Recurrent Neural Networks

CISC 7026: Introduction to Deep Learning

University of Macau

- 1. Review
- 2. Sequence Modeling
- 3. Trace Theory
- 4. Elman Networks
- 5. Backpropagation through Time
- 6. LSTM
- 7. GRU
- 8. Linear Recurrent Models
- 9. Coding

- 1. Review
- 2. Sequence Modeling
- 3. Trace Theory
- 4. Elman Networks
- 5. Backpropagation through Time
- 6. LSTM
- 7. GRU
- 8. Linear Recurrent Models
- 9. Coding

- 1. Review
- 2. Sequence Modeling
- 3. Trace Theory
- 4. Elman Networks
- 5. Backpropagation through Time
- 6. LSTM
- 7. GRU
- 8. Linear Recurrent Models
- 9. Coding

- 1. Review
- 2. Sequence Modeling
- 3. Trace Theory
- 4. Elman Networks
- 5. Backpropagation through Time
- 6. LSTM
- 7. GRU
- 8. Linear Recurrent Models
- 9. Coding

- Conv is good at some Things
- Generalizes to n dimensions, etc
- Shortcomings locality not good in all cases
- Equivariance not good in all cases
- Imagine memories created over a lifetime
 - Targeted remembering and forgetting of important events
- RNNs developed specifically for temporal problems
- Are there other ways to process sequences

Convolution is not the only way to process sequential data

We previously used convolution to model signals

We previously used convolution to model signals

Some signals, such as audio, occur over time

We previously used convolution to model signals

Some signals, such as audio, occur over time

Convolution approaches temporal data from an electrical engineering approach

We previously used convolution to model signals

Some signals, such as audio, occur over time

Convolution approaches temporal data from an electrical engineering approach

Today, we will process temporal data using a psychological approach

Convolution makes use of locality and translation equivariance properties

Convolution makes use of locality and translation equivariance properties

This makes learning more efficient, but not all problems benefit from locality and equivariance

Not local! Two events occur separated by 20 years

Not local! Two events occur separated by 20 years

Example 2: Your parent changes your diaper

Not local! Two events occur separated by 20 years

Example 2: Your parent changes your diaper

• Not equivariant! Ok if you are a baby, different meaning if you are an adult!

Not local! Two events occur separated by 20 years

Example 2: Your parent changes your diaper

• Not equivariant! Ok if you are a baby, different meaning if you are an adult!

Example 3: You kiss person A, then 5 years later marry person B

Not local! Two events occur separated by 20 years

Example 2: Your parent changes your diaper

• Not equivariant! Ok if you are a baby, different meaning if you are an adult!

Example 3: You kiss person A, then 5 years later marry person B

• Not equivariant or local! If you marry person B then kiss person A, it is different...

Not local! Two events occur separated by 20 years

Example 2: Your parent changes your diaper

• Not equivariant! Ok if you are a baby, different meaning if you are an adult!

Example 3: You kiss person A, then 5 years later marry person B

• Not equivariant or local! If you marry person B then kiss person A, it is different...

Question: Any other examples?

If your problem has local and equivariant structure, use convolution

If your problem has local and equivariant structure, use convolution

For other problems, we need something else!

If your problem has local and equivariant structure, use convolution

For other problems, we need something else!

Humans experience time and process temporal data

If your problem has local and equivariant structure, use convolution

For other problems, we need something else!

Humans experience time and process temporal data

Can we use this to come up with a new neural network architecture?

- 1. Review
- 2. Sequence Modeling
- 3. Trace Theory
- 4. Elman Networks
- 5. Backpropagation through Time
- 6. LSTM
- 7. GRU
- 8. Linear Recurrent Models
- 9. Coding

- 1. Review
- 2. Sequence Modeling
- 3. Trace Theory
- 4. Elman Networks
- 5. Backpropagation through Time
- 6. LSTM
- 7. GRU
- 8. Linear Recurrent Models
- 9. Coding

How do humans experience time?

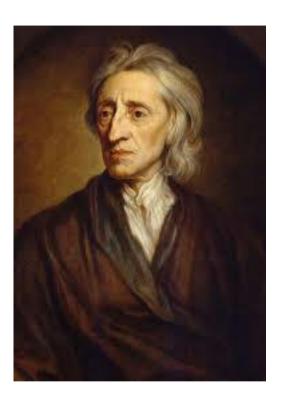
How do humans experience time?

Humans create memories

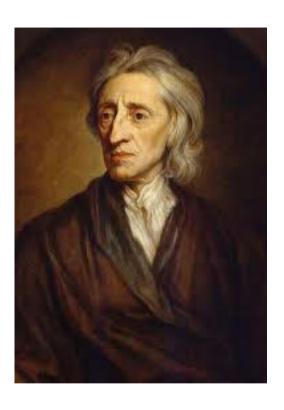
How do humans experience time?

Humans create memories

We experience time by reasoning over our memories

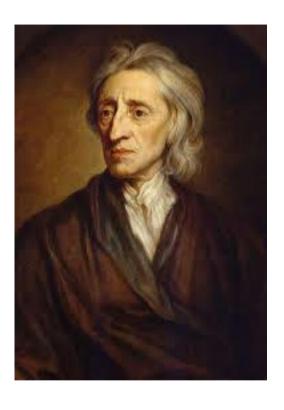


John Locke (1690) believed that conciousness and identity arise from memories



John Locke (1690) believed that conciousness and identity arise from memories

If all your memories were erased, you would be a different person



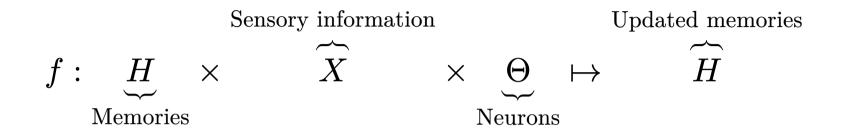
John Locke (1690) believed that conciousness and identity arise from memories

If all your memories were erased, you would be a different person

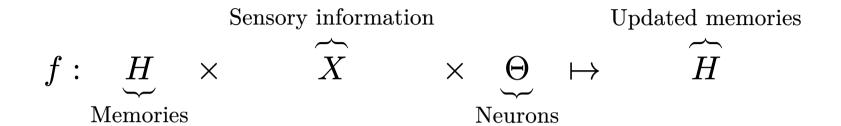
Without the ability to reason over memories, we would simply react to stimuli like bacteria

So how do we model memories in humans?

So how do we model memories in humans?



So how do we model memories in humans?



All your memories represented as a vector $h \in H$

So how do we model memories in humans?

Sensory information Updated memories
$$f: \underbrace{H} \times \overbrace{X} \times \underbrace{\Theta} \mapsto \overbrace{H}$$
 Memories Neurons

All your memories represented as a vector $m{h} \in H$

Everything you currently sense (sight, touch, sound, emotions) is a vector $\boldsymbol{x} \in X$

So how do we model memories in humans?

Sensory information Updated memories
$$f: \underbrace{H} \times \overbrace{X} \times \underbrace{\Theta} \mapsto \overbrace{H}$$

Memories Neurons

All your memories represented as a vector $m{h} \in H$

Everything you currently sense (sight, touch, sound, emotions) is a vector ${\boldsymbol x} \in X$

We update our memories following

$$\boldsymbol{s}_t = f(\boldsymbol{x}_t, \boldsymbol{h}_{t-1})$$

$$\boldsymbol{s}_t = f(\boldsymbol{x}_t, \boldsymbol{h}_{t-1})$$



After we have constructed our memories \boldsymbol{h}_t , we do not recall all information at once

After we have constructed our memories \boldsymbol{h}_t , we do not recall all information at once

Example: I ask you your favorite ice cream flavor

After we have constructed our memories \boldsymbol{h}_t , we do not recall all information at once

Example: I ask you your favorite ice cream flavor

You recall previous times you ate ice cream, but not your phone number

After we have constructed our memories \boldsymbol{h}_t , we do not recall all information at once

Example: I ask you your favorite ice cream flavor

You recall previous times you ate ice cream, but not your phone number

We should model this too

$$q: H \times X \times \Theta \mapsto Y$$

After we have constructed our memories \boldsymbol{h}_t , we do not recall all information at once

Example: I ask you your favorite ice cream flavor

You recall previous times you ate ice cream, but not your phone number

We should model this too

$$g: H \times X \times \Theta \mapsto Y$$

$$oldsymbol{y}_t = g(oldsymbol{h}_t, oldsymbol{x}_t, oldsymbol{ heta}_t)$$

$$f: H \times X \times \Theta \mapsto H; \quad g: H \times X \times \Theta \mapsto Y$$

$$f: H \times X \times \Theta \mapsto H; \quad g: H \times X \times \Theta \mapsto Y$$

$$oldsymbol{h_t} = f(oldsymbol{x}_t, oldsymbol{h}_{t-1}); \quad oldsymbol{y}_t = g(oldsymbol{x}_t, oldsymbol{h}_t)$$

$$f: H \times X \times \Theta \mapsto H; \quad g: H \times X \times \Theta \mapsto Y$$

$$oldsymbol{h_t} = f(oldsymbol{x}_t, oldsymbol{h}_{t-1}); \quad oldsymbol{y}_t = g(oldsymbol{x}_t, oldsymbol{h}_t)$$

$$h_{t+1} = f(x_{t+1}, h_t); \quad y_{t+1} = g(x_{t+1}, h_{t+1})$$

$$f: H \times X \times \Theta \mapsto H; \quad g: H \times X \times \Theta \mapsto Y$$

$$oldsymbol{h_t} = f(oldsymbol{x}_t, oldsymbol{h}_{t-1}); \quad oldsymbol{y}_t = g(oldsymbol{x}_t, oldsymbol{h}_t)$$

$$h_{t+1} = f(x_{t+1}, h_t); \quad y_{t+1} = g(x_{t+1}, h_{t+1})$$

$$m{h}_{t+2} = f(m{x}_{t+2}, m{h}_{t+1}); \quad m{y}_{t+2} = g(m{x}_{t+2}, m{h}_{t+2})$$

$$f: H \times X \times \Theta \mapsto H; \quad g: H \times X \times \Theta \mapsto Y$$

The function f is **recurrent** because it outputs a future input

$$egin{aligned} &m{h}_t = f(m{x}_t, m{h}_{t-1}); \quad m{y}_t = g(m{x}_t, m{h}_t) \ &m{h}_{t+1} = f(m{x}_{t+1}, m{h}_t); \quad m{y}_{t+1} = g(m{x}_{t+1}, m{h}_{t+1}) \ &m{h}_{t+2} = f(m{x}_{t+2}, m{h}_{t+1}); \quad m{y}_{t+2} = g(m{x}_{t+2}, m{h}_{t+2}) \ & \cdot \end{aligned}$$

Steven Morad

$$f: H \times X \times \Theta \mapsto H; \quad g: H \times X \times \Theta \mapsto Y$$

The function f is **recurrent** because it outputs a future input

$$egin{aligned} &m{h}_t = f(m{x}_t, m{h}_{t-1}); \quad m{y}_t = g(m{x}_t, m{h}_t) \ &m{h}_{t+1} = f(m{x}_{t+1}, m{h}_t); \quad m{y}_{t+1} = g(m{x}_{t+1}, m{h}_{t+1}) \ &m{h}_{t+2} = f(m{x}_{t+2}, m{h}_{t+1}); \quad m{y}_{t+2} = g(m{x}_{t+2}, m{h}_{t+2}) \ &\vdots \end{aligned}$$

If f, g are neural networks, then we call it a **recurrent neural network** (RNN)

$$f(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_1^{\top} \boldsymbol{h}_{t-1} + \boldsymbol{\theta}_2^{\top} \overline{\boldsymbol{x}}_t)$$

$$\begin{split} f(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) &= \sigma(\boldsymbol{\theta}_1^\top \boldsymbol{h}_{t-1} + \boldsymbol{\theta}_2^\top \overline{\boldsymbol{x}}_t) \\ g(\boldsymbol{h}_t, \boldsymbol{x}_t, \boldsymbol{\theta}) &= \sigma\big(\boldsymbol{\theta}_3^\top \overline{\boldsymbol{h}}_t\big) \end{split}$$

$$\begin{split} f(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) &= \sigma(\boldsymbol{\theta}_1^\top \boldsymbol{h}_{t-1} + \boldsymbol{\theta}_2^\top \overline{\boldsymbol{x}}_t) \\ g(\boldsymbol{h}_t, \boldsymbol{x}_t, \boldsymbol{\theta}) &= \sigma\big(\boldsymbol{\theta}_3^\top \overline{\boldsymbol{h}}_t\big) \end{split}$$

TODO BPTT TODO Forgetting

Add forgetting

$$f_{ ext{forget}}(m{h}_{t-1}, m{x}_t, m{ heta}) = \sigma(m{ heta}_1^ op m{\overline{x}}_t + m{ heta}_2^ op m{h}_{t-1})$$

$$f(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) = \sigma \big(\boldsymbol{\theta}_3^{\top} \overline{\boldsymbol{x}}_t + \boldsymbol{\theta}_4^{\top} \boldsymbol{h}_{t-1} \odot f_{\text{forget}}(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) \big)$$

Minimal gated unit (MGU)

Minimal gated unit (MGU)

$$\begin{split} f_{\text{forget}}(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) &= \sigma(\boldsymbol{\theta}_1^\top \overline{\boldsymbol{x}}_t + \boldsymbol{\theta}_2^\top \boldsymbol{h}_{t-1}) \\ f_2(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) &= \sigma_{\text{tanh}} \big(\boldsymbol{\theta}_3^\top \overline{\boldsymbol{x}}_t + \boldsymbol{\theta}_4^\top \big(f_{\text{forget}}(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) \odot \boldsymbol{h}_{t-1} \big) \big) \\ f(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) &= \\ \big(\big(1 - f_{\text{forget}}(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) \big) \odot \boldsymbol{h}_{t-1} + f_{\text{forget}}(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) \odot f_2(\boldsymbol{h}_{t-1}, \boldsymbol{x}_t, \boldsymbol{\theta}) \big) \end{split}$$

My PhD research focused on making RNNs more like human memory

https://openreview.net/forum?id=KTfAtro6vP