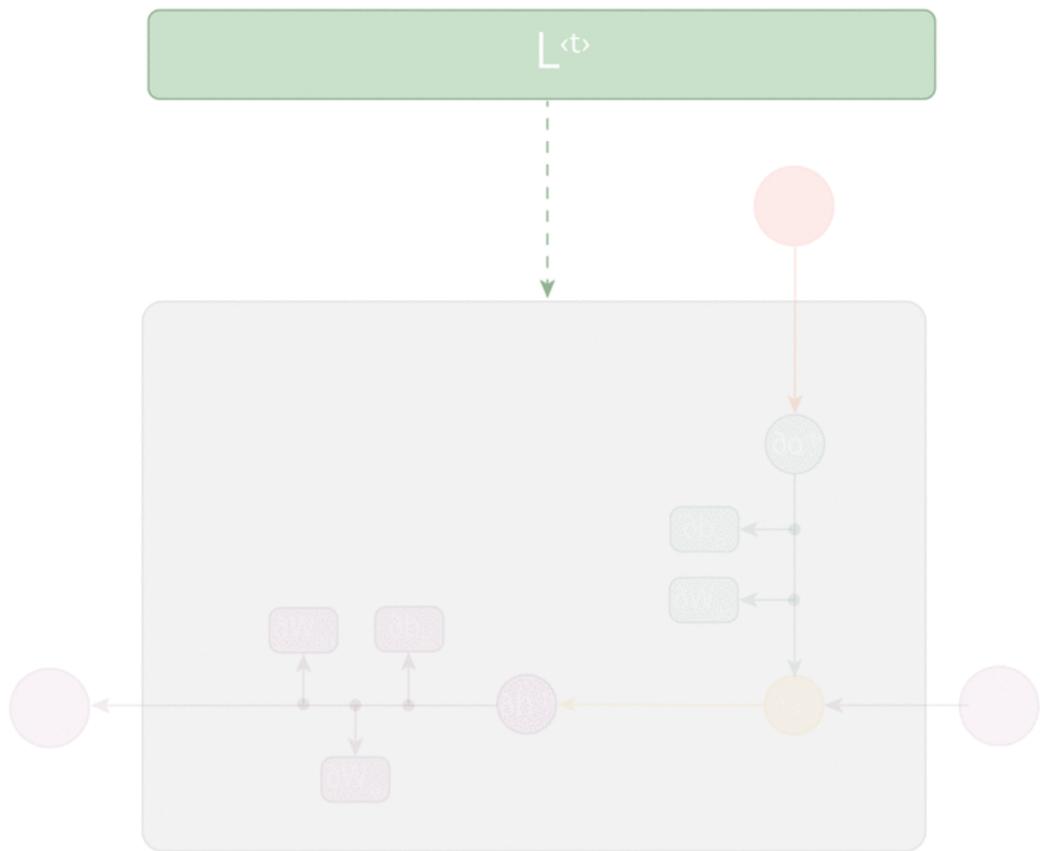
There will be shared across all time steps $t=1$ to T

2) Repeat individual steps above for each cell in time



Backward RNN

Goal → Compute $\frac{\partial L}{\partial w_{th}}, \frac{\partial L}{\partial w_{ah}}, \frac{\partial L}{\partial w_{ao}}, \frac{\partial L}{\partial b_h}, \frac{\partial L}{\partial b_o}$



Step 1) Choose some loss function based on your task. As example I'll use Cross-Entropy

$$L(y^{(t)}, \hat{y}^{(t)}) = -\sum_i^C y_i^{(t)} \log \hat{y}_i^{(t)} \quad (1)$$

Note: For any loss L , $L \in \mathbb{R}$

Removing i from summatory

$$L^{<t>}(\hat{y}^{<t>}, \hat{\bar{y}}^{<t>}) = -\hat{y}_i^{<t>} \log \hat{y}_i^{<t>} + \sum_{i \neq k}^C -\hat{y}_k^{<t>} \log \hat{y}_k^{<t>} \quad (2)$$

$L^{<t>}(\hat{y}^{<t>}, \hat{\bar{y}}^{<t>})$ = Cross-entropy loss function
for multi-class outputs at time t

C = number of classes

$y_i^{<t>}$ = True output for class i at time t

$\hat{y}_i^{<t>}$ = Predicted output for class i at time t

$y_k^{<t>}$ = True output for other class k at time t

$\hat{y}_k^{<t>}$ = Predicted output for other class k at time t

$y_i^{<t>}, y_k^{<t>} \in y^{<t>}$ and $\hat{y}_i^{<t>}, \hat{y}_k^{<t>} \in \hat{y}^{<t>}$

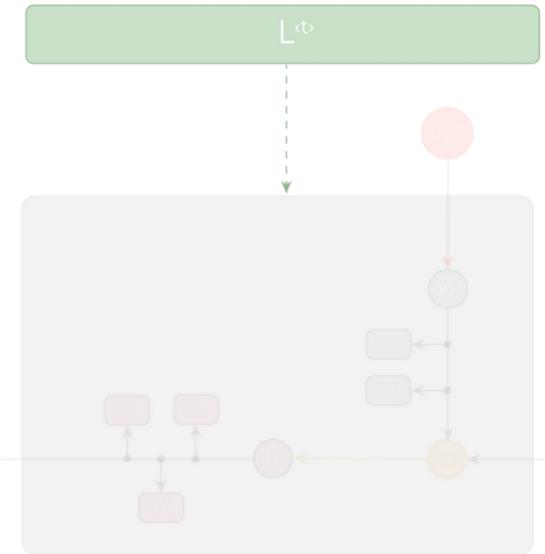
Note

(a) $y^{<t>}$ is one-hot encoded. Therefore, $y_i^{<t>}$ is 0 or 1

$$\text{eg.: } y^{<t>} = \begin{pmatrix} y_0^{<t>} \\ y_1^{<t>} \\ y_2^{<t>} \\ \vdots \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix} \begin{matrix} \xrightarrow{\text{Class 0}} \\ \xrightarrow{\text{Class 1}} \\ \xrightarrow{\text{Class 2}} \end{matrix} \therefore y^{<t>} \text{ is class t}$$

(b) $\hat{y}^{<t>}$ is normalized. Therefore, $\hat{y}_i^{<t>} \in [0, 1]$ and

$$y_i^{<t>} \in \mathbb{R}$$



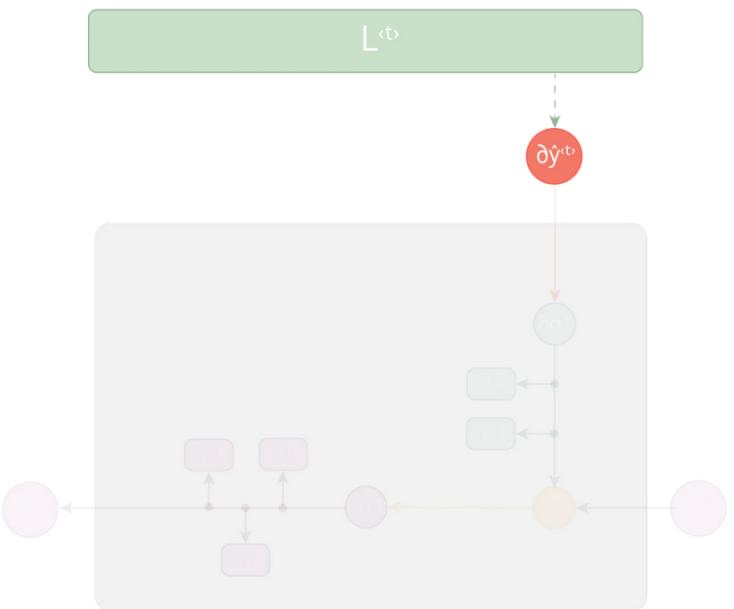
Step 2) Calculate $\frac{\partial L}{\partial \hat{y}^{t\Phi}} = \left(\frac{\partial L}{\partial \hat{y}_1^{t\Phi}}, \frac{\partial L}{\partial \hat{y}_2^{t\Phi}}, \dots, \frac{\partial L}{\partial \hat{y}_C^{t\Phi}} \right)^T$

$$\frac{\partial L}{\partial \hat{y}_i^{t\Phi}} = \frac{\partial}{\partial \hat{y}_i^{t\Phi}} \left(-y_i^{t\Phi} \log \hat{y}_i^{t\Phi} + \sum_{i \neq k} -y_k^{t\Phi} \log \hat{y}_k^{t\Phi} \right)$$

$$\frac{\partial L}{\partial \hat{y}_i^{t\Phi}} = -\frac{y_i^{t\Phi}}{\hat{y}_i^{t\Phi}} \quad (3)$$

$$\frac{\partial L}{\partial \hat{y}_k^{t\Phi}} = \frac{\partial}{\partial \hat{y}_k^{t\Phi}} \left(-y_i^{t\Phi} \log \hat{y}_i^{t\Phi} + \sum_{i \neq k} -y_k^{t\Phi} \log \hat{y}_k^{t\Phi} \right)$$

$$\frac{\partial L}{\partial \hat{y}_k^{t\Phi}} = \sum_{i \neq k} -\frac{y_k^{t\Phi}}{\hat{y}_k^{t\Phi}} \quad (4)$$

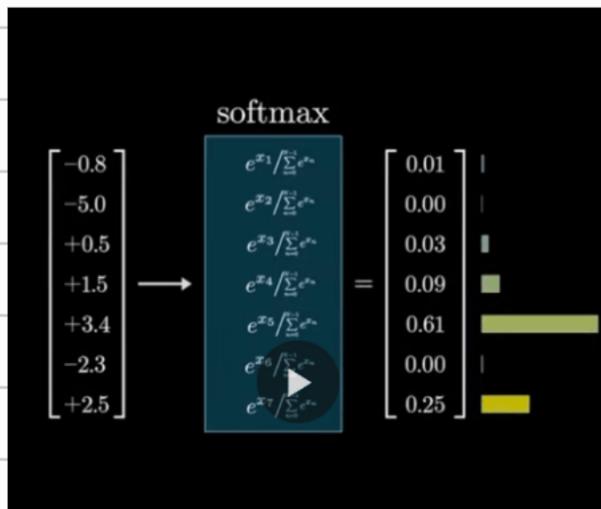


Step 3) Calculate $\frac{\partial L}{\partial \theta_i} = \left(\frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \dots, \frac{\partial L}{\partial \theta_c} \right)^T$

Remember the property of softmax that if you change any θ_u ALL \hat{y}_m where $m \in [1, c]$ (m includes i and k) will also change

Softmax

$$\hat{y}_m = \frac{e^{\theta_m}}{\sum_j e^{\theta_j}}$$



Therefore, the chain rule is $L \xrightarrow{\leftarrow} \hat{y}_i \rightarrow \theta_i$
 $\hat{y}_k \rightarrow \theta_i$

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial \theta_i} + \frac{\partial L}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial \theta_i}$$

Softmax derivative

$$\hat{y}_m = \frac{e^{\theta_m}}{\sum_j e^{\theta_j}}, \quad \text{let } S = \sum_j e^{\theta_j}$$

$$\hat{y}_m = \frac{e^{\theta_m}}{S}$$

Case 1

$$\frac{\partial \hat{y}_u^{<t>}}{\partial \theta_u^{<t>}} = \frac{e^{\theta_u^{<t>}} \cdot s - e^{\theta_u^{<t>}} \cdot e}{s^2} \\ = \frac{e^{\theta_u^{<t>}} (s - e^{\theta_u^{<t>}})}{s^2}$$

$$= \frac{e^{\theta_u^{<t>}}}{s} \cdot \left(1 - \frac{e^{\theta_u^{<t>}}}{s} \right)$$

$$\therefore \boxed{\frac{\partial \hat{y}_u^{<t>}}{\partial \theta_u^{<t>}} = \hat{y}_u^{<t>} \left(1 - \hat{y}_u^{<t>} \right)}$$

Case 2

$$\frac{\partial \hat{y}_v^{<t>}}{\partial \theta_u^{<t>}} \text{ where } v \neq u$$

$$\frac{\partial \hat{y}_v^{<t>}}{\partial \theta_u^{<t>}} = \frac{0 \cdot s - e^{\theta_v^{<t>}} \cdot e^{\theta_u^{<t>}}}{s^2} \\ = - \frac{e^{\theta_v^{<t>}}}{s} \cdot \frac{e^{\theta_u^{<t>}}}{s}$$

$$\therefore \boxed{\frac{\partial \hat{y}_v^{<t>}}{\partial \theta_u^{<t>}} = -\hat{y}_v^{<t>} \cdot \hat{y}_u^{<t>}}$$

From equations $\underline{2L} = -\frac{\hat{y}_i^{<t>}}{\hat{y}_i^{<t>}}$ (3) &

$$\frac{\underline{2L}^{<t>}}{\partial y_k^{<t>}} = \sum_{i \neq k} -\frac{\hat{y}_k^{<t>}}{\hat{y}_i^{<t>}}$$

(4) and from softmax

derivatives $\frac{\partial \hat{y}_i^{<t>}}{\partial \theta_i^{<t>}} = \hat{y}_i^{<t>} (1 - \hat{y}_i^{<t>})$ &

$$\frac{\partial \hat{y}_k^{<t>}}{\partial \theta_i^{<t>}} = -\hat{y}_k^{<t>} \cdot \hat{y}_i^{<t>}$$

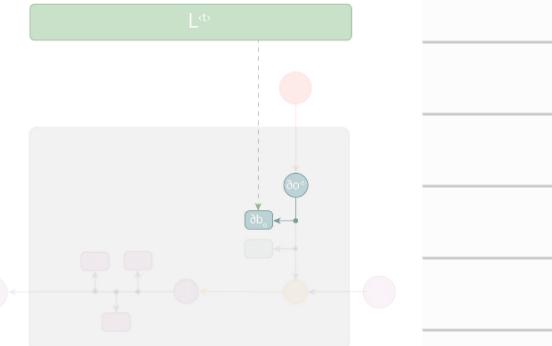
$$\begin{aligned}
 \frac{\partial L}{\partial \theta_i} &= \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial \theta_i} + \frac{\partial L}{\partial y_K} \cdot \frac{\partial y_K}{\partial \theta_i} \\
 &= \left(-\frac{y_i^{<t>}}{y_i^{<t>}} \cdot \hat{y}_i^{<t>} (1 - \hat{y}_i^{<t>}) \right) + \left(\sum_{i \neq K}^L \left(-\frac{y_K^{<t>}}{y_K^{<t>}} \right) \cdot -\hat{y}_K^{<t>} \cdot \hat{y}_i^{<t>} \right) \\
 &= -y_i^{<t>} (1 - \hat{y}_i^{<t>}) + \left(\sum_{i \neq K}^L y_K^{<t>} \cdot \hat{y}_i^{<t>} \right) \\
 &= -y_i^{<t>} + \hat{y}_i^{<t>} \cdot \hat{y}_i^{<t>} + \hat{y}_i^{<t>} \sum_{i \neq K}^L y_K^{<t>} \\
 &= -y_i^{<t>} + \hat{y}_i^{<t>} \left(y_i^{<t>} + \sum_{i \neq K}^L y_K^{<t>} \right) \\
 &\quad \underbrace{\qquad\qquad\qquad}_{=1 \text{ because } y_i^{<t>} \text{ is}} \\
 &\quad \text{a one-hot encoded vector}
 \end{aligned}$$

$$\therefore \boxed{\frac{\partial L}{\partial \theta_i} = \hat{y}_i^{<t>} - y_i^{<t>}} \quad (5)$$

Step 4) Calculate $\frac{\partial L}{\partial b} = \left(\frac{\partial L}{\partial b_1}, \frac{\partial L}{\partial b_2}, \dots, \frac{\partial L}{\partial b_C} \right)^T$

Chain-Rule

$$\begin{array}{c}
 L^{<t>} \rightarrow \hat{y}^{<t>} \rightarrow \hat{o}^{<t>} \xrightarrow{\frac{\partial \hat{o}^{<t>}}{\partial b_o}} b_o \\
 \frac{\partial L^{<t>}}{\partial b_o}
 \end{array}$$



$$\frac{\partial L^{(t)}}{\partial b_0} = \frac{2L^{(t)} \cdot 2y^{(t)} \cdot 2\theta^{(t)}}{2y^{(t)} \cdot 2\theta^{(t)} \cdot 2b_0} = \frac{2L^{(t)}}{2\theta^{(t)} \cdot 2b_0}$$

$\underbrace{\frac{\partial L^{(t)}}{\partial \theta^{(t)}}}_{2\theta^{(t)}}$

From $\sigma^{(t)} = (W_{ao} \cdot a^{(t)}) + b_0$ we want $\frac{\partial \theta^{(t)}}{\partial b_0}$

$$W_{ao} = \begin{pmatrix} w_{ao}^{11} & w_{ao}^{12} & \dots & w_{ao}^{1m} \\ w_{ao}^{21} & w_{ao}^{22} & \dots & w_{ao}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{ao}^{c1} & w_{ao}^{c2} & \dots & w_{ao}^{cm} \end{pmatrix} c \times m$$

$$a^{(t)} = \begin{pmatrix} a_1^{(t)} \\ a_2^{(t)} \\ a_3^{(t)} \\ \vdots \\ a_m^{(t)} \end{pmatrix}_{m \times 1}, \quad b_0 = \begin{pmatrix} b_{01} \\ b_{02} \\ b_{03} \\ \vdots \\ b_{0c} \end{pmatrix}_{c \times 1}, \quad \theta^{(t)} = \begin{pmatrix} \theta_1^{(t)} \\ \theta_2^{(t)} \\ \theta_3^{(t)} \\ \vdots \\ \theta_c^{(t)} \end{pmatrix}_{c \times 1}$$

$$\left. \begin{array}{l} \theta_1^{(t)} = W_{ao} \cdot a_1^{(t)} + W_{ao} \cdot a_2^{(t)} + \dots + W_{ao} \cdot a_m^{(t)} + b_{01} \\ \theta_2^{(t)} = W_{ao} \cdot a_1^{(t)} + W_{ao} \cdot a_2^{(t)} + \dots + W_{ao} \cdot a_m^{(t)} + b_{02} \\ \vdots \\ \theta_c^{(t)} = W_{ao} \cdot a_1^{(t)} + W_{ao} \cdot a_2^{(t)} + \dots + W_{ao} \cdot a_m^{(t)} + b_{0c} \end{array} \right\}$$

Let $n \in \mathbb{N}$ and $m \in [1, c]$

$$\frac{\partial \theta_n}{\partial b_m} = \frac{\partial}{\partial b_m} \left(W_{ao} \cdot a_1^{(t)} + W_{ao} \cdot a_2^{(t)} + \dots + W_{ao} \cdot a_m^{(t)} + b_{0n} \right)$$

$$= 1$$

$$\text{WLOG } \frac{\cancel{2\theta}^{<+>}}{\cancel{2b_n}^{<+>}} = 1 \text{ (unitary vector)} //$$

$$\rightarrow \frac{\cancel{2L}^{<+>}}{\cancel{2b_0}^{<+>}} = \frac{\cancel{2L}^{<+>}}{\cancel{2\theta}^{<+>}} \cdot \underbrace{\frac{\cancel{2\theta}^{<+>}}{\cancel{2b_0}^{<+>}}}_{1} : \quad \boxed{\frac{\cancel{2L}^{<+>}}{\cancel{2b_0}^{<+>}} = \frac{\cancel{2L}^{<+>}}{\cancel{2\theta}^{<+>}}} \quad (6)$$

Step 5) Calculate $\frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}}^{<+>}} = \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{11}}^{<+>}} \cdot \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{12}}^{<+>}} \cdots \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{1m}}^{<+>}}$

$L^{<+>}$

$\frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{21}}^{<+>}} \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{22}}^{<+>}} \cdots \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{2m}}^{<+>}}$

$\vdots \quad \vdots \quad \ddots \quad \vdots$

$\frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{c1}}^{<+>}} \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{c2}}^{<+>}} \cdots \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}^{cm}}^{<+>}}$

$C \times m$

Chain Rule

$$L^{<+>} \rightarrow \hat{y}^{<+>} \rightarrow o^{<+>} \xrightarrow{\frac{\cancel{2\theta}^{<+>}}{\cancel{2W_{ao}}^{<+>}}} W_{ao}$$

$\frac{\cancel{2L}^{<+>}}{\cancel{2\theta}^{<+>}}$

$$\therefore \frac{\cancel{2L}^{<+>}}{\cancel{2W_{ao}}^{<+>}} = \frac{\cancel{2L}^{<+>}}{\cancel{2\theta}^{<+>}} \cdot \frac{\cancel{2\theta}^{<+>}}{\cancel{2W_{ao}}^{<+>}} //$$

*From $\theta^{<+>} = (W_{ao} \cdot a^{<+>}) + b_o$ we want $\frac{\partial \theta}{\partial W_{ao}}$

$$W_{ao} = \begin{pmatrix} w_{ao}^{11} & w_{ao}^{12} & \dots & w_{ao}^{1m} \\ w_{ao}^{21} & w_{ao}^{22} & \dots & w_{ao}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{ao}^{c1} & w_{ao}^{c2} & \dots & w_{ao}^{cm} \end{pmatrix}_{C \times m}$$

$$a^{<+>} = \begin{pmatrix} a_1^{<+>} \\ a_2^{<+>} \\ a_3^{<+>} \\ \vdots \\ a_m^{<+>} \end{pmatrix}_{m \times 1}, \quad b_o = \begin{pmatrix} b_{o1} \\ b_{o2} \\ b_{o3} \\ \vdots \\ b_{oc} \end{pmatrix}_{C \times 1}, \quad \theta^{<+>} = \begin{pmatrix} \theta_1^{<+>} \\ \theta_2^{<+>} \\ \theta_3^{<+>} \\ \vdots \\ \theta_c^{<+>} \end{pmatrix}_{C \times 1}$$

$$\left\{ \begin{array}{l} \theta_1^{<+>} = W_{ao} \cdot a_1^{<+>} + W_{ao} \cdot a_2^{<+>} + \dots + W_{ao} \cdot a_m^{<+>} + b_{o1} \\ \theta_2^{<+>} = W_{ao} \cdot a_1^{<+>} + W_{ao} \cdot a_2^{<+>} + \dots + W_{ao} \cdot a_m^{<+>} + b_{o2} \\ \vdots \\ \theta_c^{<+>} = W_{ao} \cdot a_1^{<+>} + W_{ao} \cdot a_2^{<+>} + \dots + W_{ao} \cdot a_m^{<+>} + b_{oc} \end{array} \right.$$

Let $u, v \in \mathbb{N}$ where $u \in [1, c]$ and $v \in [1, m]$

$$\frac{\partial \theta_u}{\partial W_{au}} = \frac{\partial}{\partial W_{au}} \left(w_{au} \cdot a_1^{<+>} + w_{au} \cdot a_2^{<+>} + \dots + w_{au} \cdot a_v^{<+>} + \underbrace{w_{au} \cdot a_{u+1}^{<+>} + \dots + w_{au} \cdot a_m^{<+>}}_{a_u^{<+>}} + b_{ou} \right)$$

$$= a_{uv}^{<+>}$$

$$\text{WLOG } \frac{\partial \theta}{\partial W_{ao}} = a^{<+>} \quad \square$$

$$\frac{\frac{2L^{}}{2W_{ao}}}{2\theta^{}} = \frac{2L^{}}{2\theta^{}} \cdot \frac{2\theta^{}}{2W_{ao}} = \underbrace{\frac{2L^{}}{2\theta^{}}}_{C \times 1} \cdot \underbrace{\left(\begin{matrix} a^{} \end{matrix}\right)^T}_{m \times 1}$$

$\therefore \boxed{\frac{2L^{}}{2W_{ao}} = \frac{2L^{}}{2\theta^{}} \cdot a^{}^T} \quad (7)$

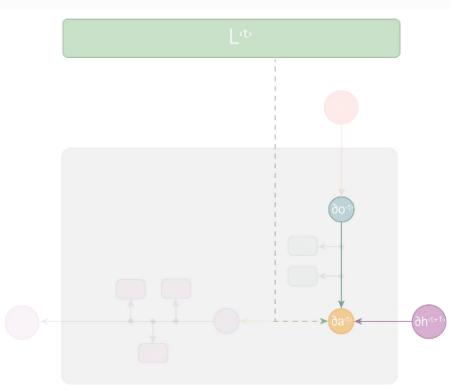
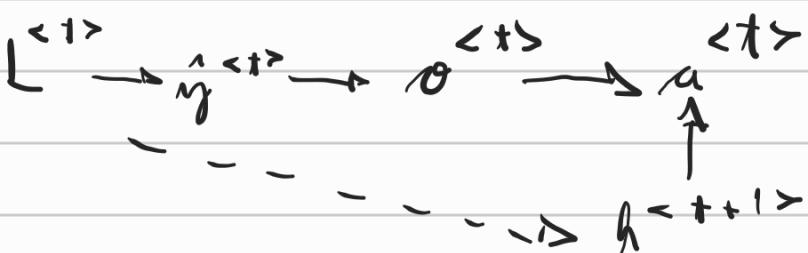
Note: $\frac{2L^{}}{2\theta^{}}, a^{}^T = \begin{pmatrix} \frac{2L^{}}{2\theta_1^{}} \\ \frac{2L^{}}{2\theta_2^{}} \\ \vdots \\ \frac{2L^{}}{2\theta_L^{}} \end{pmatrix} \cdot \left(a_1^{}, a_2^{}, \dots, a_m^{} \right)^T \quad 1 \times m$

Eg.: $\frac{2L^{}}{2W_{ao}} = \frac{2L^{}}{2\theta_1^{}} \cdot a_1^{}$

Step 6) Calculate $\frac{2L^{}}{2\theta^{}} = \begin{pmatrix} \frac{2L^{}}{2\theta_1^{}} & \frac{2L^{}}{2\theta_2^{}} & \dots & \frac{2L^{}}{2\theta_m^{}} \end{pmatrix}^T$

Chain Rule

Remember that $a^{}$ is used by $\theta^{}$ and $h^{}$ during forward propagation



Note

$h^{<+1>} \rightarrow$ contributes to $a^{<+>}$ and $a^{<+>}$ contributes to $L^{<+>} \rightarrow$, so we can derive $L^{<+>}$ in respect to $h^{<+1>} \rightarrow$ (this explain the dotted line)

$$\frac{\partial L^{<+>}}{\partial a^{<+>}} = \frac{\partial L^{<+>}}{\partial \theta^{<+>}} \cdot \frac{\partial \theta^{<+>}}{a^{<+>}} + \frac{\partial L^{<+>}}{\partial h^{<+1>}} \cdot \frac{\partial h^{<+1>}}{\partial a^{<+>}}$$

*From $\theta^{<+>} = (W_{\theta \theta} \cdot a^{<+>}) + b_\theta$ we want $\frac{\partial \theta^{<+>}}{\partial a^{<+>}}$

$$W_{\theta \theta} = \begin{pmatrix} w_{\theta \theta}^{11} & w_{\theta \theta}^{12} & \dots & w_{\theta \theta}^{1m} \\ w_{\theta \theta}^{21} & w_{\theta \theta}^{22} & \dots & w_{\theta \theta}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{\theta \theta}^{c1} & w_{\theta \theta}^{c2} & \dots & w_{\theta \theta}^{cm} \end{pmatrix} c \times m$$

$$a^{<+>} = \begin{pmatrix} a_1^{<+>} \\ a_2^{<+>} \\ a_3^{<+>} \\ \vdots \\ a_m^{<+>} \end{pmatrix}_{m \times 1}, \quad b_\theta = \begin{pmatrix} b_{\theta 1} \\ b_{\theta 2} \\ b_{\theta 3} \\ \vdots \\ b_{\theta c} \end{pmatrix}_{c \times 1}, \quad \theta^{<+>} = \begin{pmatrix} \theta_1^{<+>} \\ \theta_2^{<+>} \\ \theta_3^{<+>} \\ \vdots \\ \theta_c^{<+>} \end{pmatrix}_{c \times 1}$$

$$\left. \begin{array}{l} \theta_1^{<+>} = W_{\theta \theta} \cdot a_1^{<+>} + W_{\theta \theta} \cdot a_2^{<+>} + \dots + W_{\theta \theta} \cdot a_m^{<+>} + b_{\theta 1} \\ \theta_2^{<+>} = W_{\theta \theta} \cdot a_1^{<+>} + W_{\theta \theta} \cdot a_2^{<+>} + \dots + W_{\theta \theta} \cdot a_m^{<+>} + b_{\theta 2} \\ \vdots \\ \theta_c^{<+>} = W_{\theta \theta} \cdot a_1^{<+>} + W_{\theta \theta} \cdot a_2^{<+>} + \dots + W_{\theta \theta} \cdot a_m^{<+>} + b_{\theta c} \end{array} \right\}$$

Let $m, n \in \mathbb{N}$ where $m \in [1, c]$ and $n \in [1, m]$

$$\frac{\partial \alpha_m}{\partial a_n} = \frac{\partial}{\partial a_n} \left(w_{ao} a_1 + w_{ao} a_2 + \dots + w_{ao} a_n + \dots + w_{ao} a_m + b_o \right)$$

~~w_{ao}~~

$$= w_{ao}^m$$

$$\therefore \text{WLOG } \boxed{\frac{\partial \alpha}{\partial a} = w_{ao}} \quad \square$$

$$\rightarrow \text{From } h^{<t>} = w_{ah} \cdot x^{<t>} + w_{ah} \cdot a^{<t-1>} + b_h$$

$$\text{We want } \frac{\partial h^{<t>}}{\partial a^{<t-1>}}$$

$$w_{xh} = \begin{pmatrix} w_{xh}^{11} & w_{xh}^{12} & \dots & w_{xh}^{1m} \\ w_{xh}^{21} & w_{xh}^{22} & \dots & w_{xh}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{xh}^{m1} & w_{xh}^{m2} & \dots & w_{xh}^{mm} \end{pmatrix}_{m \times m}, \quad x^{<t>} = \begin{pmatrix} x_1^{<t>} \\ x_2^{<t>} \\ x_3^{<t>} \\ \vdots \\ x_n^{<t>} \end{pmatrix}_{m \times 1}$$

$$w_{ah} = \begin{pmatrix} w_{ah}^{11} & w_{ah}^{12} & \dots & w_{ah}^{1m} \\ w_{ah}^{21} & w_{ah}^{22} & \dots & w_{ah}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{ah}^{m1} & w_{ah}^{m2} & \dots & w_{ah}^{mm} \end{pmatrix}_{m \times m}, \quad a^{<t-1>} = \begin{pmatrix} a_1^{<t-1>} \\ a_2^{<t-1>} \\ a_3^{<t-1>} \\ \vdots \\ a_m^{<t-1>} \end{pmatrix}_{m \times 1}$$

$$h^{<t>} = \begin{pmatrix} h_1^{<t>} \\ h_2^{<t>} \\ \vdots \\ h_m^{<t>} \end{pmatrix}_{m \times 1}, \quad b_h = \begin{pmatrix} b_{h1} \\ b_{h2} \\ \vdots \\ b_{hm} \end{pmatrix}_{m \times 1}$$

$$\begin{cases} h_1^{<+>} = \left(w_{xh}^{11} \cdot x_1^{<+>} + w_{xh}^{12} \cdot x_2^{<+>} + \dots + w_{xh}^{1n} \cdot x_n^{<+>} \right) + \left(w_{ah}^{11} \cdot a_1^{<+>} + w_{ah}^{12} \cdot a_2^{<+>} + \dots + w_{ah}^{1m} \cdot a_m^{<+>} \right) + b_{h1} \\ h_2^{<+>} = \left(w_{xh}^{21} \cdot x_1^{<+>} + w_{xh}^{22} \cdot x_2^{<+>} + \dots + w_{xh}^{2n} \cdot x_n^{<+>} \right) + \left(w_{ah}^{21} \cdot a_1^{<+>} + w_{ah}^{22} \cdot a_2^{<+>} + \dots + w_{ah}^{2m} \cdot a_m^{<+>} \right) + b_{h2} \\ \vdots \\ h_m^{<+>} = \left(w_{xh}^{m1} \cdot x_1^{<+>} + w_{xh}^{m2} \cdot x_2^{<+>} + \dots + w_{xh}^{mn} \cdot x_n^{<+>} \right) + \left(w_{ah}^{m1} \cdot a_1^{<+>} + w_{ah}^{m2} \cdot a_2^{<+>} + \dots + w_{ah}^{mm} \cdot a_m^{<+>} \right) + b_{hm} \end{cases}$$

Let $u, v \in \mathbb{N}$, $u \in [1, m]$ and $v \in [1, n]$

$$\begin{aligned} \frac{\partial h_u^{<+>}}{\partial a_v^{<+>}} &= \frac{\partial}{\partial a_v^{<+>}} \left(\left(w_{xh}^{u1} \cdot x_1^{<+>} + w_{xh}^{u2} \cdot x_2^{<+>} + \dots + w_{xh}^{un} \cdot x_n^{<+>} \right) + \dots + \left(w_{ah}^{uv} \cdot a_v^{<+>} + \dots + w_{ah}^{um} \cdot a_m^{<+>} \right) + b_{hu} \right) \\ &+ \left(w_{ah}^{u1} \cdot a_1^{<+>} + w_{ah}^{u2} \cdot a_2^{<+>} + \dots + w_{ah}^{uv} \cdot a_v^{<+>} + \dots + w_{ah}^{um} \cdot a_m^{<+>} \right) + b_{ahu} \\ &= w_{ah}^{uve} \end{aligned}$$

$$\therefore \text{WLOG } \frac{\partial h_u^{<+>}}{\partial a_v^{<+>}} = w_{ah}^{uve} \quad \square$$

$$\xrightarrow{t+1} \frac{\partial h_u^{<+>}}{\partial a_v^{<+>}} = w_{ah}^{uve} \quad \square$$

$$\frac{\partial L^{<+>}}{\partial a_v^{<+>}} = \frac{\partial L^{<+>}}{\partial a_v^{<+>}} \cdot \frac{\partial a_v^{<+>}}{a_v^{<+>}} + \frac{\partial L^{<+>}}{\partial h_u^{<+>}} \cdot \frac{\partial h_u^{<+>}}{\partial a_v^{<+>}} \quad \therefore$$

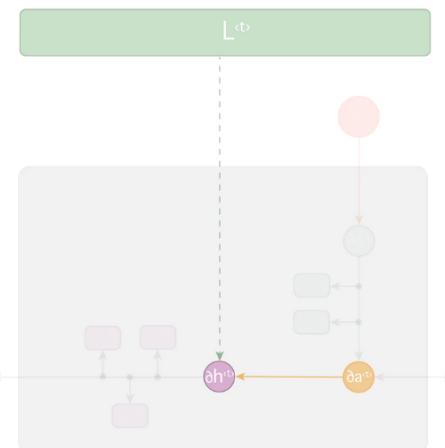
$$\frac{\partial L^{<+>}}{\partial a_v^{<+>}} = \underbrace{w_{ao}^T \cdot \frac{\partial L^{<+>}}{\partial a_v^{<+>}}}_{m \times c} + \underbrace{w_{ah}^T \cdot \frac{\partial L^{<+>}}{\partial h_u^{<+>}}}_{m \times m} \quad (8)$$

Step 7) Calculate $\frac{\partial L^{<t>}}{\partial h^{<t>}} = \left(\frac{\partial L^{<t>}}{\partial h_1^{<t>}}, \frac{\partial L^{<t>}}{\partial h_2^{<t>}}, \dots, \frac{\partial L^{<t>}}{\partial h_m^{<t>}} \right)^T$

Chain - Rule

$$L^{<t>} \rightarrow \hat{y}^{<t>} \rightarrow o^{<t>} \rightarrow a^{<t>} \rightarrow h^{<t>}$$

$$\frac{\partial L^{<t>}}{\partial h^{<t>}} = \frac{\partial L^{<t>}}{\partial o^{<t>}} \cdot \frac{\partial o^{<t>}}{\partial a^{<t>}} \cdot \frac{\partial a^{<t>}}{\partial h^{<t>}}$$



From $a^{<t>} = \tanh(h^{<t>})$ we want $\frac{\partial a^{<t>}}{\partial h^{<t>}}$

$$a^{<t>} = \begin{pmatrix} a_1^{<t>} \\ a_2^{<t>} \\ a_3^{<t>} \\ \vdots \\ a_m^{<t>} \end{pmatrix}_{m \times 1} \quad h^{<t>} = \begin{pmatrix} h_1^{<t>} \\ h_2^{<t>} \\ h_3^{<t>} \\ \vdots \\ h_m^{<t>} \end{pmatrix}_{m \times 1}$$

$$\begin{cases} a_1^{<t>} = \tanh(h_1^{<t>}) \\ a_2^{<t>} = \tanh(h_2^{<t>}) \\ \vdots \\ a_m^{<t>} = \tanh(h_m^{<t>}) \end{cases}$$

Let $u \in \mathbb{N}$ where $u \in [1, m]$

$$\frac{\partial a_u^{<t>}}{\partial h_u^{<t>}} = \frac{\partial}{\partial} \left(\tanh(h_u^{<t>}) \right)$$

$$= 1 - (a_u^{<t>})^2$$

$$\therefore \text{WLOG } \frac{\partial a^{<t>}}{\partial h^{<t>}} = 1 - (a^{<t>})^2$$

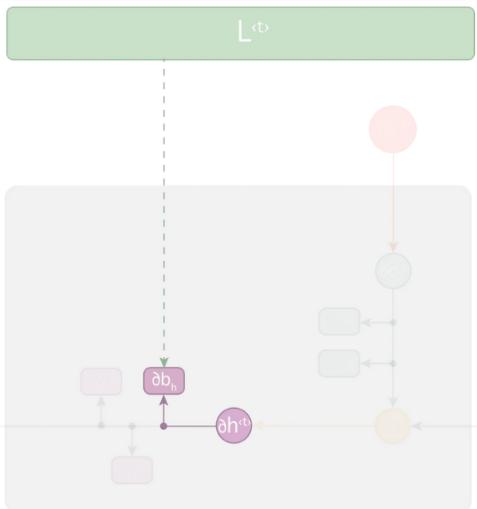
$$\boxed{\frac{\partial L^{<t>}}{\partial h^{<t>}} = \frac{\partial L^{<t>}}{\partial a^{<t>}} \cdot \left(1 - (a^{<t>})^2 \right)} \quad (9)$$

Step 8) Calculate $\frac{\partial L^{<t>}}{\partial b_h} = \left(\frac{\partial L^{<t>}}{\partial b_{h_1}}, \frac{\partial L^{<t>}}{\partial b_{h_2}}, \dots, \frac{\partial L^{<t>}}{\partial b_{h_m}} \right)$

Chain rule

$$L^{<t>} \rightarrow \hat{y}^{<t>} \rightarrow a^{<t>} \rightarrow h^{<t>} \rightarrow b_h$$

$$\frac{\partial L^{<t>}}{\partial b_h} = \frac{\partial L^{<t>}}{\partial h^{<t>}} \cdot \frac{\partial h^{<t>}}{\partial b_h}$$



From $h^{<t>} = W_{xh} \cdot x^{<t>} + W_{ah} \cdot a^{<t-1>} + b_h$

We want $\frac{\partial h^{<t>}}{\partial b_h}$

$$W_{xh} = \begin{pmatrix} w_{xh}^{11} & w_{xh}^{12} & \dots & w_{xh}^{1m} \\ w_{xh}^{21} & w_{xh}^{22} & \dots & w_{xh}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{xh}^{m1} & w_{xh}^{m2} & \dots & w_{xh}^{mm} \end{pmatrix}_{m \times n}, \quad x^{<t>} = \begin{pmatrix} x_1^{<t>} \\ x_2^{<t>} \\ x_3^{<t>} \\ \vdots \\ x_n^{<t>} \end{pmatrix}_{n \times 1}$$

$$W_{ah} = \begin{pmatrix} w_{ah}^{11} & w_{ah}^{12} & \dots & w_{ah}^{1m} \\ w_{ah}^{21} & w_{ah}^{22} & \dots & w_{ah}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{ah}^{m1} & w_{ah}^{m2} & \dots & w_{ah}^{mm} \end{pmatrix}_{m \times m}, \quad a^{<t-1>} = \begin{pmatrix} a_1^{<t-1>} \\ a_2^{<t-1>} \\ a_3^{<t-1>} \\ \vdots \\ a_m^{<t-1>} \end{pmatrix}_{m \times 1}$$

$$h^{<t>} = \begin{pmatrix} h_1^{<t>} \\ h_2^{<t>} \\ \vdots \\ h_m^{<t>} \end{pmatrix}_{m \times 1}, \quad b_h = \begin{pmatrix} b_{h_1} \\ b_{h_2} \\ \vdots \\ b_{h_m} \end{pmatrix}_{m \times 1}$$

$$\begin{cases} h_1^{<+>} = \left(w_{xh}^{11} \cdot x_1^{<+>} + w_{xh}^{12} \cdot x_2^{<+>} + \dots + w_{xh}^{1n} \cdot x_n^{<+>} \right) + \left(w_{ah}^{11} \cdot a_1^{<+>} + w_{ah}^{12} \cdot a_2^{<+>} + \dots + w_{ah}^{1m} \cdot a_m^{<+>} \right) + b_{h1} \\ h_2^{<+>} = \left(w_{xh}^{21} \cdot x_1^{<+>} + w_{xh}^{22} \cdot x_2^{<+>} + \dots + w_{xh}^{2n} \cdot x_n^{<+>} \right) + \left(w_{ah}^{21} \cdot a_1^{<+>} + w_{ah}^{22} \cdot a_2^{<+>} + \dots + w_{ah}^{2m} \cdot a_m^{<+>} \right) + b_{h2} \\ \vdots \\ h_m^{<+>} = \left(w_{xh}^{m1} \cdot x_1^{<+>} + w_{xh}^{m2} \cdot x_2^{<+>} + \dots + w_{xh}^{mn} \cdot x_n^{<+>} \right) + \left(w_{ah}^{m1} \cdot a_1^{<+>} + w_{ah}^{m2} \cdot a_2^{<+>} + \dots + w_{ah}^{mm} \cdot a_m^{<+>} \right) + b_{hm} \end{cases}$$

Let $n \in \mathbb{N}$ where $n \in [1, m]$

$$\begin{aligned} \frac{\partial h_n^{<+>}}{\partial b_{hn}} &= \frac{\partial}{\partial b_{hn}} \left(\left(w_{xh}^{n1} \cdot x_1^{<+>} + w_{xh}^{n2} \cdot x_2^{<+>} + \dots + w_{xh}^{nn} \cdot x_n^{<+>} \right) + \left(w_{ah}^{n1} \cdot a_1^{<+>} + w_{ah}^{n2} \cdot a_2^{<+>} + \dots + w_{ah}^{nm} \cdot a_m^{<+>} \right) + b_{hn} \right) \\ &+ \left(w_{ah}^{n1} \cdot a_1^{<+>} + w_{ah}^{n2} \cdot a_2^{<+>} + \dots + w_{ah}^{nm} \cdot a_m^{<+>} \right) + b_{hn} = 1 \end{aligned}$$

\downarrow

$$\therefore \text{WLOG } \frac{\partial h_n^{<+>}}{\partial b_{hn}} = 1 \text{ (unitary vector)}$$

\rightarrow $\frac{\partial L^{<+>}}{\partial b_n} = \frac{\partial L^{<+>}}{\partial h^{<+>}}$ (10)

Step 9) Compute $\frac{\partial L^{<t>}}{\partial W_{xh}}$

$$\frac{\partial L^{<t>}}{\partial W_{xh}} = \begin{pmatrix} \frac{\partial L}{\partial W_{xh}^{11}} & \frac{\partial L}{\partial W_{xh}^{12}} & \dots & \frac{\partial L}{\partial W_{xh}^{1n}} \\ \frac{\partial L}{\partial W_{xh}^{21}} & \frac{\partial L}{\partial W_{xh}^{22}} & \dots & \frac{\partial L}{\partial W_{xh}^{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial W_{xh}^{m1}} & \frac{\partial L}{\partial W_{xh}^{m2}} & \dots & \frac{\partial L}{\partial W_{xh}^{mn}} \end{pmatrix}_{m \times n}$$

Chain Rule

$$L^{<t>} \rightarrow \hat{y}^{<t>} \rightarrow a^{<t>} \rightarrow a^{<t>} \rightarrow h^{<t>} \rightarrow W_{xh}$$

$$\frac{\partial L^{<t>}}{\partial W_{xh}} = \frac{\partial L^{<t>}}{\partial h^{<t>}} \cdot \frac{\partial h^{<t>}}{\partial W_{xh}}$$

From $h^{<t>} = W_{xh} \cdot x^{<t>} + W_{ah} \cdot a^{<t-1>} + b_h$
 we want $\frac{\partial h^{<t>}}{\partial W_{xh}}$

$$W_{xh} = \begin{pmatrix} W_{xh}^{11} & W_{xh}^{12} & \dots & W_{xh}^{1n} \\ W_{xh}^{21} & W_{xh}^{22} & \dots & W_{xh}^{2n} \\ \vdots & \vdots & \ddots & \vdots \\ W_{xh}^{m1} & W_{xh}^{m2} & \dots & W_{xh}^{mn} \end{pmatrix}_{m \times n}, \quad x^{<t>} = \begin{pmatrix} x_1^{<t>} \\ x_2^{<t>} \\ x_3^{<t>} \\ \vdots \\ x_n^{<t>} \end{pmatrix}_{n \times 1}$$

$$W_{ah} = \begin{pmatrix} w_{ah}^{11} & w_{ah}^{12} & \dots & w_{ah}^{1m} \\ w_{ah}^{21} & w_{ah}^{22} & \dots & w_{ah}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{ah}^{m1} & w_{ah}^{m2} & \dots & w_{ah}^{mm} \end{pmatrix}_{m \times m}, \quad a^{\langle t-1 \rangle} = \begin{pmatrix} a_1^{\langle t-1 \rangle} \\ a_2^{\langle t-1 \rangle} \\ a_3^{\langle t-1 \rangle} \\ \vdots \\ a_m^{\langle t-1 \rangle} \end{pmatrix}_{m \times 1}$$

$$h^{\langle t \rangle} = \begin{pmatrix} h_1^{\langle t \rangle} \\ h_2^{\langle t \rangle} \\ \vdots \\ h_m^{\langle t \rangle} \end{pmatrix}_{m \times 1}, \quad b_h = \begin{pmatrix} b_{h1} \\ b_{h2} \\ \vdots \\ b_{hm} \end{pmatrix}_{m \times 1}$$

$$\left\{ \begin{array}{l} h_1^{\langle t \rangle} = (w_{xh}^{11} \cdot x_1^{\langle t \rangle} + w_{xh}^{12} \cdot x_2^{\langle t \rangle} + \dots + w_{xh}^{1m} \cdot x_m^{\langle t \rangle}) + (w_{ah}^{11} \cdot a_1^{\langle t-1 \rangle} + w_{ah}^{12} \cdot a_2^{\langle t-1 \rangle} + \dots + w_{ah}^{1m} \cdot a_m^{\langle t-1 \rangle}) + b_{h1} \\ h_2^{\langle t \rangle} = (w_{xh}^{21} \cdot x_1^{\langle t \rangle} + w_{xh}^{22} \cdot x_2^{\langle t \rangle} + \dots + w_{xh}^{2m} \cdot x_m^{\langle t \rangle}) + (w_{ah}^{21} \cdot a_1^{\langle t-1 \rangle} + w_{ah}^{22} \cdot a_2^{\langle t-1 \rangle} + \dots + w_{ah}^{2m} \cdot a_m^{\langle t-1 \rangle}) + b_{h2} \\ \vdots \\ h_m^{\langle t \rangle} = (w_{xh}^{m1} \cdot x_1^{\langle t \rangle} + w_{xh}^{m2} \cdot x_2^{\langle t \rangle} + \dots + w_{xh}^{mm} \cdot x_m^{\langle t \rangle}) + (w_{ah}^{m1} \cdot a_1^{\langle t-1 \rangle} + w_{ah}^{m2} \cdot a_2^{\langle t-1 \rangle} + \dots + w_{ah}^{mm} \cdot a_m^{\langle t-1 \rangle}) + b_{hm} \end{array} \right.$$

Let $m, n \in \mathbb{N}$ where $m \in [1, m]$ and $n \in [1, n]$

$$\begin{aligned} \frac{\partial h_n}{\partial w_{xh}^{mn}} &= \frac{\partial}{\partial w_{xh}^{mn}} \left((w_{xh}^{11} \cdot x_1^{\langle t \rangle} + w_{xh}^{12} \cdot x_2^{\langle t \rangle} + \dots + w_{xh}^{1m} \cdot x_m^{\langle t \rangle}) + (w_{ah}^{11} \cdot a_1^{\langle t-1 \rangle} + w_{ah}^{12} \cdot a_2^{\langle t-1 \rangle} + \dots + w_{ah}^{1m} \cdot a_m^{\langle t-1 \rangle}) + b_{hn} \right) \\ &+ \left((w_{xh}^{21} \cdot x_1^{\langle t \rangle} + w_{xh}^{22} \cdot x_2^{\langle t \rangle} + \dots + w_{xh}^{2m} \cdot x_m^{\langle t \rangle}) + (w_{ah}^{21} \cdot a_1^{\langle t-1 \rangle} + w_{ah}^{22} \cdot a_2^{\langle t-1 \rangle} + \dots + w_{ah}^{2m} \cdot a_m^{\langle t-1 \rangle}) + b_{hn} \right) \\ &= x_n^{\langle t \rangle} \end{aligned}$$

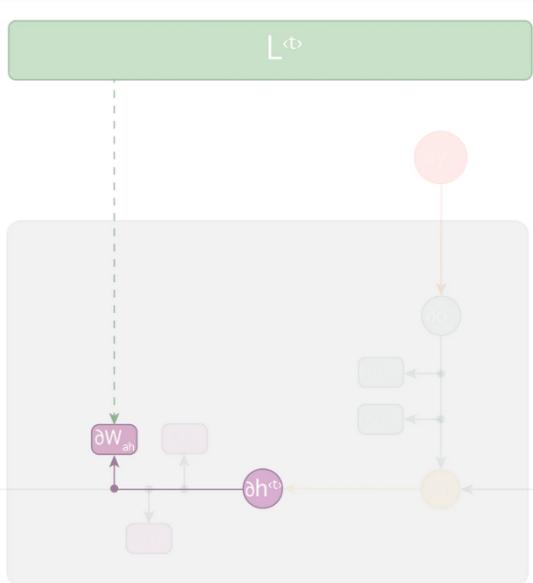
$$\therefore \text{WLOG } \frac{\partial h}{\partial w_{xh}}^{\langle t \rangle} = x^{\langle t \rangle}$$

$$\frac{\partial L^{(+)}}{\partial w_{xh}} = \frac{\partial l^{(+)}}{\partial w_{xh}} \cdot x^{(+)T} \quad (11)$$

Note :

$$\left(\begin{array}{c} \frac{\partial L}{\partial h_1} \\ \frac{\partial L}{\partial h_2} \\ \vdots \\ \frac{\partial L}{\partial h_m} \end{array} \right) \cdot \underbrace{\left(\begin{array}{c} x_1^{(t)} \\ x_2^{(t)} \\ \vdots \\ x_n^{(t)} \end{array} \right)}_{m \times 1} = \frac{\partial L}{\partial w_{12}} \quad m \times n$$

Step 10) Compute $\frac{\partial L}{\partial W_{ab}}^{(t)}$



$$\begin{array}{cccc}
 \frac{2L}{2w_{ah}^{m1}} & \frac{2L}{2w_{ah}^{m2}} & \dots & \frac{2L}{2w_{ah}^{mm}} \\
 \frac{2L}{2w_{ah}^{11}} & \frac{2L}{2w_{ah}^{12}} & & \\
 & & \ddots & \\
 \frac{2L}{2w_{ah}^{21}} & \frac{2L}{2w_{ah}^{22}} & \dots & \frac{2L}{2w_{ah}^{2m}} \\
 & & \vdots & \\
 & & & \vdots \\
 \frac{2L}{2w_{ah}^{m1}} & \frac{2L}{2w_{ah}^{m2}} & \dots & \frac{2L}{2w_{ah}^{mm}}
 \end{array}$$

Chain Rule

$$L^{<t>} \rightarrow \hat{y}^{<t>} \rightarrow o^{<t>} \rightarrow a^{<t>} \rightarrow h^{<t>} \rightarrow W_{ah}$$

$$\frac{\partial L^{<t>}}{\partial W_{ah}} = \frac{\partial L^{<t>}}{\partial h^{<t>}} \cdot \frac{\partial h^{<t>}}{\partial W_{ah}}$$

* From $h^{<t>} = W_{xh} \cdot x^{<t>} + W_{ah} * a^{<t>} + b_h$
 We want $\frac{\partial h^{<t>}}{\partial W_{ah}}$

$$W_{xh} = \begin{pmatrix} w_{xh}^{11} & w_{xh}^{12} & \dots & w_{xh}^{1m} \\ w_{xh}^{21} & w_{xh}^{22} & \dots & w_{xh}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{xh}^{m1} & w_{xh}^{m2} & \dots & w_{xh}^{mm} \end{pmatrix}_{m \times m}, \quad x^{<t>} = \begin{pmatrix} x_1^{<t>} \\ x_2^{<t>} \\ x_3^{<t>} \\ \vdots \\ x_m^{<t>} \end{pmatrix}_{m \times 1}$$

$$W_{ah} = \begin{pmatrix} w_{ah}^{11} & w_{ah}^{12} & \dots & w_{ah}^{1m} \\ w_{ah}^{21} & w_{ah}^{22} & \dots & w_{ah}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{ah}^{m1} & w_{ah}^{m2} & \dots & w_{ah}^{mm} \end{pmatrix}_{m \times m}, \quad a^{<t-1>} = \begin{pmatrix} a_1^{<t-1>} \\ a_2^{<t-1>} \\ a_3^{<t-1>} \\ \vdots \\ a_m^{<t-1>} \end{pmatrix}_{m \times 1}$$

$$h^{<t>} = \begin{pmatrix} h_1^{<t>} \\ h_2^{<t>} \\ \vdots \\ h_m^{<t>} \end{pmatrix}_{m \times 1}, \quad b_h = \begin{pmatrix} b_h_1 \\ b_h_2 \\ \vdots \\ b_h_m \end{pmatrix}_{m \times 1}$$

$$\left\{ \begin{array}{l} h_1 = (w_{xh}^{11} \cdot x_1 + w_{xh}^{12} \cdot x_2 + \dots + w_{xh}^{1m} \cdot x_m) + (w_{ah}^{11} \cdot a_1 + w_{ah}^{12} \cdot a_2 + \dots + w_{ah}^{1m} \cdot a_m) + b_{h1} \\ h_2 = (w_{xh}^{21} \cdot x_1 + w_{xh}^{22} \cdot x_2 + \dots + w_{xh}^{2m} \cdot x_m) + (w_{ah}^{21} \cdot a_1 + w_{ah}^{22} \cdot a_2 + \dots + w_{ah}^{2m} \cdot a_m) + b_{h2} \\ \vdots \\ h_m = (w_{xh}^{m1} \cdot x_1 + w_{xh}^{m2} \cdot x_2 + \dots + w_{xh}^{mm} \cdot x_m) + (w_{ah}^{m1} \cdot a_1 + w_{ah}^{m2} \cdot a_2 + \dots + w_{ah}^{mm} \cdot a_m) + b_{hm} \end{array} \right.$$

Let $m, n \in \mathbb{N}$ where $m, n \in [1, m]$

$$\begin{aligned} \frac{\partial h_n}{\partial w_{ah}^{mn}} &= \frac{\partial}{\partial w_{ah}^{mn}} \left((w_{xh}^{m1} \cdot x_1 + w_{xh}^{m2} \cdot x_2 + \dots + w_{xh}^{mn} \cdot x_n) + (w_{ah}^{m1} \cdot a_1 + w_{ah}^{m2} \cdot a_2 + \dots + w_{ah}^{mn} \cdot a_m) + b_{hn} \right) \\ &+ \left(\cancel{(w_{ah}^{m1} \cdot a_1 + w_{ah}^{m2} \cdot a_2 + \dots + w_{ah}^{mn} \cdot a_n)} + \dots + \cancel{(w_{ah}^{m1} \cdot a_1 + w_{ah}^{m2} \cdot a_2 + \dots + w_{ah}^{mn} \cdot a_m)} + b_{hn} \right) \\ &= a_n^{m-1} \\ \therefore \text{WLOG } \frac{\partial h}{\partial w_{ah}} &= a^{m-1} \end{aligned}$$

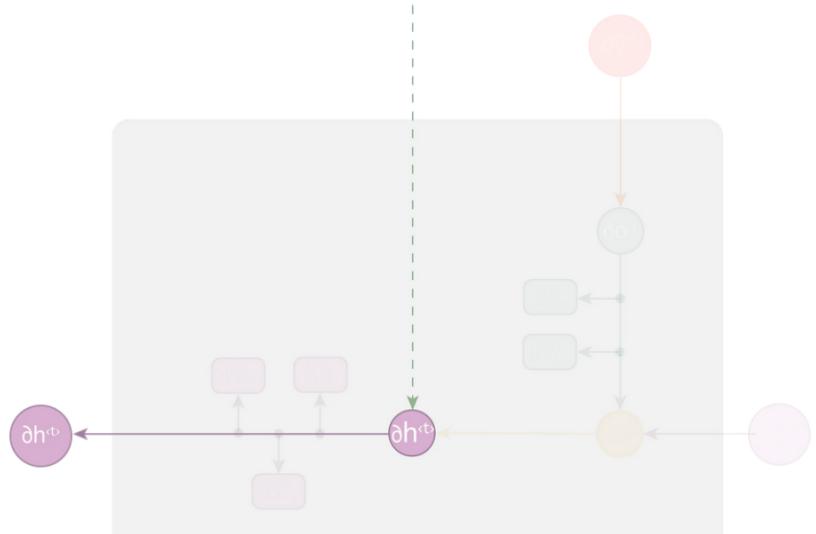
$$\boxed{\frac{\partial L}{\partial w_{ah}} = \frac{\partial L}{\partial h^{m-1}} \cdot a^{m-1}^T} \quad (12)$$

Step 11) Pass $\frac{\partial L}{\partial h^{m-1}}$ to previous cell

Therefore,

$$\boxed{\frac{\partial L}{\partial h} = \frac{\partial L}{\partial h^{m-1}}} \quad (13)$$

L^{fb}



Equations

Backward

$$\frac{\partial L^{<t>}}{\partial \theta^{<t>}} = \hat{y}^{<t>} - y^{<t>} \quad (\text{a})$$

$$\frac{\partial L^{<t>}}{\partial h^{<t+1>}} = \frac{\partial L^{<t+1>}}{\partial h^{<t+1>}} \quad (\text{b})$$

$$\frac{\partial L^{<t>}}{\partial a^{<t>}} = W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial \theta^{<t>}} + W_{ah}^T \cdot \frac{\partial L^{<t>}}{\partial h^{<t+1>}} \quad (\text{c})$$

$$\frac{\partial L^{<t>}}{\partial h^{<t>}} = \frac{\partial L^{<t>}}{\partial a^{<t>}} \cdot (1 - a^{<t>}^2) \quad (\text{d})$$

$$\frac{\partial L}{\partial W_{ao}} = \sum_{t=1}^T \frac{\partial L^{<t>}}{\partial \theta^{<t>}} \cdot (a^{<t>})^T \quad (\text{e})$$

$$\frac{\partial L}{\partial b_{ao}} = \sum_{t=1}^T \frac{\partial L^{<t>}}{\partial \theta^{<t>}} \quad (\text{f})$$

$$\frac{\partial L}{\partial W_{ah}} = \sum_{t=1}^T \frac{\partial L^{<t>}}{\partial h^{<t>}} \cdot (a^{<t-1>})^T \quad (\text{g})$$

$$\frac{\partial L}{\partial b_{ah}} = \sum_{t=1}^T \frac{\partial L^{<t>}}{\partial h^{<t>}} \cdot (x^{<t>})^T \quad (\text{h})$$

$$\frac{\partial L}{\partial b_h} = \sum_{t=1}^T \frac{\partial L^{<t>}}{\partial h^{<t>}} \quad (\text{i})$$

$$W_{ao} \leftarrow W_{ao} - \eta \frac{\partial L}{2W_{ao}}$$

$$b_o \leftarrow b_o - \eta \frac{\partial L}{2b_o}$$

$$W_{ah} \leftarrow W_{ah} - \eta \frac{\partial L}{2W_{ah}}$$

$$W_{xh} \leftarrow W_{xh} - \eta \frac{\partial L}{2W_{xh}}$$

$$b_h \leftarrow b_h - \eta \frac{\partial L}{2b_h}$$

$$L = \sum_{t=0}^T L^{<t>}$$

Forward

$$h^{<t>} = (W_{xh} \cdot x^{<t>}) + (W_{ah} \cdot a^{<t-1>}) + b_h \quad (a)$$

$$a^{<t>} = \tanh(h^{<t>}) \quad (b)$$

$$o^{<t>} = (W_{ao} \cdot a^{<t>}) + b_o \quad (c)$$

$$\hat{y}^{<t>} = \text{softmax}(o^{<t>}) \quad (d)$$

$x^{<t>}$ = Input vector in time t

$h^{<t>}$ = Hidden state

w_{xh} = Matrix that is multiplied by x and generates $h^{<t>}$

w_{ah} = Matrix that is multiplied by a and generates h

$a^{<t>}$ = Activation in time t

b^h = Hidden bias

$\theta^{<t>}$ = Non-normalized output in time t

w_{ao} = Matrix that is multiplied by a and generates o

b_o = Output bias

$\hat{o}^{<t>}$ = Normalized output in time t

$$w_{xh} = \begin{pmatrix} w_{xh}^{11} & w_{xh}^{12} & \dots & w_{xh}^{1n} \\ w_{xh}^{21} & w_{xh}^{22} & \dots & w_{xh}^{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{xh}^{m1} & w_{xh}^{m2} & \dots & w_{xh}^{mn} \end{pmatrix}_{m \times n}, \quad x^{<t>} = \begin{pmatrix} x_1^{<t>} \\ x_2^{<t>} \\ x_3^{<t>} \\ \vdots \\ x_n^{<t>} \end{pmatrix}_{n \times 1}$$

$$w_{ah} = \begin{pmatrix} w_{ah}^{11} & w_{ah}^{12} & \dots & w_{ah}^{1m} \\ w_{ah}^{21} & w_{ah}^{22} & \dots & w_{ah}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ w_{ah}^{m1} & w_{ah}^{m2} & \dots & w_{ah}^{mm} \end{pmatrix}_{m \times m}, \quad a^{<t>} = \begin{pmatrix} a_1^{<t>} \\ a_2^{<t>} \\ a_3^{<t>} \\ \vdots \\ a_m^{<t>} \end{pmatrix}_{m \times 1}$$

$$h^{<t>} = \begin{pmatrix} h_1^{<t>} \\ h_2^{<t>} \\ \vdots \\ h_m^{<t>} \end{pmatrix}_{m \times 1}, \quad b_h = \begin{pmatrix} b_h^1 \\ b_h^2 \\ \vdots \\ b_h^m \end{pmatrix}_{m \times 1}$$

$$W_{AO} = \begin{pmatrix} W_{AO}^{11} & W_{AO}^{12} & \dots & W_{AO}^{1m} \\ W_{AO}^{21} & W_{AO}^{22} & \dots & W_{AO}^{2m} \\ \vdots & \vdots & \ddots & \vdots \\ W_{AO}^{c1} & W_{AO}^{c2} & \dots & W_{AO}^{cm} \end{pmatrix}_{C \times m}$$

$$b_O = \begin{pmatrix} b_{O1} \\ b_{O2} \\ b_{O3} \\ \vdots \\ b_{Oc} \end{pmatrix}_{C \times 1}, \quad O^{<t>} = \begin{pmatrix} O_1^{<t>} \\ O_2^{<t>} \\ O_3^{<t>} \\ \vdots \\ O_c^{<t>} \end{pmatrix}_{C \times 1} \quad | \quad \hat{y}^{<t>} = \begin{pmatrix} \hat{y}_1^{<t>} \\ \hat{y}_2^{<t>} \\ \hat{y}_3^{<t>} \\ \vdots \\ \hat{y}_c^{<t>} \end{pmatrix}_{C \times 1}$$

Exploding & Vanishing gradient

$$\begin{cases} \frac{\partial L^{<t>}}{\partial a^{<t>}} = W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial o^{<t>}} + W_{ah}^T \cdot \frac{\partial L^{<t+1>}}{\partial h^{<t+1>}} \\ \frac{\partial L^{<t>}}{\partial h^{<t>}} = \frac{\partial L^{<t>}}{\partial a^{<t>}} \cdot (1 - a^{<t>^2}) \end{cases}$$

$$\frac{\partial L^{<t>}}{\partial a^{<t>}} = W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial o^{<t>}} + W_{ah}^T \left(\frac{\partial L^{<t+1>}}{\partial a^{<t+1>}} \cdot (1 - a^{<t+1>^2}) \right)$$

$$= W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial o^{<t>}} + W_{ah}^T \left(1 - a^{<t+1>^2} \right) \cdot \frac{\partial L^{<t+1>}}{\partial a^{<t+1>}}$$

$$= W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial o^{<t>}} + W_{ah}^T \left(1 - a^{<t+1>^2} \right) \cdot \left(W_{ao}^T \cdot \frac{\partial L^{<t+1>}}{\partial o^{<t+1>}} + W_{ah}^T \cdot \frac{\partial L^{<t+2>}}{\partial h^{<t+2>}} \right)$$

$$= W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial o^{<t>}} + W_{ao}^T \cdot \frac{\partial L^{<t+1>}}{\partial o^{<t+1>}} \left[W_{ah}^T \left(1 - a^{<t+1>^2} \right) \right] + \left(W_{ah}^T \right)^2 \left(1 - a^{<t+1>^2} \right) \cdot \frac{\partial L^{<t+2>}}{\partial h^{<t+2>}}$$

$$= W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial o^{<t>}} + W_{ao}^T \cdot \frac{\partial L^{<t+1>}}{\partial o^{<t+1>}} \left[W_{ah}^T \left(1 - a^{<t+1>^2} \right) \right] + \left(W_{ah}^T \right)^2 \left(1 - a^{<t+1>^2} \right) \left(1 - a^{<t+2>^2} \right) \cdot \frac{\partial L^{<t+3>}}{\partial a^{<t+3>}}$$

$$= W_{ao}^T \cdot \frac{\partial L^{<t>}}{\partial o^{<t>}} + W_{ao}^T \cdot \frac{\partial L^{<t+1>}}{\partial o^{<t+1>}} \left[W_{ah}^T \left(1 - a^{<t+1>^2} \right) \right] + \left(W_{ah}^T \right)^2 \left(1 - a^{<t+1>^2} \right) \left(1 - a^{<t+2>^2} \right) \cdot$$

$$\left(W_{ao}^T \cdot \frac{\partial L^{<t+2>}}{\partial o^{<t+2>}} + W_{ah}^T \cdot \frac{\partial L^{<t+3>}}{\partial h^{<t+3>}} \right)$$

$$= W_{ao}^T \cdot \frac{\partial L^{<+>}}{\partial a^{<+>}} + W_{ao}^T \cdot \frac{\partial L^{<++>}}{\partial a^{<++>}} \left[W_{ah}^T (1 - a^{<++>^2}) \right] + \frac{W_{ao}^T \cdot \partial L^{<++>}}{\partial a^{<++>}} \left[\left(W_{ah}^T \right)^2 (1 - a^{<++>^2}) (1 - a^{<++>^2}) \right]$$

$$+ (W_{ah}^T)^3 (1 - a^{<++>^2}) (1 - a^{<++>^2}) (1 - a^{<++>^2}) \cdot \frac{\partial L^{<++3>}}{\partial a^{<++3>}}$$

⋮

$$\frac{\partial L^{<+>}}{\partial a^{<+>}} = W_{ao}^T \left[\frac{\partial L^{<+>}}{\partial a^{<+>}} + \sum_{i=1}^{T'} \frac{\partial L^{<++i>}}{\partial a^{<++i>}} \cdot (W_{ah}^T) \cdot \prod_{j=1}^{i-1} 1 - (a^{<++j>})^2 \right]$$

$$+ (W_{ah}^T)^{T'} \prod_{i=1}^{T'} 1 - (a^{<++i>})^2$$

Where T' is the last timestamp

* Note that if T' is too big and :

1) W_{ah} and W_{ao} have large values, so $\frac{\partial L^{<+>}}{\partial a^{<+>}}$ will be high and the gradient will overflow (explode)

2) W_{ah} and W_{ao} have small values, so $\frac{\partial L^{<+>}}{\partial a^{<+>}}$ will be high and the gradient will underflow (vanish)

Therefore, if a RNN has much cells, it will be hard to learn due to the exploding and vanishing gradient problem

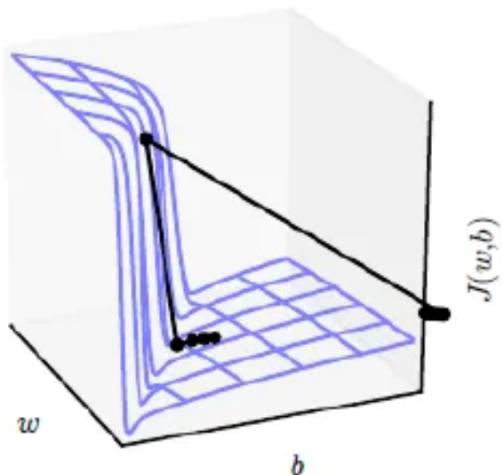
Gradient Clipping: Places a predefined threshold on gradients to prevent them from getting too large. This approach does not change the direction of the gradient, only change its length

if $\|g\| > \text{threshold}$:

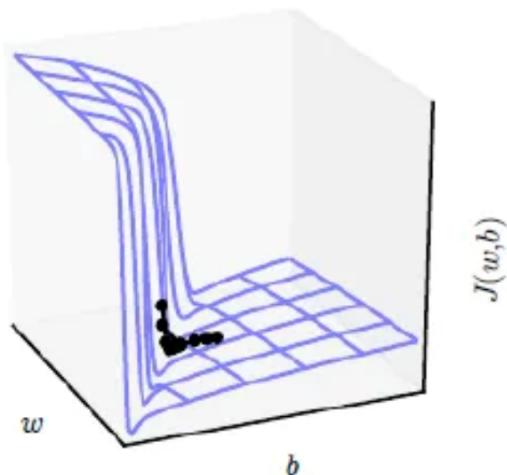
$$g \leftarrow \frac{\text{threshold} \cdot g}{\|g\|}$$

Where g is the gradient and $\|g\|$ is the norm of the gradient

Without clipping



With clipping



Identity RNN: All weights are initialized to identity matrix and all activation function are set to ReLU. This forces the network to stay close to identity function and the problem of vanishing/exploding gradient is solved