# Sequence Modeling: Recurrent and Recursive Nets

ADVANCED ARTIFICIAL INTELLIGENCE

JUCHEOL MOON

1

# How are you?

| DETECT LANGUAGE | FRENCH | ENGLISH | SPANISH | ⌄ | ⇄ | HINDI | ENGLISH | FRENCH | ⌄ |
|---|---|---|---|---|---|---|---|---|---|
| how are you? | | | | ✕ | | क्या हाल है? ⊘ | | | ☆ |

| DETECT LANGUAGE | ENGLISH | SPANISH | FRENCH | ⌄ | ⇄ | FRENCH | KOREAN | ENGLISH | ⌄ |
|---|---|---|---|---|---|---|---|---|---|
| how are you? | | | | ✕ | | 어떻게 지내? | | | ☆ |

| DETECT LANGUAGE | ENGLISH | SPANISH | FRENCH | ⌄ | ⇄ | CHINESE (SIMPLIFIED) | KOREAN | FRENCH | ⌄ |
|---|---|---|---|---|---|---|---|---|---|
| how are you? | | | | ✕ | | 你好吗? ⊘ | | | ☆ |

| DETECT LANGUAGE | ENGLISH | SPANISH | FRENCH | ⌄ | ⇄ | CHINESE (SIMPLIFIED) | KOREAN | FRENCH | ⌄ |
|---|---|---|---|---|---|---|---|---|---|
| how are you? | | | | ✕ | | Comment allez-vous? | | | ☆ |

🎤 🔊     12/5000 ⌨ ⌄     🔊     ⎘ ✎ ⤴
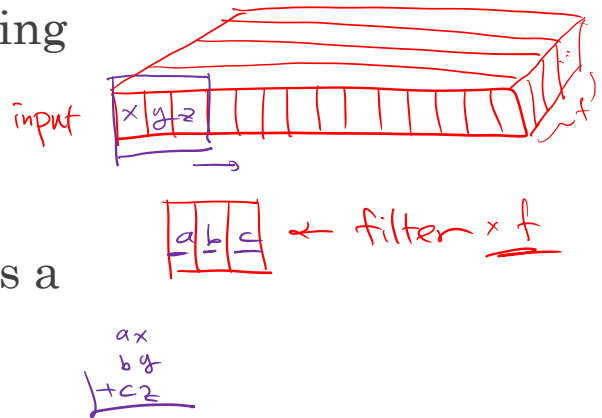
# Recurrent neural networks (RNN)

- Neural networks for processing sequential data

- Sharing parameters across different parts of a model

- The use of convolution across a 1-D temporal sequence
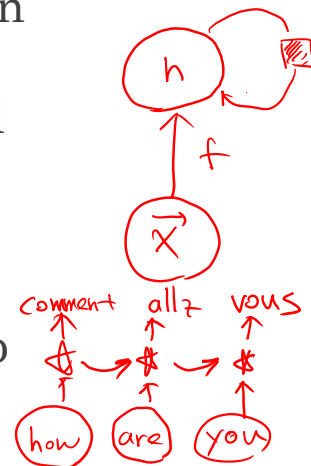  - The convolution operation allows a network to share parameters across time

*Handwritten annotations:* input; x y z; a b c ← filter × f; ax / by / +cz

3

# Recurrent neural networks (RNN)

- Recurrent networks share parameters in a different way.
  - Each member of the output is a function of the previous members of the output.
  - Each member of the output is produced using the same update rule applied to the previous outputs.

- RNNs as operating on a sequence that contains vectors $\vec{x}^{(t)}$ with the time step index $t$ ranging from 1 to $\tau$.
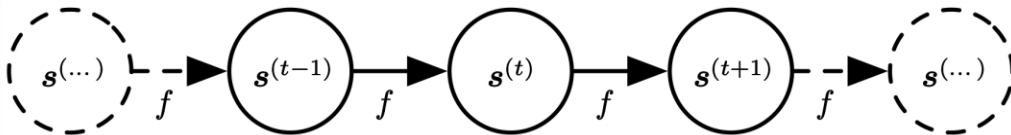
*Handwritten annotations:* h; f; $\vec{x}$; comment allz vous; how are you

4

# Unfolding Computational Graphs

- The classical form of a dynamical system

  - $$\vec{s}^{(t)} = f(\vec{s}^{(t-1)}; \vec{\theta})$$

- where $s^t$ is called the state of the system

- For a finite number of time steps $\tau$, the graph can be unfolded by applying the definition $\tau - 1$ times.

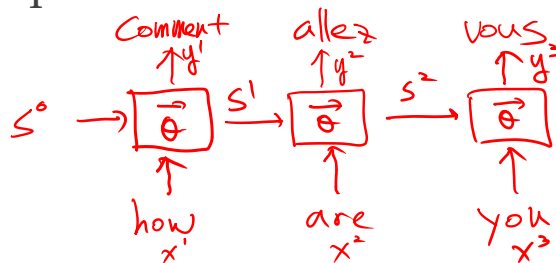  - $$\vec{s}^{(3)} = f(\vec{s}^{x(2)}; \vec{\theta}) = f\left(f(\vec{s}^{x(1)}; \vec{\theta}); \vec{\theta}\right)$$

# Unfolding Computational Graphs

- Many recurrent neural networks use equation

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

- where the state is the hidden units

- The network typically learns to use $h(t)$ as a kind of lossy summary of the task-relevant aspects of the past sequence of inputs up to $t$.
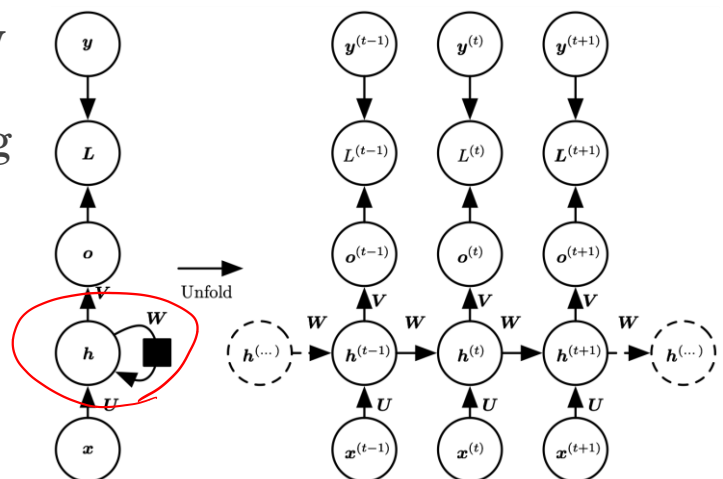
# Unfolding Computational Graphs

- The unfolding process thus introduces two major advantages:
  - Regardless of the sequence length, the learned model always has the same input size, because it is specified in terms of transition from one state to another state, rather than specified in terms of a variable-length history of states
  - It is possible to use the same transition function $f$ with the same parameters at every time step
- These two factors make it possible to learn a single model $f$ that operates on all time steps and all sequence lengths

# Recurrent Neural Networks

- Recurrent Hidden Units
  - A loss $L$ measures how far each $\vec{o}$ is from the corresponding training target $\vec{y}$.
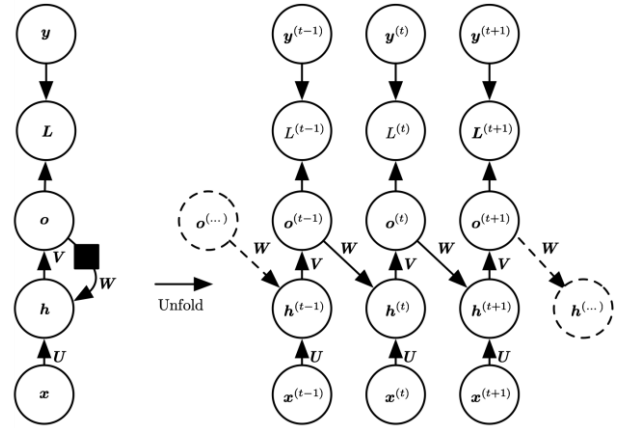  - The figure does not specify the choice of activation function for the hidden units.

# Recurrent Neural Networks

- Recurrence through only the Output
  - An RNN whose recurrence is the feedback connection from the output to the hidden layer.
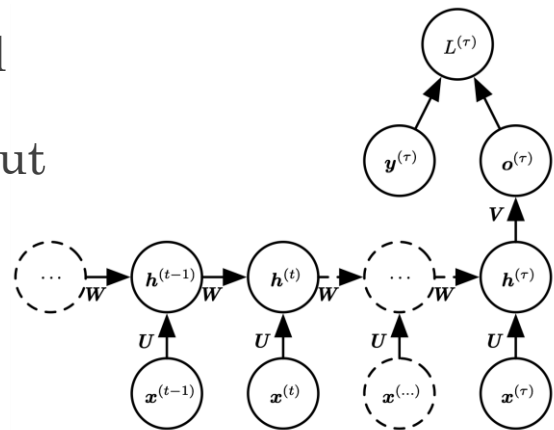
# Recurrent Neural Networks

- Sequence Input, Single Output
  - network can be used to summarize a sequence and produce a fixed-size representation used as input for further processing
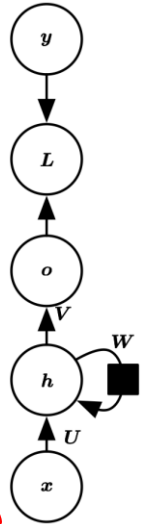
# Forward propagation

- Assumptions
  - The hyperbolic tangent activation function
  - The output $\vec{o}$ as giving the unnormalized log probabilities of each possible value of the discrete variable
  - Applying the softmax operation as a post-processing step to obtain a vector $\hat{\vec{y}}$ of normalized probabilities over the output

$$\underbrace{how}_{\vec{x}^{(1)}} \quad \underbrace{are}_{\vec{x}^{(2)}} \quad \underbrace{you}_{\vec{x}^{(3)}} \quad \xrightarrow{\text{one-hot-encode}}$$

$$\vec{x}^{(1)} = [\,0\ 0\ 0\ \cdots\ 1\ \cdots\ 0\ 0\,]$$
$$\vec{x}^{(2)} = [\,0\ 0\ \cdots\ 1\ \cdots\cdots\cdots\ 0\ 0\,]$$
$$\vec{x}^{(3)} = [\,0\ 0\ \cdots\cdots\ 1\ \cdots\ 0\ 0\,]$$

# Forward propagation

- Forward propagation begins with a specification of the initial state $\vec{h}^{(0)}$

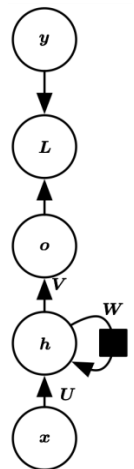$$\vec{a}^{(t)} = \vec{b} + \vec{W}\,\vec{h}^{(t-1)} + \vec{U}\,\vec{x}^{(t)}$$

$$\vec{h}^{(t)} = \tanh\left(\vec{a}^{(t)}\right)$$

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

$$\hat{\vec{y}}^{(t)} = \text{softmax}\left(\vec{o}^{(t)}\right)$$

$$\begin{cases} \vec{b}\ \&\ \vec{c} \quad \text{bias vectors} \\ \vec{W},\ \vec{U},\ \vec{V} \quad \text{weigh matrices} \end{cases}$$

# Forward propagation

- An example of a recurrent network that maps an input sequence to an output sequence of the same length

- The total loss for a given sequence of $\vec{x}$ values paired with a sequence of $\vec{y}$ values would then be just the sum of the losses over all the time steps

$$\vec{y} = [\overset{\in^0}{\hat{y}_1}, \overset{\in^0}{\underline{\hat{y}_2}} \cdots , \hat{y}_n]$$

$$L = \sum_t L^{(t)} = \sum_t L(\vec{y}^{(t)}, \overset{\wedge}{\vec{y}}^{(t)})$$

$$= \sum_t \sum_k y_k^{(t)} \log \hat{y}_k^{(t)} = -\sum_t \sum_k \underset{y_i^{(t)}=1}{y_k^{(t)}} \log \frac{e^{O_k^{(t)}}}{\sum_j e^{O_j^{(t)}}}$$

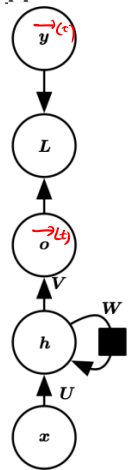$$= \sum_t \left(-O_i^{(t)} + \log \sum_j e^{O_j^{(t)}}\right) \qquad (i \text{ is the target})$$

13

13

# Computing the Gradient

- Computing the gradient of this loss function with respect to the parameters is an expensive operation
  - The back-propagation algorithm applied to the unrolled graph with $O(\tau)$ cost is called back-propagation through time (BPTT)

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

$$= L^{(1)} + L^{(2)} + \cdots + L^{(t)} + \cdots + L^{(\tau)} \qquad \frac{\partial L}{\partial L^{(t)}} = 1$$

$$\left(\nabla_{\vec{O}^{(t)}} L\right)_i = \frac{\partial L}{\partial O_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \cdot \frac{\partial L^{(t)}}{\partial O_i^{(t)}} = \frac{\partial L^{(t)}}{\partial O_i^{(t)}}$$
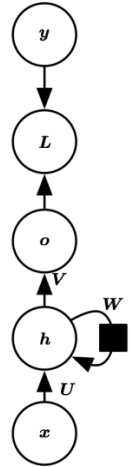


14

14

7

# Computing the Gradient

$$\frac{d \log f(x)}{dx} = \frac{f'(x)}{f(x)}$$

• $L^{(t)} = -o_i^{(t)} + \log \sum_j e^{o_j^{(t)}}$

$$\frac{\partial L^{(t)}}{\partial o_i^{(t)}} = -1 + \frac{\frac{\partial}{\partial o_i^{(t)}}\left(e^{o_1^{(t)}} + e^{o_2^{(t)}} + \cdots + e^{o_i^{(t)}} + \cdots\right)}{\sum_j e^{o_j^{(t)}}}$$

$$= -1 + \frac{e^{o_i^{(t)}}}{\underbrace{\sum_j e^{o_j^{(t)}}}_{\hat{y}_i^{(t)}}}$$

$$= -1 + \hat{y}_i^{(t)}$$

15

# Computing the Gradient

• At the final time step, $\vec{h}^{(t)}$ only has $\vec{o}^{(t)}$ as a descendent ($\tau$)

• $\left(\nabla_{\vec{h}^{(\tau)}} L\right)_i = \frac{\partial L}{\partial h_i^{(\tau)}} = \frac{\partial L}{\partial L^{(\tau)}} \sum_j \frac{\partial L^{(\tau)}}{\partial o_j^{(\tau)}} \frac{\partial o_j^{(\tau)}}{\partial h_i^{(\tau)}}$

$$\vec{o}^{(t)} = \vec{c} + V \vec{h}^{(t)} \qquad \frac{\partial L^{(\tau)}}{\partial o_j^{(\tau)}} = -1 + \hat{y}_j^{(\tau)}$$

$$\frac{\partial \vec{o}^{(\tau)}}{\partial h_i^{(\tau)}} = \frac{\partial\left(V_{j1} h_1^{(\tau)} + V_{j2} h_2^{(\tau)} + \cdots + V_{jn} h_n^{(\tau)}\right)}{\partial h_i^{(\tau)}} = V_{ji}$$

$$\nabla_{\vec{h}^{(\tau)}} L = V^{T} \nabla_{\vec{o}^{(\tau)}} L$$

16

8

# Computing the Gradient
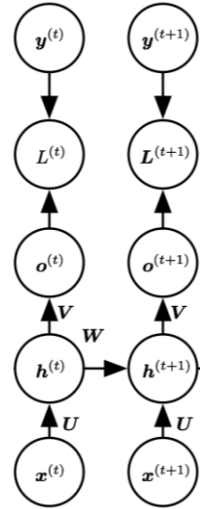
$\tanh' = 1 - \tanh^2$

- From $t = \tau - 1$ down to $t = 1$
  - $\left(\nabla_{\vec{h}^{(t)}} L\right)_i = \dfrac{\partial L}{\partial h_i^{(t)}} = \sum_j \dfrac{\partial h_j^{(t+1)}}{\partial h_i^{(t)}} \cdot \dfrac{\partial L}{\partial h_j^{(t+1)}} + \sum_j \dfrac{\partial o_j^{(t)}}{\partial h_i^{(t)}} \cdot \dfrac{\partial L}{\partial o_j^{(t)}}$

- $\vec{h}^{(t+1)} = \tanh\left(\vec{b} + \overrightarrow{W}\vec{h}^{(t)} + \vec{U}\vec{x}^{(t+1)}\right)$

- $h_j^{(t+1)} = \tanh\left(b_j + \sum_k W_{jk} h_k^{(t)} + \left(\vec{U}\vec{x}^{(t-1)}\right)_j\right)$

$\dfrac{\partial(W_{j1} + W_{j2} + \cdots + W_{ji} + \cdots)}{\partial h_i}$

- $\dfrac{\partial h_j^{(t+1)}}{\partial h_i^{(t)}} = \left(1 - \left(h_j^{(t+1)}\right)^2\right)\left(W_{ji}\right)$
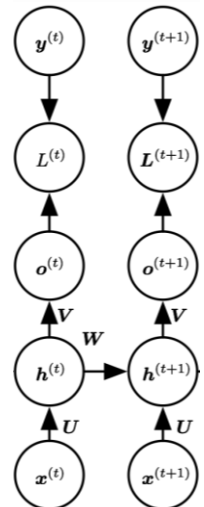
---

# Computing the Gradient

- From $t = \tau - 1$ down to $t = 1$
  - $\left(\nabla_{\vec{h}^{(t)}} L\right)_i = \sum_j \dfrac{\partial h_j^{(t+1)}}{\partial h_i^{(t)}} \dfrac{\partial L}{\partial h_j^{(t+1)}} + \sum_j \dfrac{\partial o_j^{(t)}}{\partial h_i^{(t)}} \cdot \dfrac{\partial L}{\partial o_j^{(t)}}$

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

$\left(\text{———}\right)\left(\big|\big|\right) = \left(\cdot\right)$

- $o_j^{(t)} = c_j + \sum_k V_{jk} \vec{h}_k^{(t)}$

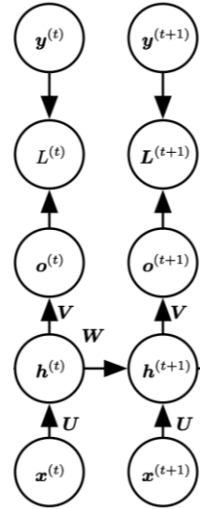- $\dfrac{\partial o_j^{(t)}}{\partial h_i^{(t)}} = V_{ji}$

# Computing the Gradient

- From $t = \tau - 1$ down to $t = 1$

- $$\frac{\partial \vec{h}^{(t+1)}}{\partial \vec{h}^{(t)}} = \begin{bmatrix} \dfrac{\partial h_1^{(t+1)}}{\partial h_1^{(t)}} & \cdots & \dfrac{\partial h_1^{(t+1)}}{\partial h_n^{(t)}} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial h_n^{(t+1)}}{\partial h_1^{(t)}} & \cdots & \dfrac{\partial h_n^{(t+1)}}{\partial h_1^{(t)}} \end{bmatrix}$$

- $$\nabla_{\vec{h}^{(t)}} L = \left( \frac{\partial \vec{h}^{(t+1)}}{\partial \vec{h}^{(t)}} \right)^T \left( \nabla_{\vec{h}^{(t+1)}} L \right) + \left( \frac{\partial \vec{o}^{(t)}}{\partial \vec{h}^{(t)}} \right)^T \left( \nabla_{\vec{o}^{(t)}} L \right)$$

$$= \vec{W}^T \left( \nabla_{\vec{h}^{(t+1)}} L \right) \operatorname{diag}\left( 1 - \left( \vec{h}^{(t+1)} \right)^2 \right) + \vec{V}^T \left( \nabla_{\vec{o}^{(t)}} L \right)$$

19

19

---

# Computing the Gradient

- $\vec{a}^{(t)} = \vec{b} + \overrightarrow{W}\vec{h}^{(t-1)} + \vec{U}\vec{x}^{(t)}$
- $\vec{h}^{(t)} = \tanh(\vec{a}^{(t)})$
- $\vec{o}^{(t)} = \vec{c} + \vec{V}\vec{h}^{(t)}$

- From $\nabla_{\vec{o}^{(t)}} L$ and $\nabla_{\vec{h}^{(t)}} L$

- $\nabla_{\vec{c}} L = \sum_t \nabla_{\vec{o}^{(t)}} L$

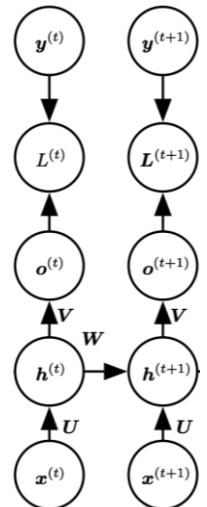- $\nabla_{\vec{b}} L = \sum_t \operatorname{diag}\left( 1 - \left( \vec{h}^{(t)} \right)^2 \right) \nabla_{\vec{h}^{(t)}} L$

- $\nabla_{\vec{V}} L = \sum_t \left( \nabla_{\vec{o}^{(t)}} L \right) \vec{h}^{(t)T}$

- $\nabla_{\overrightarrow{W}} L = \sum_t \operatorname{diag}\left( 1 - \left( \vec{h}^{(t)} \right)^2 \right) \left( \nabla_{\vec{h}^{(t)}} L \right) \vec{h}^{(t-1)T}$

- $\nabla_{\vec{U}} L = \sum_t \operatorname{diag}\left( 1 - \left( \vec{h}^{(t)} \right)^2 \right) \left( \nabla_{\vec{h}^{(t)}} L \right) \vec{x}^{(t)T}$
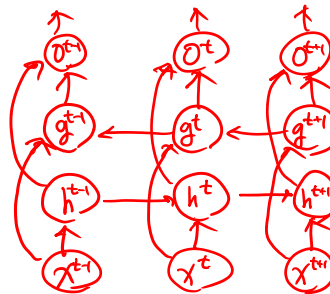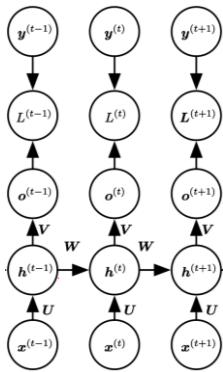
20

20

# Bidirectional RNNs

- All of the recurrent networks we have considered up to now have a "causal" structure.
  - The state at time $t$ only captures information from the past, $x^{(1)}, \ldots, x^{(t-1)}$, and the present input $x^{(t)}$.
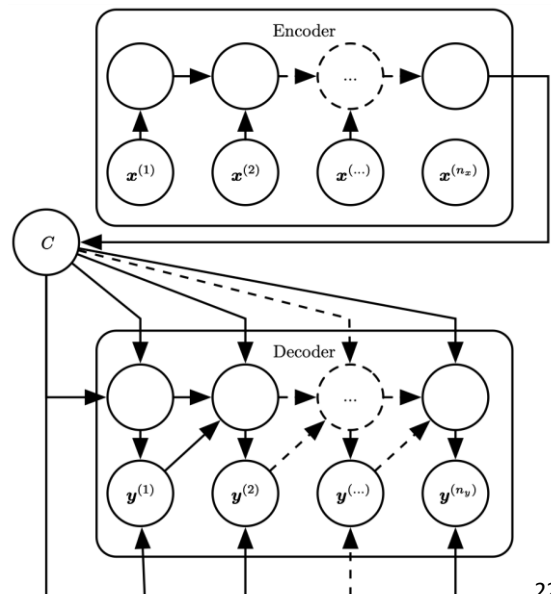
# Sequence-to-sequence

- How are you?
  - ↔ Comment allez vous?
  - ↔ ¿cómo estás?

- RNN can be trained to map an input sequence to an output sequence which is not necessarily of the same length

- Context variable $C$ represents a semantic summary of the input sequence and is given as input to the decoder.

# Encoder-Decoder

- The idea
  - (1) an encoder or reader or input RNN processes the input sequence. The encoder emits the context $C$, usually as a simple function of its final hidden state.
  - (2) a decoder or writer or output RNN is conditioned on that fixed-length vector to generate the output sequence

- One clear limitation of this architecture is when the context $C$ output by the encoder RNN has a dimension that is too small to properly summarize a long sequence
  - make $C$ a variable-length sequence

# Deep Recurrent Networks

- Three components of RNN
  - from the input to the hidden state
  - from the previous hidden state to the next hidden state
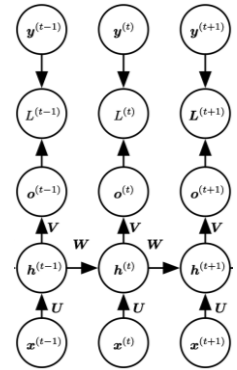  - from the hidden state to the output

# The Challenge of Long-Term Dependencies

- The basic problem is that gradients propagated over many stages tend to either vanish or explode.
  - Simple case w/o activation
- $\vec{h}^{(t)} = \vec{W}^{\top}\vec{h}^{(t-1)} = \left(\vec{W}^{t}\right)^{\top}\left(\vec{h}^{(0)}\right) = \vec{Q}^{\top}\vec{\Lambda}^{t}\vec{Q}\,\vec{h}^{(0)}$

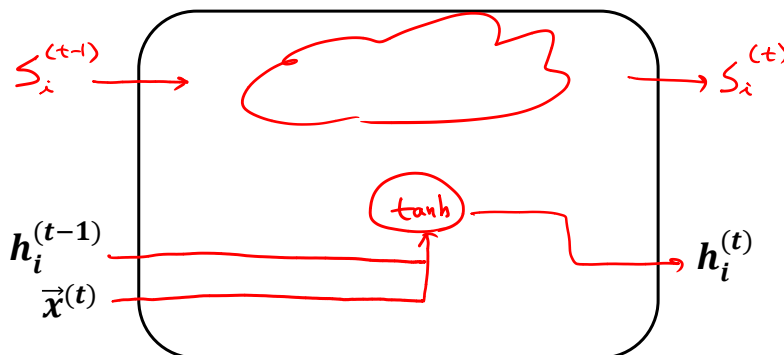- The eigenvalues are raised to the power of $t$ causing eigenvalues with magnitude less than one to decay to zero and eigenvalues with magnitude greater than one to explode

# Long Short-Term Memory (LSTM)

- Self-loops to produce paths where the gradient can flow for long durations
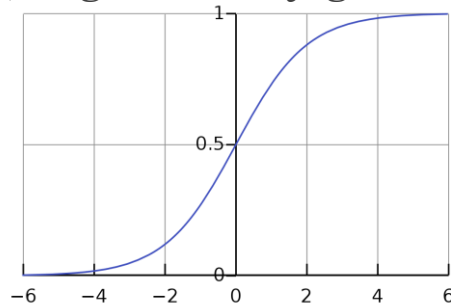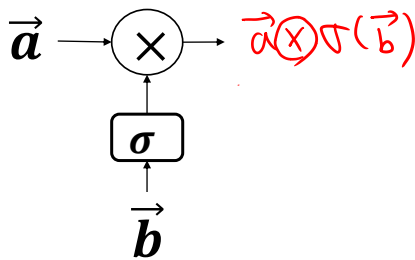- $\vec{h}^{(t)} = \tanh\left(\vec{b} + \vec{W}\,\vec{h}^{(t-1)} + \vec{U}\vec{x}^{(t)}\right)$

# Long Short-Term Memory (LSTM)

- The key to LSTM is the cell (internal) state $\vec{s}^{(t)}$

- The cell state runs straight down the entire chain, with only some minor linear interactions

- The LSTM does have the ability to remove or add information to the cell state, regulated by gates
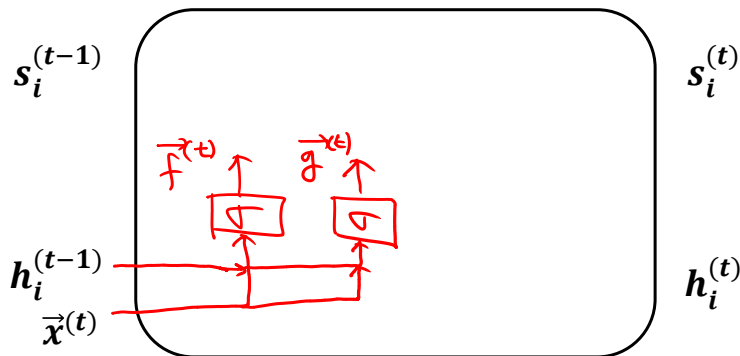


$\vec{a} \rightarrow \times \rightarrow \vec{a} \otimes \sigma(\vec{b})$

$\sigma$

$\vec{b}$

27

# Forget/input gate

- $f_i^{(t)} = \sigma\left( b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} \right)$

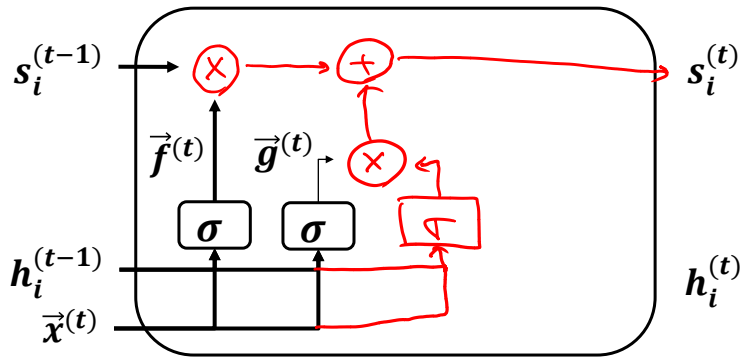- $g_i^{(t)} = \sigma\left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right)$



$s_i^{(t-1)}$ $\qquad$ $s_i^{(t)}$

$\vec{f}^{(t)}$ $\qquad$ $\vec{g}^{(t)}$

$\sigma$ $\qquad$ $\sigma$

$h_i^{(t-1)}$ $\qquad$ $h_i^{(t)}$

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

28

# Cell state

$$s_i^{(t)} = f_i^{(t)} S_i^{(t-1)} + g_i^{(t)} \sigma\left(b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)}\right)$$
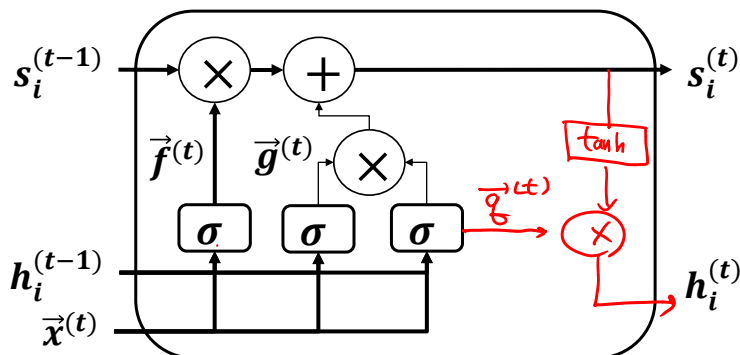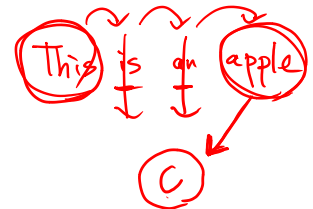
# Output gate

$$q_i^{(t)} = \sigma\left(b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)}\right)$$

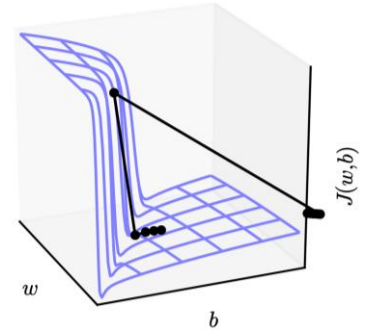$$h_i^{(t)} = \tanh\left(S_i^{(t)}\right) q_i^{(t)}$$

This is an apple

C

# Gradient Clipping

▪Recurrent net over many time steps tend to have derivatives that can be either very large in magnitude.

▪A simple type of solution has been in use by practitioners for many years

$$\text{if } \quad \|\vec{g}\| > \upsilon \leftarrow \text{hyper-param.}$$

$$\vec{g} \leftarrow \upsilon \cdot \frac{\vec{g}}{\|\vec{g}\|}$$
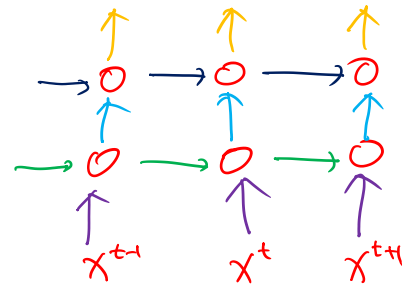
# Regularization

▪Naïve dropout
  ▪Use different masks at different time steps
  ▪No dropout on the recurrent layers

▪Variational dropout
  ▪Same mask at each time step and recurrent layers