

CS2109S Tutorial 8

Recurrent Neural Networks

(AY 25/26 Semester 2)

April 2, 2026

(Prepared by Benson)

Contents

Recurrent Neural Networks

- Q1. Model Selection for Sequential Data
- Q2. RNN Tracing
- Q3. Bidirectional RNNs
- Q4. 1D Convolution vs. RNNs

Why CNNs/RNNs?

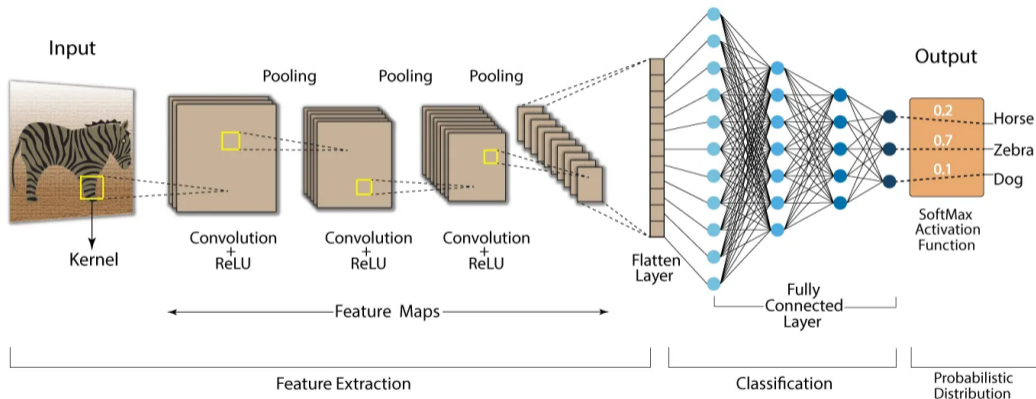
Should we simply whack a *very very very* deep Neural Network?

Problem context is important.

Idea: Exploit characteristics in the problem domain.

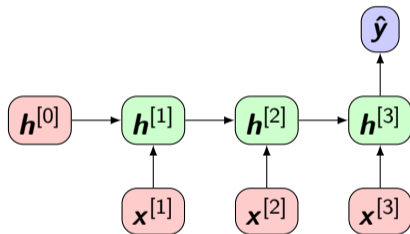
Why CNNs/RNNs?

Convolutional Neural Networks: Exploiting **locality**.



Why CNNs/RNNs?

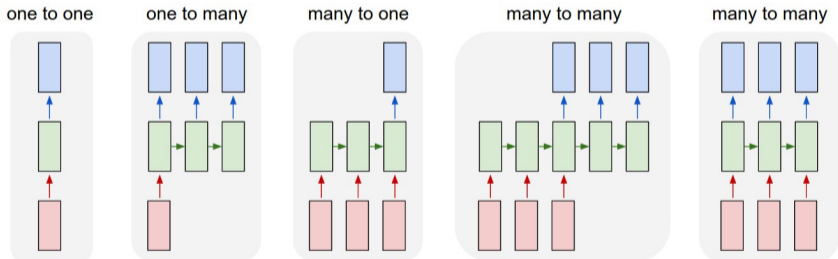
Recurrent Neural Networks: Exploiting **sequential nature**.



Q1. Model Selection for Sequential Data

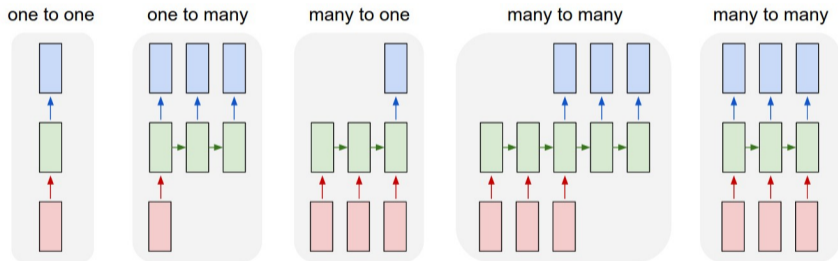
What type of model should we use?

- (a) Image Captioning One-to-many model
- (b) Stock market prediction Many-to-one model
- (c) Language Translation Many-to-many model



Q1. Model Selection for Sequential Data

Exercise: For music generation with a prompt (e.g. the style), what type of RNN model should we use?



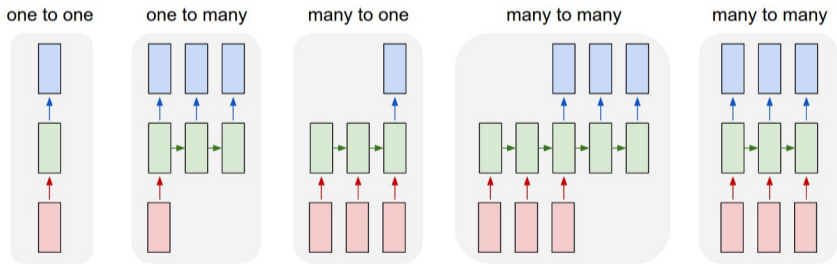
Q2. RNN Tracing

- (a) We would like to use simple RNN model to predict the last word in the following incomplete sentence made of 2 words:

“I love...”

What type of RNN model is needed?

- Many-to-one.

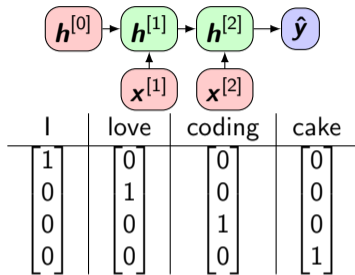


Q2. RNN Tracing

(b) Compute the hidden states $\mathbf{h}^{[1]}$, $\mathbf{h}^{[2]}$.

$$\begin{aligned}\mathbf{h}^{[1]} &= (\mathbf{W}^{[xh]})^\top \mathbf{x}^{[1]} + (\mathbf{W}^{[hh]})^\top \mathbf{h}^{[0]} \\ &= \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}\end{aligned}$$

$$\begin{aligned}\mathbf{h}^{[2]} &= (\mathbf{W}^{[xh]})^\top \mathbf{x}^{[2]} + (\mathbf{W}^{[hh]})^\top \mathbf{h}^{[1]} \\ &= \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}\end{aligned}$$



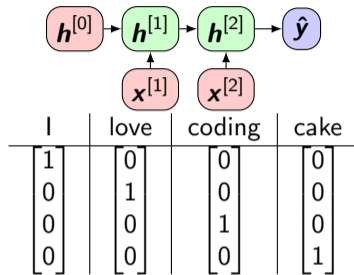
$$\mathbf{W}^{[xh]} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{W}^{[hh]} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$
$$\mathbf{W}^{[hy]} = \begin{bmatrix} 1 & 0 & 2 & -1 \\ 0 & 1 & 1 & 0 \end{bmatrix}, \quad \mathbf{h}^{[0]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Q2. RNN Tracing

(c) Compute the most probable next word.

$$\begin{aligned} \hat{y} &= \text{softmax} \left((W^{[hy]})^\top h^{[2]} \right) \\ &= \text{softmax} \left(\begin{pmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 2 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \right) \\ &= \text{softmax} \left(\begin{bmatrix} 1 \\ 2 \\ 4 \\ -1 \end{bmatrix} \right) = \begin{bmatrix} 0.04 \\ 0.11 \\ 0.84 \\ 0.01 \end{bmatrix} \end{aligned}$$

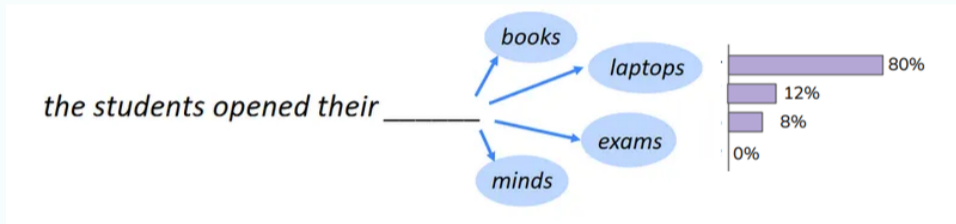
\therefore The most probable next word is "coding".



$$W^{[xh]} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad W^{[hh]} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$$

$$W^{[hy]} = \begin{bmatrix} 1 & 0 & 2 & -1 \\ 0 & 1 & 1 & 0 \end{bmatrix}, \quad h^{[0]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Extra: Greedy Search vs Sampling



Think: LLMs like ChatGPT allow you to specify a **temperature**. How does temperature affect the sampling?

Q2. RNN Tracing

- (d) What happens if we input “love I” instead?
 - ▶ The order of computing the hidden states will change.
 - ▶ The predictions will be different.

- (e) What is the potential issue of using one-hot encoding to represent the input vector?
 - ▶ High dimensional.
 - ▶ Two semantically similar words (e.g., “trend” and “pattern”) have zero similarity in one-hot space.





Extra: Token Embeddings

Demo: GPT4 Tokenizer

Write an email apologizing
to Sarah for the tragic
gardening mishap. Explain
how it happened.

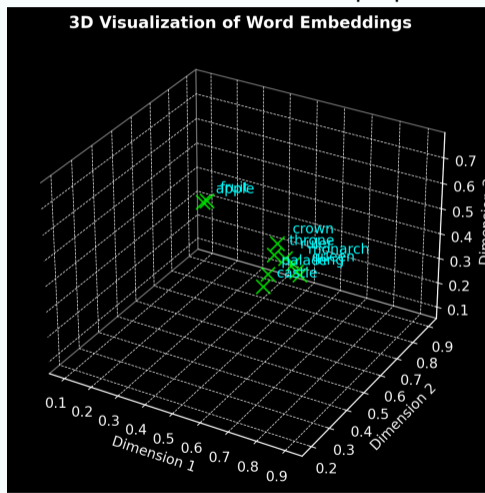
[10930, 448, 3719,
39950, 6396, 316, 32145,
395, 290, 62374, 66241,
80785, 403, 13, 115474,
1495, 480, 12570, 13]

Token Embeddings

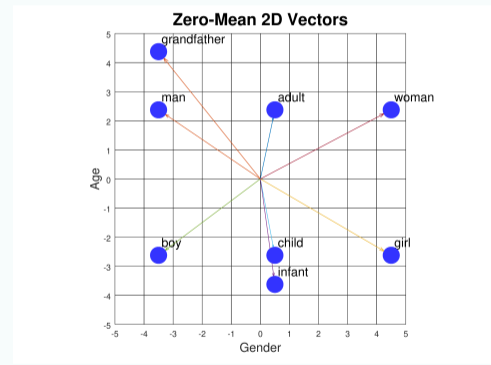
10930 
448 
3719 
39950 
⋮

Extra: Token Embeddings

Think: What are some ideal properties of token embeddings?



Capture semantic similarities



Meaningful linear substructures

Q3. Bidirectional RNNs

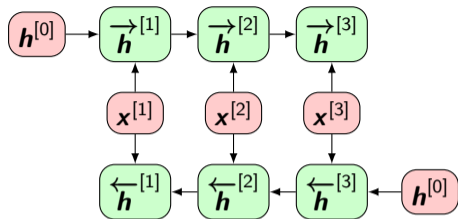
(a) Compute the hidden states

$$\vec{h}^{[1]}, \vec{h}^{[2]}, \vec{h}^{[3]} \text{ and } \overleftarrow{h}^{[1]}, \overleftarrow{h}^{[2]}, \overleftarrow{h}^{[3]}.$$

$$\begin{aligned} \vec{h}^{[1]} &= (\vec{W}^{[xh]})^\top \mathbf{x}^{[1]} + (\vec{W}^{[hh]})^\top \mathbf{h}^{[0]} \\ &= [1 \ 2 \ 3] [1 \ 0 \ 0]^\top + [1] [0] = [1] \end{aligned}$$

$$\begin{aligned} \vec{h}^{[2]} &= (\vec{W}^{[xh]})^\top \mathbf{x}^{[2]} + (\vec{W}^{[hh]})^\top \vec{h}^{[1]} \\ &= [1 \ 2 \ 3] [0 \ 1 \ 0]^\top + [1] [1] = [3] \end{aligned}$$

$$\begin{aligned} \vec{h}^{[3]} &= (\vec{W}^{[xh]})^\top \mathbf{x}^{[3]} + (\vec{W}^{[hh]})^\top \vec{h}^{[2]} \\ &= [1 \ 2 \ 3] [0 \ 0 \ 1]^\top + [1] [3] = [6] \end{aligned}$$



	not	very	good
$\vec{W}^{[xh]}$	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

$$\vec{W}^{[xh]} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \overleftarrow{W}^{[xh]} = \begin{bmatrix} 0 \\ 2 \\ 3 \end{bmatrix}$$

$$\vec{W}^{[hh]} = [1], \quad \overleftarrow{W}^{[hh]} = [1], \quad \mathbf{h}^{[0]} = [0]$$

Q3. Bidirectional RNNs

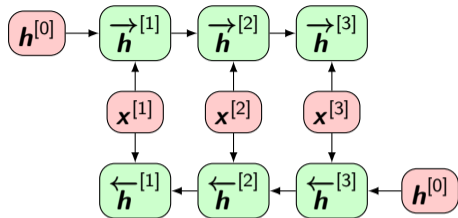
(a) Compute the hidden states

$$\vec{h}^{[1]}, \vec{h}^{[2]}, \vec{h}^{[3]} \text{ and } \overleftarrow{h}^{[1]}, \overleftarrow{h}^{[2]}, \overleftarrow{h}^{[3]}.$$

$$\begin{aligned} \overleftarrow{h}^{[3]} &= (\overleftarrow{W}^{[xh]})^\top \mathbf{x}^{[3]} + (\overleftarrow{W}^{[hh]})^\top \mathbf{h}^{[0]} \\ &= [0 \ 2 \ 3] [0 \ 0 \ 1]^\top + [1] [0] = [3] \end{aligned}$$

$$\begin{aligned} \overleftarrow{h}^{[2]} &= (\overleftarrow{W}^{[xh]})^\top \mathbf{x}^{[2]} + (\overleftarrow{W}^{[hh]})^\top \overleftarrow{h}^{[3]} \\ &= [0 \ 2 \ 3] [0 \ 1 \ 0]^\top + [1] [3] = [5] \end{aligned}$$

$$\begin{aligned} \overleftarrow{h}^{[1]} &= (\overleftarrow{W}^{[xh]})^\top \mathbf{x}^{[1]} + (\overleftarrow{W}^{[hh]})^\top \overleftarrow{h}^{[2]} \\ &= [0 \ 2 \ 3] [1 \ 0 \ 0]^\top + [1] [5] = [5] \end{aligned}$$



	not	very	good
$\vec{W}^{[xh]}$	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

$$\vec{W}^{[xh]} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \overleftarrow{W}^{[xh]} = \begin{bmatrix} 0 \\ 2 \\ 3 \end{bmatrix}$$

$$\vec{W}^{[hh]} = [1], \overleftarrow{W}^{[hh]} = [1], \mathbf{h}^{[0]} = [0]$$

Q3. Bidirectional RNNs

- (b) Compare the representation of the word “very” in the standard RNN and the bidirectional RNN.

▶ Standard RNN: $\mathbf{h}^{[2]} = [3]$.

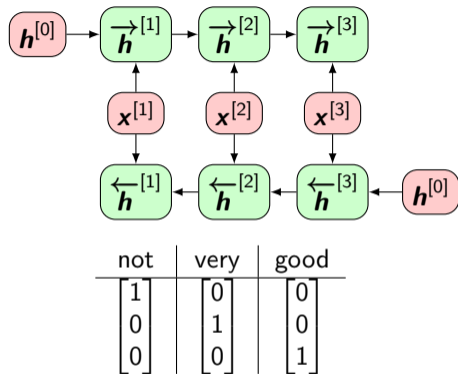
▶ Bidirectional RNN:

$$\mathbf{h}_{bi}^{[2]} = \begin{bmatrix} \vec{\mathbf{h}}^{[2]} \\ \leftarrow{\mathbf{h}}^{[2]} \end{bmatrix} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}.$$

Uses both **left** and **right** contexts.

- (c) Advantages and limitations of bidirectional RNN?

- ▶ Advantage: Right context would be helpful for some tasks, e.g. sentiment analysis, part-of-speech (POS) tagging.
- ▶ Disadvantage: Increases computational cost.



$$\vec{\mathbf{W}}^{[xh]} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \leftarrow{\mathbf{W}}^{[xh]} = \begin{bmatrix} 0 \\ 2 \\ 3 \end{bmatrix}$$

$$\vec{\mathbf{W}}^{[hh]} = [1], \quad \leftarrow{\mathbf{W}}^{[hh]} = [1], \quad \mathbf{h}^{[0]} = [0]$$

Q4. 1D Convolution vs. RNNs

Consider a sequence of sensor readings from a patient's heart beat data:

$$\mathbf{x} = [2 \quad 4 \quad 1 \quad -1 \quad 0].$$

- (a) Calculate the output after going through a single kernel $\mathbf{w} = [1 \quad -1]$.
 - ▶ $[-2 \quad 3 \quad 2 \quad -1]$.
- (b) What is the purpose of the filter?
 - ▶ Detect sudden drops in the sensor reading.
- (c) Compare this with a RNN, in terms of parallelization and long-term context.
 - ▶ RNNs have sequential dependencies while CNNs can be more effectively parallelized.
 - ▶ CNNs with $k = 2$ structurally fails to retrieve long-term contexts, while RNNs can theoretically succeed.

That's it!

