# Natural Language Processing
## CSE 325/425

Sihong Xie

<u>Lecture 13:</u>

- Recurrent neural networks (RNN)

(Not Recursive)
for trees

# Language model review

Bi-gram $\quad P(w_t|w_{t-1}) = \dfrac{\text{Count } w_t}{\text{Count}[w_{t-1}, w_t]}$

$C(\text{``race''})$

n-gram $\quad P(w_t|[w_{t-1}, \ldots, w_{t-n+1}]) = \dfrac{\text{Count }[w_t, \ldots, w_{t-n+1}] + \varepsilon}{\text{Count }[w_{t-1}, \ldots, w_{t-n+1}] + |V|\varepsilon}$

$\geq C(\text{``race a''})$

$\geq C(\text{``race a car''})$

$\geq$

Two issues

Laplacian Smooth

- Data sparsity: the occurrences of many $[w_t, \ldots, w_{t-n+1}]$ are zeros.

- Model complexity: number of parameters increases exponentially in *n*.

**The two issues are related**: if we want longer range dependencies, we increase *n*,

then both the data sparsity and model complexity become worse.

*The students walked in the room and asked the ___?___ about the quiz questions.*

$w?$

RNN:
$Pr(w \mid \text{``The students --- --- asked the''}; \bar{\theta})$

tri-gram: $Pr(w \mid \text{``asked the''})$

If Count (``asked the'') $= 0$ on training corpus

# Address the issues using neural network

Don't store the *n*-grams, but use a fixed-size model to predict the *n*-grams.

- model complexity is fixed.
- no data sparsity issue (no *n*-gram is computed)
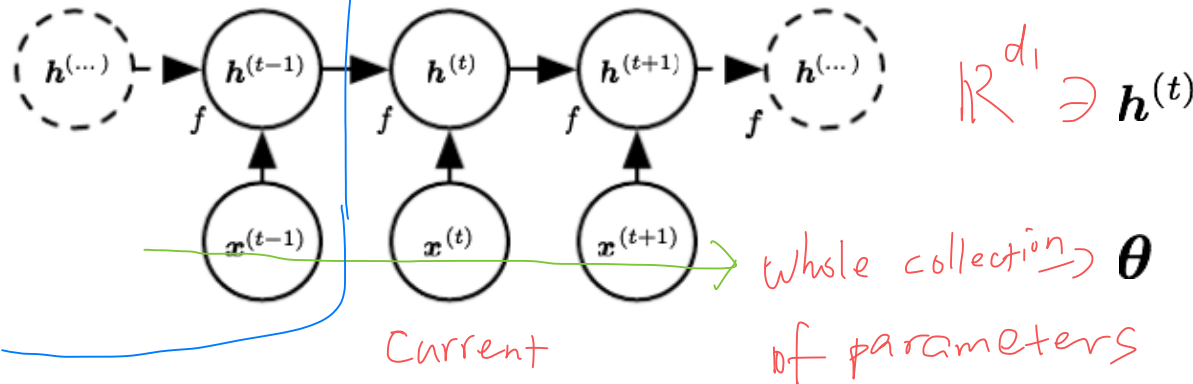- can be generalized to unseen sequences.

**Recurrent Neural Networks (RNN)**

*history before* $x^{(t)}$

$$= f\big(f(h^{(t-2)}, x^{(t-1)}; \theta)$$

$$h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta) \quad x^{(t)}; \theta)$$

"*Recurrent*" because the next **h** is compute by the same function that calculates the previous **h**.



$\mathbb{R}^{d_1} \ni h^{(t)}$

state summarizing what has happened before time step *t*. (theoretically speaking)

*whole collection* $\theta$

*of parameters*

a single model specifying how to transit to the next state (independent of *t*).

*Current Input*

$\theta$ *is given*

$Pr(W_t \mid W_{t-1} \cdots W_{t-n+1})$

$d_0, d_1$ : *Integers of your choice*

$x^{(t)} \in \mathbb{R}^{d_0}$ (*e.g. word embedding by Glove for the t-th word*)

# A running example

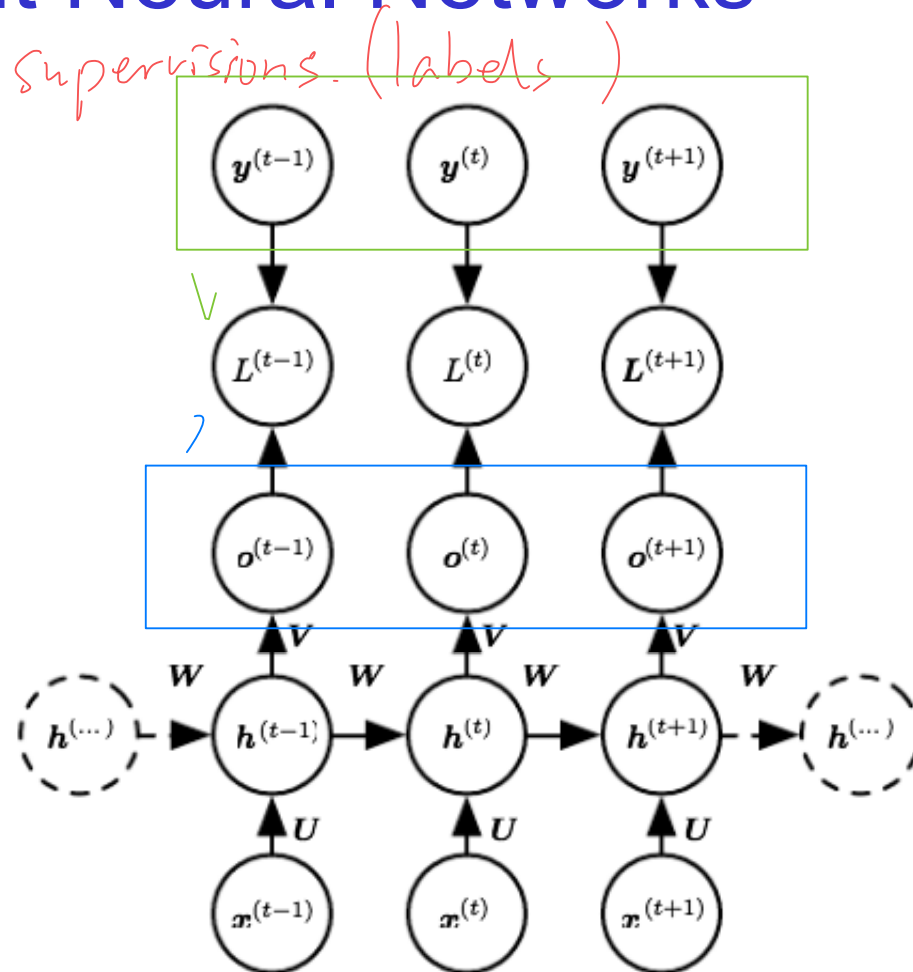# Recurrent Neural Networks

Supervisions. (labels)

Output sequence
(e.g., POS tags)

Loss function

Output units

Hidden states

Input sequence
(e.g., sentence)



**Training data:**

$$\{x^{(1)}, \ldots, x^{(\tau)}\}, \{y^{(1)}, \ldots, y^{(\tau)}\}$$

$\tau$ : /tau/

**Trainable parameters:**

$$\theta = \{U, W, V, b, c\}$$

$V$ : maps from $h$ to $o$

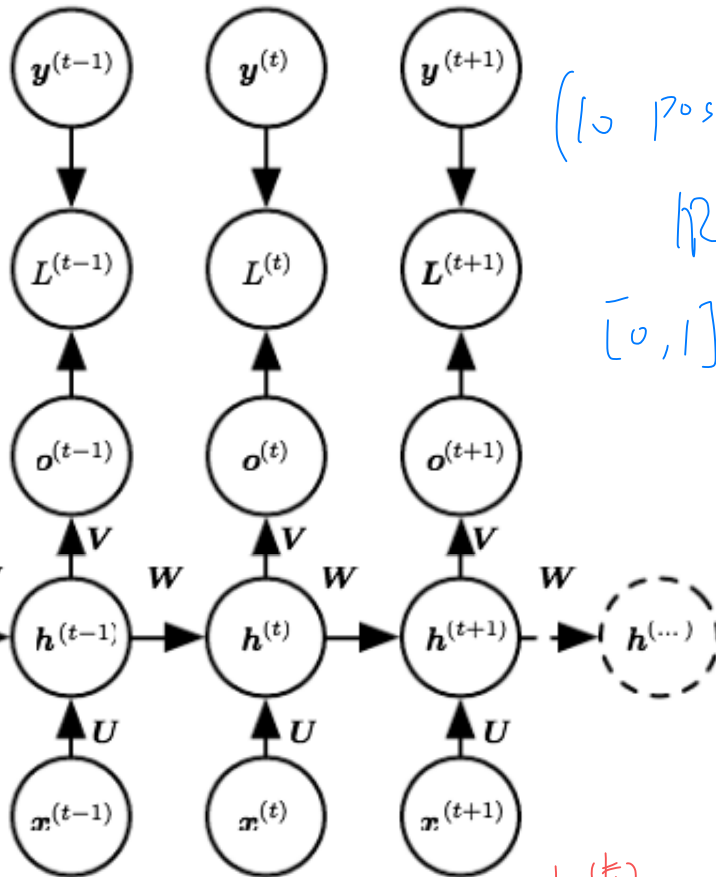$W$ : maps from $h$ to $h$

$U$ : maps from $x$ to $h$

$b, c$ : biases

# RNN forward pass

ground Truth $\vec{y}^{(t)}$

Prob mass

one hot encoding

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)},$$
(10 pos tags) $$h^{(t)} = \tanh(a^{(t)}),$$ activation function
$\mathbb{R}^{10} \ni$ $$o^{(t)} = c + Vh^{(t)},$$
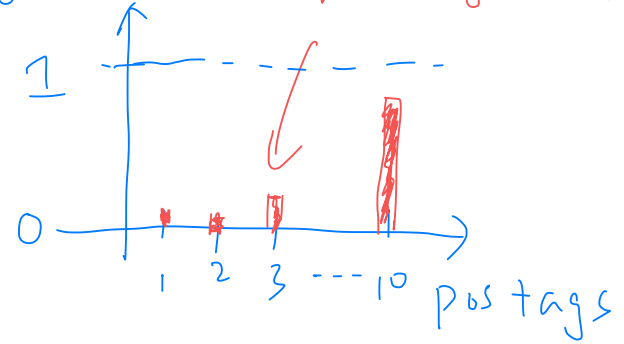$[0,1]^{10} \ni$ $$\hat{y}^{(t)} = \text{softmax}(o^{(t)}),$$ probability distribution

prediction $\hat{\vec{y}}^{(t)}$

$\text{Prob}(y^{(t)} = 3 | \cdots)$

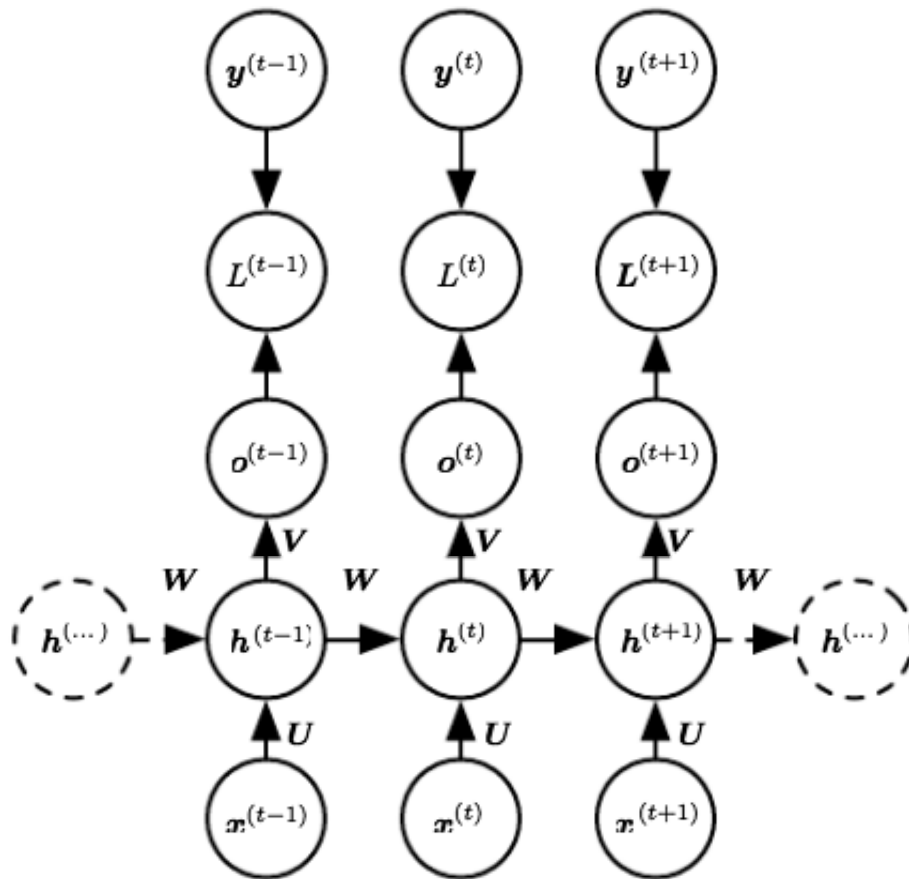$$L^{(t)} = L(\vec{y}^{(t)}, \hat{\vec{y}}^{(t)}) = -\log \Pr(y^{(t)} = 3 | \cdots)$$

# RNN forward pass



Negative log likelihood (NLL) loss, or the "perplexity"

$$L\left(\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(\tau)}\}, \{\boldsymbol{y}^{(1)}, \ldots, \boldsymbol{y}^{(\tau)}\}\right)$$

$$= \sum_t L^{(t)}$$

the ground truth label.

$$= -\sum_t \log p_{\text{model}}\left(y^{(t)} \mid \{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(t)}\}\right)$$

$$\hat{\boldsymbol{y}}^{(t)} = \text{softmax}(\boldsymbol{o}^{(t)})$$

# A running example

# RNN back propagation



**BPTT** (Back Propagation through Time)

- Used for gradient descent training;

- A special name for RNN back-propagation;

- Need all information in the forward pass, making BPTT sequential and hard to parallelize.

- $\theta = \{U, W, V, b, c\}$ used in all steps.
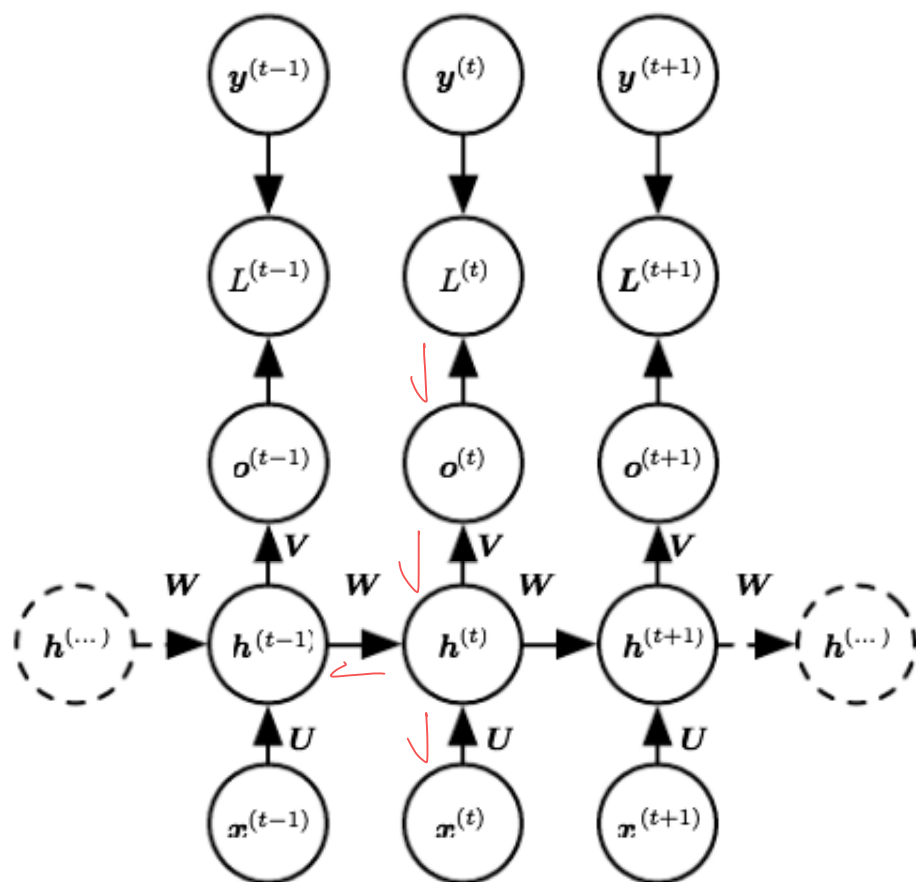
- Two derivative rules applied:

$$\nabla_x(f(x) + g(x)) = \nabla_x f(x) + \nabla_x g(x)$$

$$\nabla_x(f(g(x))) = \nabla_g f(g(x)) \times \nabla_x g(x)$$

Jacobian    vector

time

# RNN back propagation



$$L = \sum_{t=1}^{T} L^{(t)} = \text{Sum of all local losses}$$

**BPTT** (Back Propagation through Time)

- focus on each step t (the final gradient is the sum of all gradients at each step.

$$\frac{\partial L}{\partial L^{(t)}} = 1$$

$$\begin{bmatrix} 0.1 \\ 0.2 \\ \vdots \\ 0.01 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$(\nabla_{o^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i=y^{(t)}}$$

(an element of $\nabla_{o^{(t)}} L$ )

*error vector over tags.*

$$L^{(t)} = -\log Pr(y^{(t)} \mid x^{(1)}, \ldots, x^{(t)}; \theta)$$

$$= -\log \text{softmax}(o^{(t)})$$

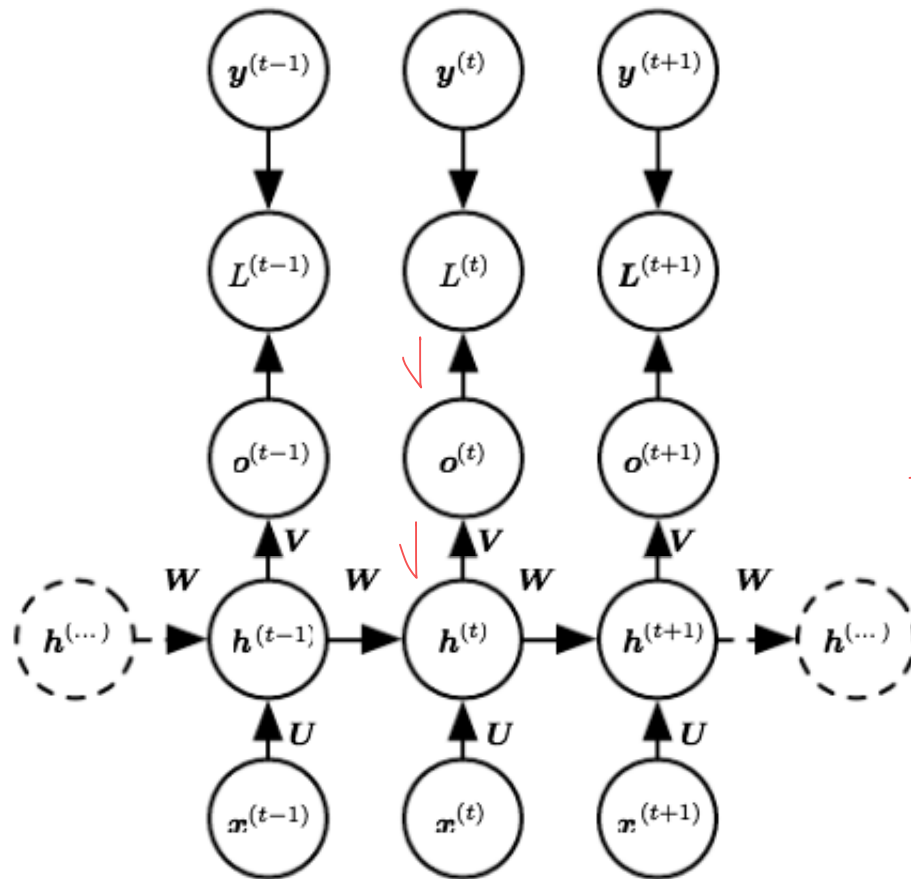$$= -\log \frac{\exp(o^{(t)}_{\hat{y}^{(t)}})}{\sum_{j=1}^{10} \exp(o^{(t)}_j)}$$

# RNN backprop (base case)



**BPTT** (Back Propagation through Time)

- Base case: at the final step $\tau$

$$\nabla_{\boldsymbol{h}^{(\tau)}} L = \boldsymbol{V}^\top \nabla_{\boldsymbol{o}^{(\tau)}} L$$

since $\boldsymbol{o}^{(\tau)} = \boldsymbol{V}\boldsymbol{h}^{(\tau)} + \boldsymbol{c}$

Handwritten annotations:

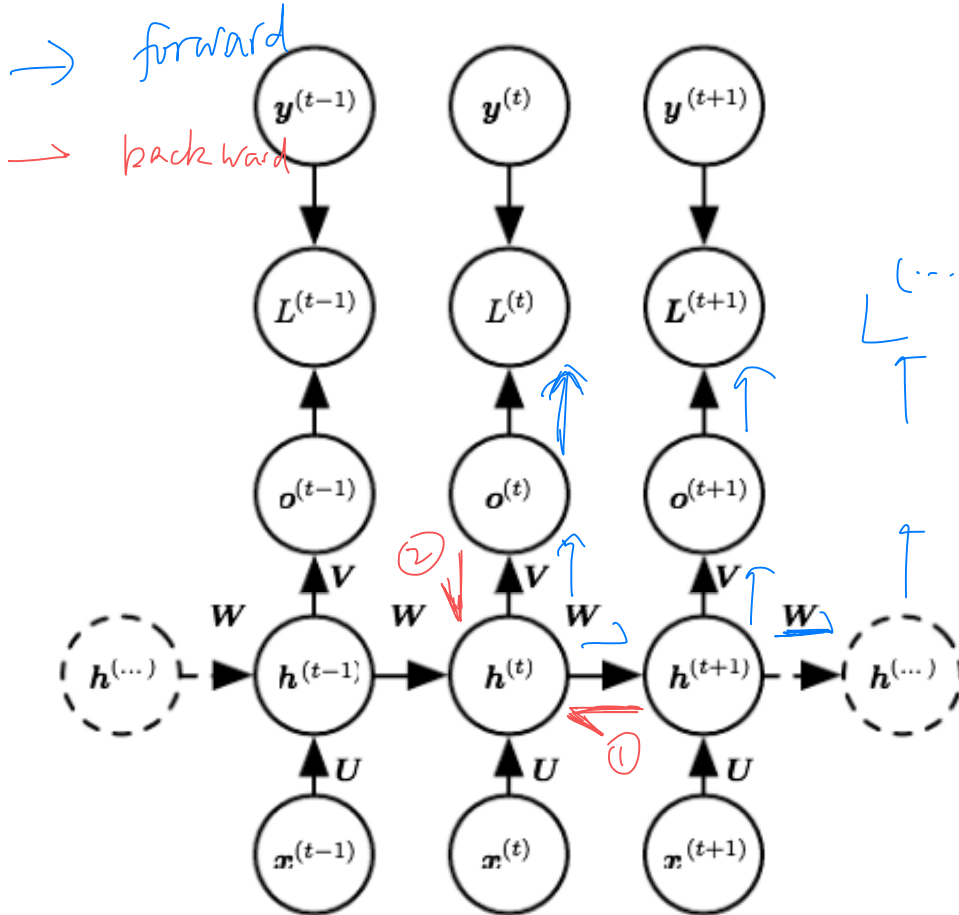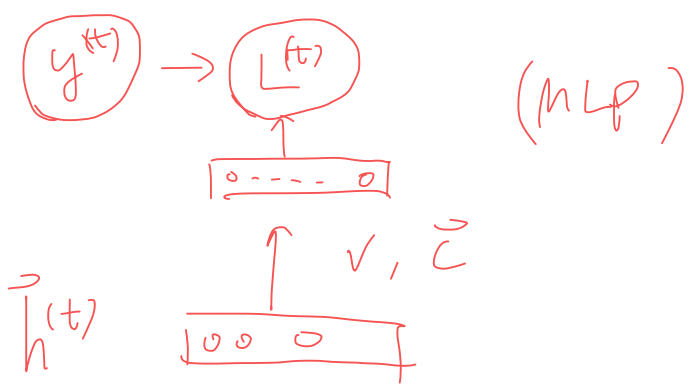$$\mathbb{R}^2 \ni \frac{\partial L}{\partial h^{(\tau)}} = \begin{bmatrix} \frac{\partial L}{\partial h_1^{(\tau)}} \\ \frac{\partial L}{\partial h_2^{(\tau)}} \end{bmatrix} = \begin{bmatrix} V_{:1}^T \frac{\partial L}{\partial o^{(\tau)}} \\ V_{:2}^T \frac{\partial L}{\partial o^{(\tau)}} \end{bmatrix} = V^T \frac{\partial L}{\partial o^{(\tau)}}$$

$$= V^T (\hat{y}^{(t)} - y^{(t)})$$

$$\frac{\partial L}{\partial h_1^{(\tau)}} = \sum_{j=1}^{|o|} V_{j1} \frac{\partial h}{\partial o^{(\tau)}_j} = V_{:1}^T \frac{\partial L}{\partial o^{(\tau)}}$$

# RNN backprop (recurrent)

$$\vec{o}^{(t)} = softmax(\vec{c} + V\vec{h}^{(t)})$$

$$y^{(t)} \rightarrow L^{(t)} \qquad (MLP)$$

$$V, \vec{c}$$

$$\vec{h}^{(t)}$$

→ forward

→ backward



**BPTT** (Back Propagation through Time)

- Recursively, at any step $1 \le t < \tau$

$$\nabla_{\boldsymbol{h}^{(t)}} L = \left(\frac{\partial \boldsymbol{h}^{(t+1)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \left(\frac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$

$$= \boldsymbol{W}^{\top} \text{diag}\left(1 - \left(\boldsymbol{h}^{(t+1)}\right)^2\right) (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \boldsymbol{V}^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$

$$\boldsymbol{o}^{(t)} = V\boldsymbol{h}^{(t)} + \boldsymbol{c}$$

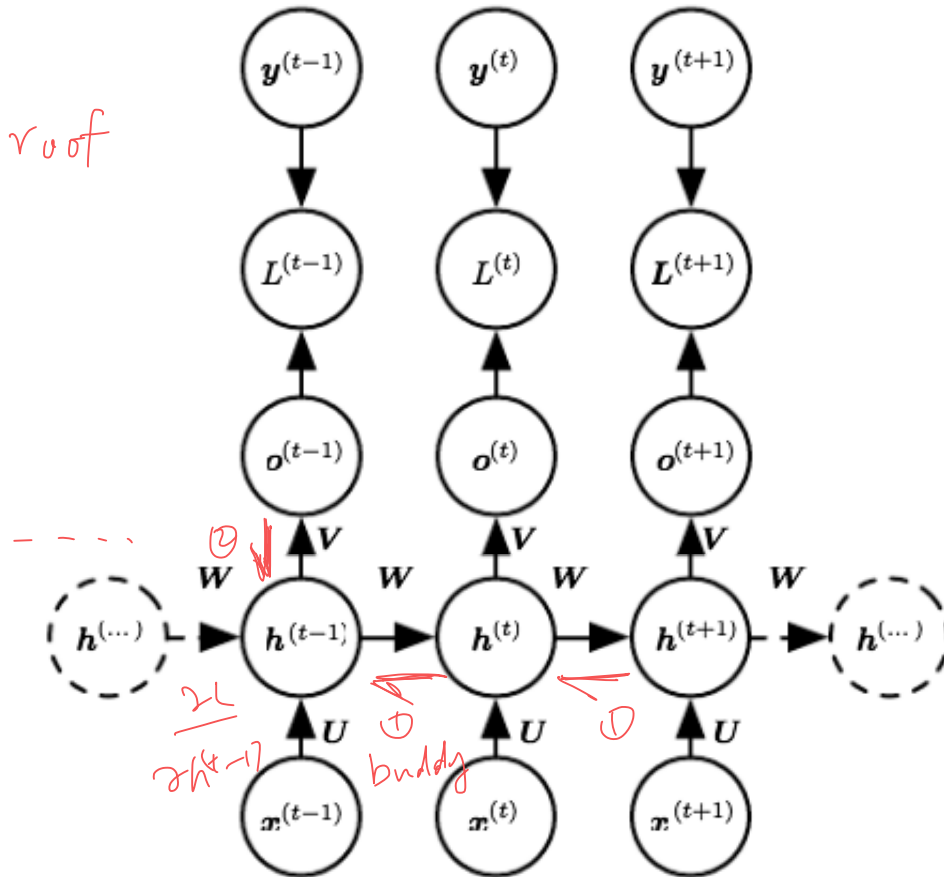$$\frac{\partial L^{(t)}}{\partial O^{(t)}} = \hat{\vec{y}}^{(t)} - \vec{y}^{(t)} \quad (\text{error vector})$$

$$0 \le t < \tau$$

$$\frac{\partial L^{(t)}}{\partial h^{(t)}} = V^{\top} (\hat{\vec{y}}^{(t)} - \vec{y}^{(t)})$$

# RNN backprop (recurrent)

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

$$= 1 - \tanh^2(x)$$

$$\text{diag}\left(\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} \tanh^2(h_1^{(t+1)}) \\ \tanh^2(h_2^{(t+1)}) \\ \tanh^2(h_3^{(t+1)}) \end{bmatrix}\right)$$

$$\mathbb{1} - h^{2(t+1)}$$

$$= \begin{bmatrix} 1 - \tanh^2(h_1^{(t+1)}) & 0 & 0 \\ 0 & 1 - \tanh^2(h_2^{(t+1)}) & \vdots \\ 0 & & \end{bmatrix} \quad 1 - \tanh h^2$$



**BPTT** (Back Propagation through Time)

- Recursively, at any step $1 \le t < \tau$

for $t = \tau - 1, \cdots, 1$

$$\nabla_{\boldsymbol{h}^{(t)}} L = \left(\frac{\partial \boldsymbol{h}^{(t+1)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{h}^{(t+1)}} L) + \left(\frac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{h}^{(t)}}\right)^{\top} (\nabla_{\boldsymbol{o}^{(t)}} L)$$

$$= \boldsymbol{W}^{\top} \text{diag}\left(1 - \left(\boldsymbol{h}^{(t+1)}\right)^2\right)(\nabla_{\boldsymbol{h}^{(t+1)}} L) + \boldsymbol{V}^{\top}(\nabla_{\boldsymbol{o}^{(t)}} L)$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad \boldsymbol{h}^{(t+1)} = \tanh(\boldsymbol{a}^{(t+1)}) \qquad \text{(element-wise)}$$

$$h_i^{(t+1)} = \tanh\left(a_i^{(t+1)}\right) \quad \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\boldsymbol{a}^{(t+1)} = W\boldsymbol{h}^{(t)} + U\boldsymbol{x}^{(t+1)} + \boldsymbol{b}$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \cdots$$
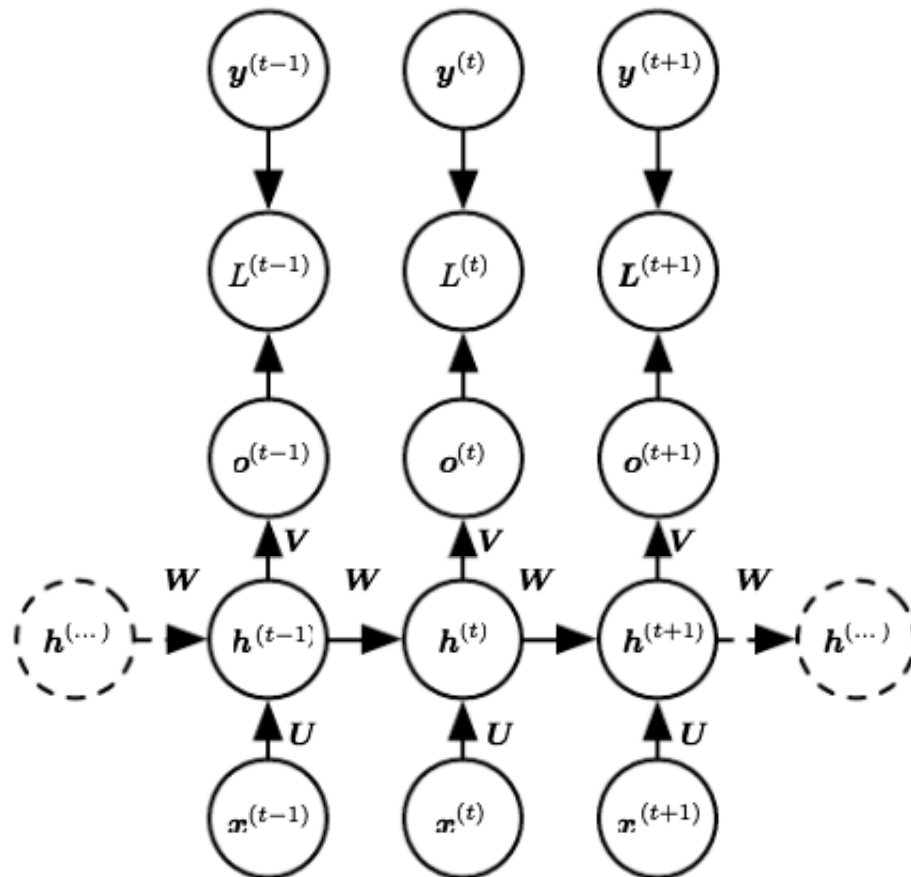
# RNN backprop (params)



**BPTT** (Back Propagation through Time)

- at any step $1 \leq t < \tau$

  since $\boldsymbol{o}^{(t)} = V\boldsymbol{h}^{(t)} + \boldsymbol{c}$

$$\nabla_{\boldsymbol{V}} L^{(t)} = (\nabla_{\boldsymbol{o}^{(t)}} L^{(t)})\boldsymbol{h}^{(t)\top} \qquad \nabla_{\boldsymbol{c}} L^{(t)} = \nabla_{\boldsymbol{o}^{(t)}} L^{(t)}$$
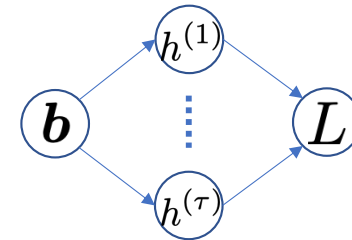
accumulate over time

$$\nabla_{\boldsymbol{V}} L = \sum_t \sum_i \left( \frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\boldsymbol{V}^{(t)}} o_i^{(t)} = \sum_t (\nabla_{\boldsymbol{o}^{(t)}} L) \boldsymbol{h}^{(t)\top}$$

$$\nabla_{\boldsymbol{c}} L = \sum_t \left( \frac{\partial \boldsymbol{o}^{(t)}}{\partial \boldsymbol{c}} \right)^{\top} \nabla_{\boldsymbol{o}^{(t)}} L = \sum_t \nabla_{\boldsymbol{o}^{(t)}} L$$

# RNN backprop (params)



**BPTT** (Back Propagation through Time)
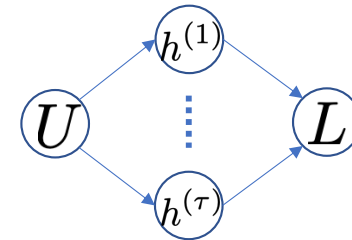
- at any step $1 \leq t < \tau$

since $\boldsymbol{h}^{(t)} = \tanh(\boldsymbol{a}^{(t)})$

$$\boldsymbol{a}^{(t)} = \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)} + \boldsymbol{b}$$

$$\nabla_{\boldsymbol{b}} L = \sum_t \left(\frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{b}^{(t)}}\right)^{\top} \nabla_{\boldsymbol{h}^{(t)}} L = \sum_t \text{diag}\left(1 - \left(\boldsymbol{h}^{(t)}\right)^2\right) \nabla_{\boldsymbol{h}^{(t)}} L$$
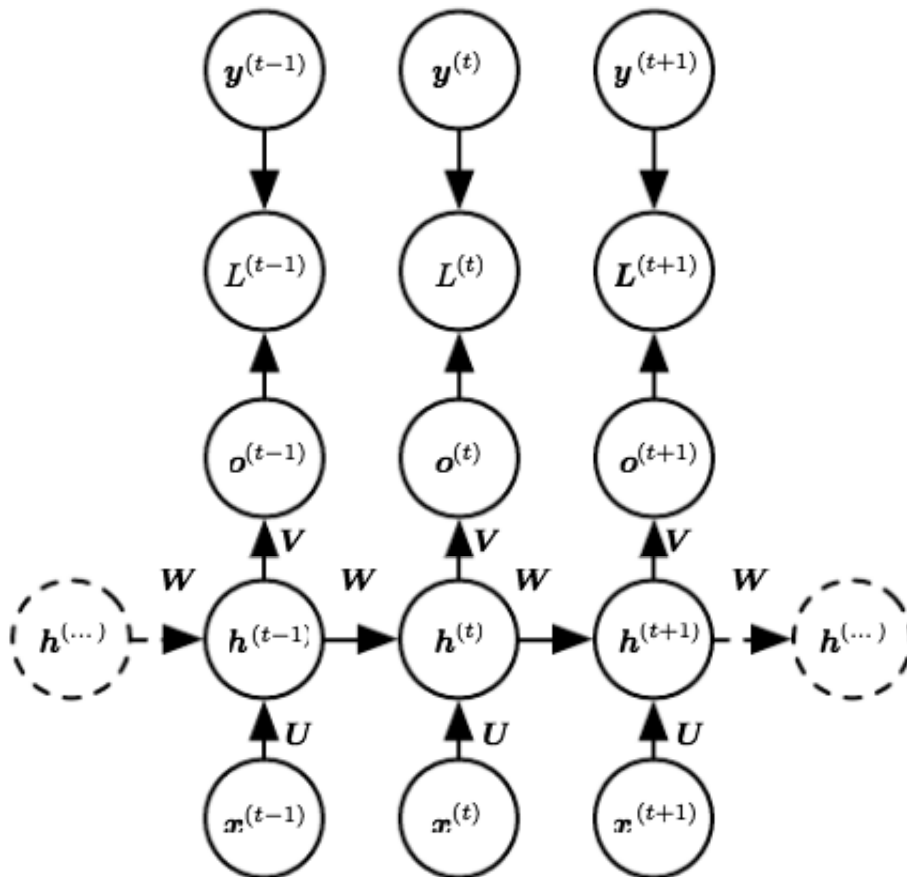
# RNN backprop (params)



**BPTT** (Back Propagation through Time)

- at any step $1 \leq t < \tau$

  since $\boldsymbol{h}^{(t)} = \tanh(\boldsymbol{a}^{(t)})$

  $$\boldsymbol{a}^{(t)} = \boldsymbol{W}\boldsymbol{h}^{(t-1)} + \boldsymbol{U}\boldsymbol{x}^{(t)} + \boldsymbol{b}$$

$$\nabla_{\boldsymbol{U}} L = \sum_t \text{diag}\left(1 - \left(\boldsymbol{h}^{(t)}\right)^2\right)(\nabla_{\boldsymbol{h}^{(t)}} L)\,\boldsymbol{x}^{(t)\top}$$

(We have done $\nabla_{\boldsymbol{h}^{(t-1)}} L$ , and leave $\nabla_{\boldsymbol{W}} L$ as an exercise.)