# Sequence models

**Lecture 11 — CS 577 Deep Learning**

Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology

October 30, 2024

- Seq model
- Flash Attention

Systems (ML + optimizing for memory bandwidth)

memory-bound
moving things
from
memory
to compute
(GPU)
is the
bottleneck

- Andre Bauer /
  time series
    for AI
  Minxuan Zhou
  Hardware topics
    Tailored hardware
  In-memory compute

- Matt Smith
  CV for ecology   understanding insects   biodiversity
    review for final

# Sequence models

**Lecture 11 — CS 577 Deep Learning**

Instructor: Yutong Wang

Computer Science
Illinois Institute of Technology
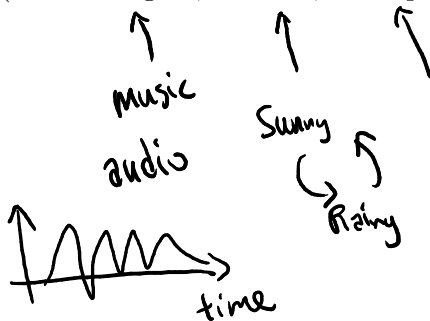
October 30, 2024

# Topics

- Temporal convolutional networks
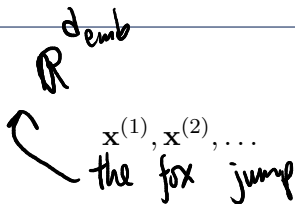- Recurrent neural networks
- Transformers

# Sequences

dependency across time

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$$

- Example: Time series (acoustic signal, weather, stock price,...)

music
audio

Sunny
Rainy

time

# Sequences

$$\mathbb{R}^{d_{emb}}$$

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$$

the fox jump

- Example: Time series (acoustic signal, weather, stock price,...)
- Example: Sequence of word/token embeddings

# Sequences

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots$$

- Example: Time series (acoustic signal, weather, stock price,...)
- Example: Sequence of word/token embeddings
- Task: denoising (make a signal less noisy)    clean up audio

## Sequences

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots$$

- Example: Time series (acoustic signal, weather, stock price,...)
- Example: Sequence of word/token embeddings

- Task: denoising (make a signal less noisy)
- Task: forecasting (what is the weather like next week?)

## Sequences

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$$

- Example: Time series (acoustic signal, weather, stock price,...)
- Example: Sequence of word/token embeddings

- Task: denoising (make a signal less noisy)
- Task: forecasting (what is the weather like next week?)
- Task: classification (is this news article about sport or about technology?)

*document classif.*

*Is generative next word prediction*

# Sequences

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots$$

- Example: Time series (acoustic signal, weather, stock price,...)
- Example: Sequence of word/token embeddings

- Task: denoising (make a signal less noisy)
- Task: forecasting (what is the weather like next week?)
- Task: classification (is this news article about sport or about technology?)
- Task: regression (how favorable is this product review?)

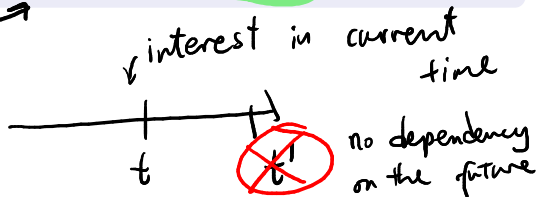Sentiment analysis        $0 \longrightarrow 100$

# Causal models

put "causal rail guard" in training loop implicitly make the model causal

$$f(\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(t)}; \theta) = \mathbf{y}^{(t)}$$

### Causality assumption

The output $\mathbf{y}^{(t)}$ does not depend on future values of $\mathbf{x}^{(t')}$ where $t' > t$

interest in current time

no dependency on the future

t

# Prediction/forecasting/generative models

input space = output space

↑ when is the not true?

- Let $\hat{\mathbf{x}}^{(t+1)} := \mathbf{y}^{(t)}$

-
$$\underbrace{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots, \mathbf{x}^{(C)}}_{\text{observed}} \overset{f}{\to} \underbrace{\hat{\mathbf{x}}^{(C+1)}, \hat{\mathbf{x}}^{(C+2)}, \ldots}_{\text{predicted/generated}}$$

- Next: for simplicity, let's consider $x^{(t)} \in \mathbb{R}$. (So drop the bold fontface.)

Context = ↑ time observed max up to

# Prediction/forecasting/generative models

- Let $\hat{\mathbf{x}}^{(t+1)} := \mathbf{y}^{(t)}$

- $$\underbrace{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots, \mathbf{x}^{(C)}}_{\text{observed}} \quad \xrightarrow{f} \quad \underbrace{\hat{\mathbf{x}}^{(C+1)}, \hat{\mathbf{x}}^{(C+2)}, \ldots}_{\text{predicted/generated}}$$

- Next: for simplicity, let's consider $x^{(t)} \in \mathbb{R}$. (So drop the bold fontface.)

Scalar

Later: how this leads to simplest non-trivial RNNs

# Autoregressive model

$$AR(5)$$

- For $t \geq s$

$s := $ filter-size

$\hat{x}^{(t+1)} = y^{(t)} = f(x^{(1)}, \ldots, x^{(t)}; w_0, \ldots, w_{s-1}, b)$

$s = 5$

Goal

$$= \sum_{\tau=0}^{s-1} w_\tau x^{(t-\tau)} + b$$

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

$x^{(t)}$

dot product

weight / slope

## Autoregressive model

- For $t \geq s$

$$\hat{x}^{(t+1)} = y^{(t)} = f(x^{(1)}, \ldots, x^{(t)}; w_0, \ldots, w_{s-1}, b)$$

$$= \sum_{\tau=0}^{s-1} w_\tau x^{(t-\tau)} + b$$

- Loss

$$\frac{1}{C} \sum_{t=1}^{C} (x^{(t+1)} - f(x^{(1)}, \ldots, x^{(t)}; w_0, \ldots, w_{s-1}, b))^2$$

↙ MSE

predict / forecast

next time data arrives

# Autoregressive model

- For $t \geq s$

$$y^{(t)} = f(x^{(1)}, \ldots, x^{(t)}; w_0, \ldots, w_{s-1}, b)$$
$$= \sum_{\tau=0}^{s-1} w_\tau x^{(t-\tau)} + b$$

- Loss

$$\frac{1}{C} \sum_{t=1}^{C} (x^{(t+1)} - f(x^{(1)}, \ldots, x^{(t)}; w_0, \ldots, w_{s-1}, b))^2$$

- Solve for $w_0, \ldots, w_{s-1}, b$

  $\hookrightarrow$ np. array

## Autoregressive model

- For $t \geq s$

$$y^{(t)} = f(x^{(1)}, \ldots, x^{(t)}; w_0, \ldots, w_{s-1}, b)$$
$$= \sum_{\tau=0}^{s-1} w_\tau x^{(t-\tau)} + b$$

- Loss

$$\frac{1}{C} \sum_{t=1}^{C} (x^{(t+1)} - f(x^{(1)}, \ldots, x^{(t)}; w_0, \ldots, w_{s-1}, b))^2$$

- Solve for $w_0, \ldots, w_{s-1}, b$
- Next: seq2col

im2col

# Seq2col: in class exercise 1

In class exercise

- Complete the `seq2col` function
- Discuss with your neighbors regarding questions in the "Fibonacci" block

# Seq2col: in class exercise 1 discussion

$$1 \quad 1 \quad 2 \quad 3 \quad 5 \quad 8$$

$$x^{(t+1)} = x^{(t)} + x^{(t-1)}$$

How do you interpret the coefficients when we set `filter_size = 3`?

```
1 filter_size = 3
2 def fibonacci(n):
3 #    [...]
4
5 fibonnaci_seq = list(fibonacci(20))
6
7 X, y = seq2col(fibonnaci_seq[:10], filter_size)
8
9 X_tilde = np.hstack([X, np.ones((X.shape[0], 1))])
10 w = np.linalg.pinv(X_tilde) @ y
11 np.round(w,5)
12 # array([-0.,  1.,  1.,  0.])
```

most recent

least

w

b

# Another way to look at autoregressive models

- Drop the bias term for simplicity

→ to RNN (recurrent neural network)

$$y^{(t)} = w_0 x^{(t)} + \underbrace{\sum_{\tau=1}^{s-1} w_\tau x^{(t-\tau)}}_{\text{weight}} = w_0 x^{(t)} + \underbrace{\begin{bmatrix} w_1 \\ \vdots \\ w_{s-1} \end{bmatrix}^\top \begin{bmatrix} x^{(t-1)} \\ \vdots \\ x^{(t-(s-1))} \end{bmatrix}}_{\text{dot product}}$$

# Another way to look at autoregressive models

- Drop the bias term for simplicity

$$y^{(t)} = w_0 x^{(t)} + \sum_{\tau=1}^{s-1} w_\tau x^{(t-\tau)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ \vdots \\ w_{s-1} \end{bmatrix}^\top \begin{bmatrix} x^{(t-1)} \\ \vdots \\ x^{(t-(s-1))} \end{bmatrix}$$

- History up to time $t-1$ (aka hidden state at time $t-1$)

$$y^{(t)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ \vdots \\ w_{s-1} \end{bmatrix}^\top \tilde{\mathbf{h}}^{(t-1)}$$

def !!

current input + weight @ hidden_state

# Another way to look at autoregressive models

- Drop the bias term for simplicity

$$y^{(t)} = w_0 x^{(t)} + \sum_{\tau=1}^{s-1} w_\tau x^{(t-\tau)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ \vdots \\ w_{s-1} \end{bmatrix}^\top \begin{bmatrix} x^{(t-1)} \\ \vdots \\ x^{(t-(s-1))} \end{bmatrix}$$

- History up to time $t-1$ (aka hidden state at time $t-1$)

*Instead of modeling $y^{(t)}$*

$$y^{(t)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ \vdots \\ w_{s-1} \end{bmatrix}^\top \tilde{\mathbf{h}}^{(t-1)} \quad \leftarrow \quad \textit{hidden state}$$

- Next: directly model the hidden state

# Recurrent unit (with linear activation)

- Recurrent unit

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{w_{in}}x^{(t)}$$

hidden
state

input

# Recurrent unit (with linear activation)

- Recurrent unit

$d_h$

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{w_{in}}x^{(t)}$$

- $\mathbf{h}^{(t)} \in \mathbb{R}^{d_h}$ (hidden) states at time $t$

# Recurrent unit (with linear activation)

- Recurrent unit

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{w_{in}}x^{(t)}$$

- $\mathbf{h}^{(t)} \in \mathbb{R}^{d_{\mathtt{h}}}$ (hidden) states at time $t$
- $d_{\mathtt{h}}$ hidden dimension

# Recurrent unit (with linear activation)

- Recurrent unit

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W}_{\texttt{rec}}\mathbf{h}^{(t-1)} + \mathbf{w}_{\texttt{in}}x^{(t)}$$

- $\mathbf{h}^{(t)} \in \mathbb{R}^{d_{\mathbf{h}}}$ (hidden) states at time $t$
- $d_{\mathbf{h}}$ hidden dimension
- Read out $\mathbf{w}_{\texttt{out}} \in \mathbb{R}^{d_{\mathbf{h}}}$

$$y^{(t)} = \mathbf{w}_{\texttt{out}}^{\top}\mathbf{h}^{(t)}$$

final layer of your model

classification "head"

regression "head"

hidden-state

# Recurrent unit (with linear activation)

- Recurrent unit

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{w_{in}}x^{(t)}$$

- $\mathbf{h}^{(t)} \in \mathbb{R}^{d_\mathtt{h}}$ (hidden) states at time $t$
- $d_\mathtt{h}$ hidden dimension
- Read out $\mathbf{w_{out}} \in \mathbb{R}^{d_\mathtt{h}}$

$$y^{(t)} = \mathbf{w_{out}}^{\top}\mathbf{h}^{(t)}$$

- $\mathbf{W_{rec}} \in \mathbb{R}^{d_\mathtt{h} \times d_\mathtt{h}}$, $\mathbf{w_{in}} \in \mathbb{R}^{d_\mathtt{h}}$ and $\mathbf{w_{out}} \in \mathbb{R}^{d_\mathtt{h}}$ are parameters

# Recovering the autoregressive model

- Autoregressive (AR) model (with filter size 4)

$$y^{(t)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^\top \tilde{\mathbf{h}}^{(t-1)} \quad \text{where} \quad \tilde{\mathbf{h}}^{(t-1)} = \begin{bmatrix} x^{(t-1)} \\ x^{(t-2)} \\ x^{(t-3)} \end{bmatrix}$$

- Recurrent unit: how can we choose the parameters to recover the AR model?

Want this ↘

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}} \mathbf{h}^{(t-1)} + \mathbf{w_{in}} x^{(t)}$$

$$h^{(t)} := \begin{bmatrix} x^{(t)} \\ x^{(t-1)} \\ x^{(t-2)} \\ x^{(t-3)} \end{bmatrix} \qquad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

# Recovering the autoregressive model

- Autoregressive (AR) model (with filter size 4)

$$y^{(t)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^\top \tilde{\mathbf{h}}^{(t-1)} \quad \text{where} \quad \tilde{\mathbf{h}}^{(t-1)} = \begin{bmatrix} x^{(t-1)} \\ x^{(t-2)} \\ x^{(t-3)} \end{bmatrix}$$

- Recurrent unit: how can we choose the parameters to recover the AR model?

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}} \mathbf{h}^{(t-1)} + \mathbf{w_{in}} x^{(t)}$$

Want

$$h^{(t)} = \begin{pmatrix} x^{(t)} \\ x^{(t-1)} \\ x^{(t-2)} \\ x^{(t-3)} \end{pmatrix}$$

Assume
this is
achieved
up to
time t-1

then
$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$
$$\begin{bmatrix} x^{(t-1)} \\ x^{(t-2)} \\ x^{(t-3)} \\ x^{(t-4)} \end{bmatrix}$$

drop this

# Recovering the autoregressive model

- Autoregressive (AR) model (with filter size 4)

$$y^{(t)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^\top \tilde{\mathbf{h}}^{(t-1)} \quad \text{where} \quad \tilde{\mathbf{h}}^{(t-1)} = \begin{bmatrix} x^{(t-1)} \\ x^{(t-2)} \\ x^{(t-3)} \end{bmatrix}$$

- Recurrent unit: how can we choose the parameters to recover the AR model?

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{w_{in}}x^{(t)}$$

upshot

$$W_{out} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$W_{out}^\top h(t)$

recovers autoreg

choose

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 & \\ & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

ensure

$$h^{(t)} = \begin{bmatrix} x^{(t)} \\ x^{(t-1)} \\ x^{(t-2)} \end{bmatrix}$$

# Recovering the autoregressive model

*a single recurrent unit is already expressive*

- Autoregressive (AR) model (with filter size 4)

$$y^{(t)} = w_0 x^{(t)} + \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}^\top \tilde{\mathbf{h}}^{(t-1)} \quad \text{where} \quad \tilde{\mathbf{h}}^{(t-1)} = \begin{bmatrix} x^{(t-1)} \\ x^{(t-2)} \\ x^{(t-3)} \end{bmatrix}$$

- Recurrent unit: how can we choose the parameters to recover the AR model?

$$\mathbf{h}^{(t)} = f(x^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W}_{\texttt{rec}} \mathbf{h}^{(t-1)} + \mathbf{w}_{\texttt{in}} x^{(t)}$$

AR(4) can be implemented as a Rec Unit with $d_n = 4$

# Recurrent unit (with linear activation)

- Recurrent unit

$$\mathbf{h}^{(t)} = f(\mathbf{x}^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{W_{in}}\mathbf{x}^{(t)}$$

- $\mathbf{W_{rec}}$ and $\mathbf{W_{in}}$ are parameters
- $\mathbf{h}^{(t)}$ (hidden) states at time $t$
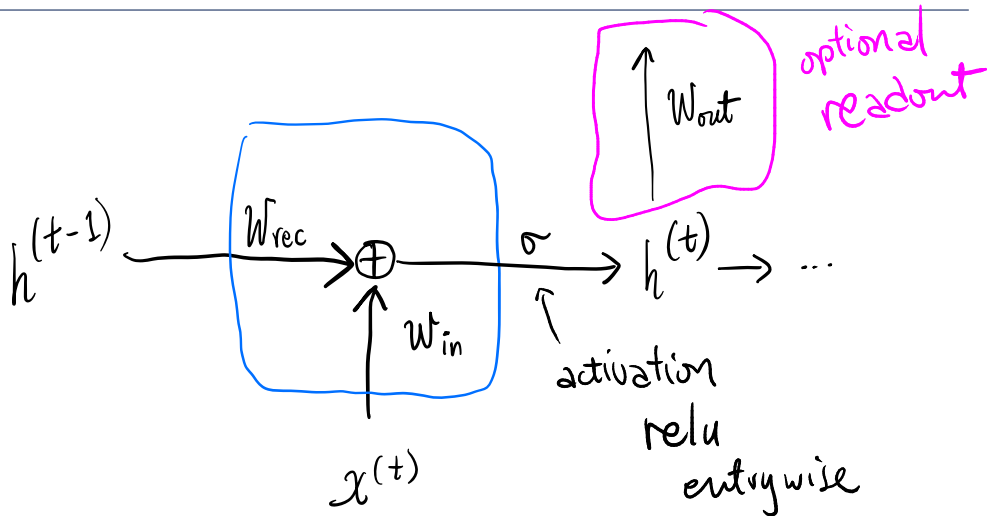
can be vectorial

# Recurrent unit (with linear activation)

- Recurrent unit

$$\mathbf{h}^{(t)} = f(\mathbf{x}^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \sigma(\mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{W_{in}}\mathbf{x}^{(t)})$$

- $\mathbf{W_{rec}}$ and $\mathbf{W_{in}}$ are parameters
- $\mathbf{h}^{(t)}$ (hidden) states at time $t$
- $\sigma$ can be the relu, hyperbolic tangent, or anything

more
trad

# Recurrent unit

# Recurrent unit

# Recurrent unit (with linear activation)

- Recurrent unit

$$\mathbf{h}^{(t)} = f(\mathbf{x}^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \sigma(\mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{W_{in}}\mathbf{x}^{(t)})$$

```
1 dim_hidden = 6
2 dim_input = 1
3 dim_output = 1
4
5 Wrec = torch.randn(dim_hidden, dim_hidden, requires_grad=True)
6 Win  = torch.randn(dim_hidden, dim_input, requires_grad=True)
7 Wout = torch.randn(dim_hidden, dim_output, requires_grad=True)
8
9 def recurrent_unit(x, h_prev):
10     # YOUR CODE HERE
11     raise NotImplementedError
```
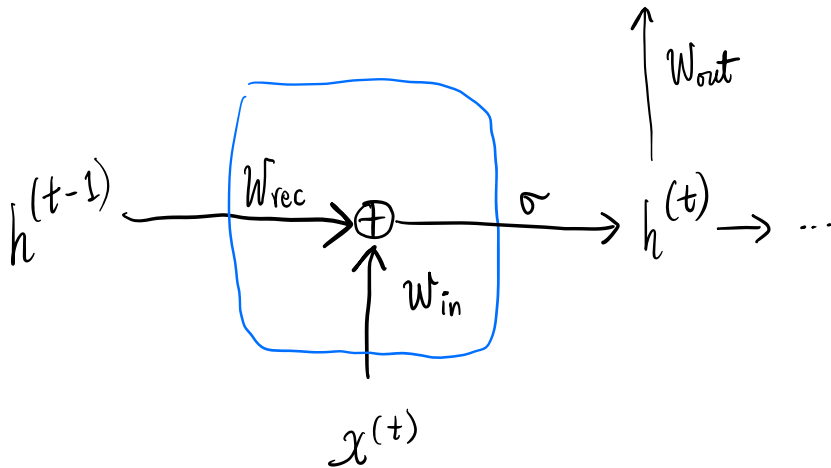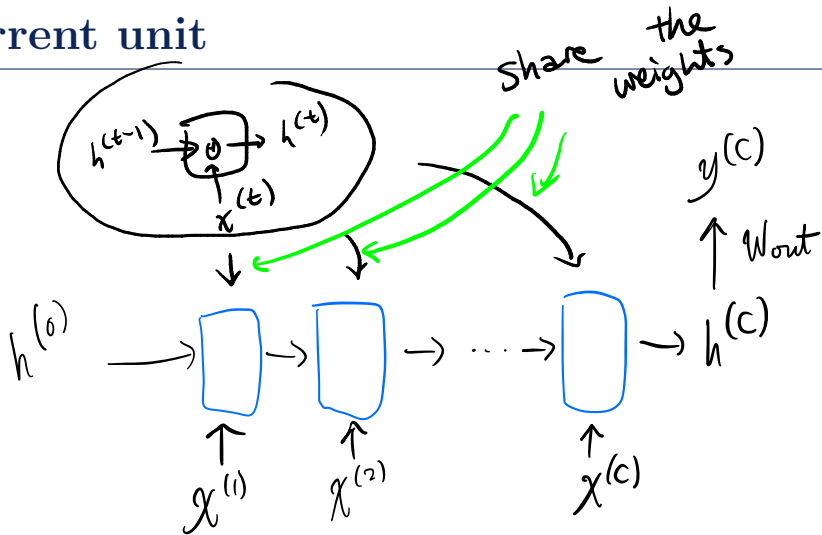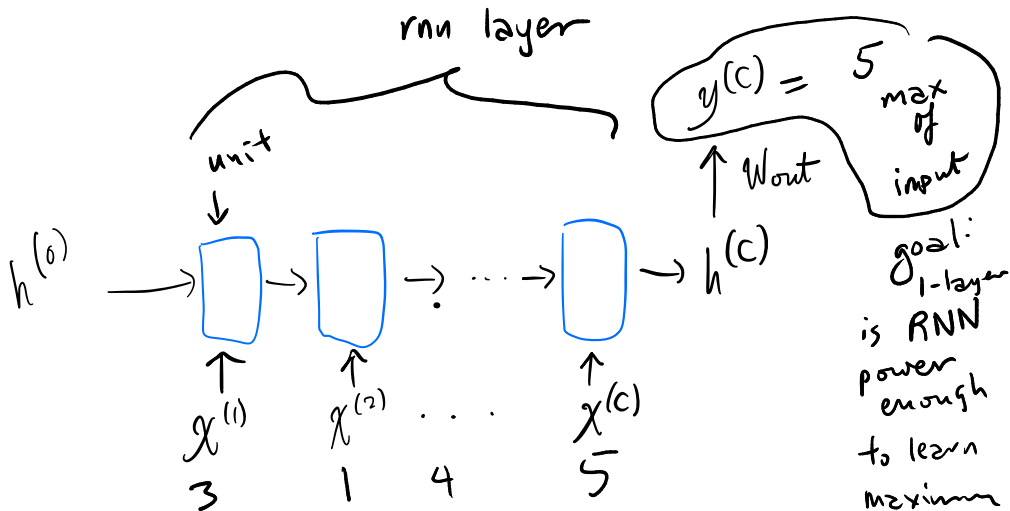
# Recurrent unit

# Computing the maximum of a sequence



rnn layer

unit

$h^{(0)}$

$y^{(c)} = 5$ max of input

$W_{out}$

$h^{(c)}$

$x^{(1)}$    $x^{(2)}$   ...   $x^{(c)}$

3     1    4      5

goal: 1-layer is RNN power enough to learn maximum

# Computing the maximum of a sequence

```
1  minibatch_size = 2**8
2  seq_length = 5
3  for epoch in range(num_epochs):
4      loss = 0
5      for _ in range(minibatch_size):
6          xs = torch.tensor(np.random.randn(seq_length, 1), dtype=torch.float32)
7          max_xs = torch.max(xs).item()
8          yhat = forward(xs, Wrec, Win, Wout)
9          loss += torch.abs(yhat - max_xs) # try the MSE loss!
10     loss = loss / minibatch_size
11
12     optimizer.zero_grad()
13     loss.backward()
14     optimizer.step()
```

"seq" = xs = [6, 2, ···]

max(seq)

yhat

(yhat - max_xs) ** 2

# In class exercise 2

torch.Relu

$$\mathbf{h}^{(t)} = f(\mathbf{x}^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \sigma(\mathbf{W_{rec}}\mathbf{h}^{(t-1)} + \mathbf{W_{in}}\mathbf{x}^{(t)})$$

given

```
1  def recurrent_unit(x, h_prev): # returns h_next
2      # YOUR CODE HERE
3      raise NotImplementedError
4
5
6  def forward(seq, Wrec, Win, Wout): # returns yhat
7      # YOUR CODE HERE
8      # hint:  use 'for x in seq: ...'
9      raise NotImplementedError
```

Wrec, Win

task 1

initialize $h^{(0)}$

# Vectorize the minibatch?

*triple for loop* (handwritten)
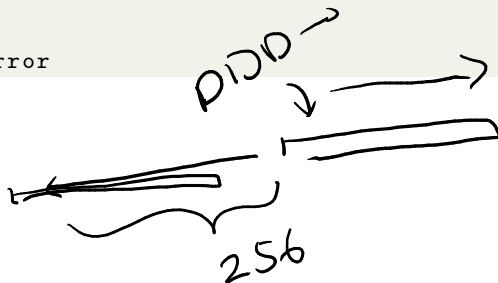
```
1  minibatch_size = 2**8          256
2  seq_length = 5
3  for epoch in range(num_epochs):
4      loss = 0
5      # CAN WE GET RID OF THE FOLLOWING FOR loop?
6      for _ in range(minibatch_size):
7          xs = torch.tensor(np.random.randn(seq_length, 1), dtype=torch.
   float32)
8          max_xs = torch.max(xs).item()
9          yhat = forward(xs, Wrec, Win, Wout)
10         loss += torch.abs(yhat - max_xs)
11     loss = loss / minibatch_size
12
13     optimizer.zero_grad()
14     loss.backward()
15     optimizer.step()
```

*Vectorized the minibatch* (handwritten)

```
1  def forward(seqs,Wrec, Win, Wout):
2      """
3      INPUT
4      seqs - a (minibatch_size, seq_length, dim_input) tensor
5
6      RETURN
7      yhat - a (minibatch_size, ) tensor
8      """
9      # YOUR CODE HERE
10     raise NotImplementedError
```

1

Possible to
achieve
zero error
(challenge)

DDD

256

# In class exercise 3

*Expressivity*
↳ universal
↳ approximate

gap

achievable via SGD adam

```
1  def forward(seqs,Wrec, Win, Wout):
2      """
3      INPUT
4      seqs - a (minibatch_size, seq_length, dim_input) tensor
5
6      RETURN
7      yhat - a (minibatch_size, ) tensor
8      """
9      # YOUR CODE HERE
10     raise NotImplementedError
```

there exists a set of weight $(W_{rec}, W_{in}, W_{out})$
st. we can get 0 error (forever)
upshot: RNNs are very expressive "universal approximation"

# An annoying loop

```python
def forward(seqs,Wrec, Win, Wout):
    """
    INPUT
    seqs - a (minibatch_size, seq_length, dim_input) tensor

    RETURN
    yhat - a (minibatch_size, ) tensor
    """
    h = torch.zeros(seqs.shape[0], dim_hidden)

    for x in seqs.transpose(0, 1): # CAN WE GET RID OF THIS?
        h = recurrent_unit(x, h, Wrec,Win)
    # [...]
```

*Vectorize over time?*

- No, at least not yet (active research area)

# An annoying loop

"( # of param )" square

not SGD
not adam
need 2nd deriv
Newton's

```python
def forward(seqs,Wrec, Win, Wout):
    """
    INPUT
    seqs - a (minibatch_size, seq_length, dim_input) tensor

    RETURN
    yhat - a (minibatch_size, ) tensor
    """
    h = torch.zeros(seqs.shape[0], dim_hidden)

    for x in seqs.transpose(0, 1): # CAN WE GET RID OF THIS?
        h = recurrent_unit(x, h, Wrec,Win)
    # [...]
```

Catch

- No, at least not yet (active research area)
- "Parallelizing non-linear sequential models over the sequence length" (Lim et al., 2024 ICLR)

σ ≠ identity

# An annoying loop

$$\frac{\partial f}{\partial \theta_i}$$

first order

$i \in \{\text{all params}\}$

$$(\# \text{param})^2 = \frac{\partial^2 f}{\partial \theta_i \partial \theta_j}$$

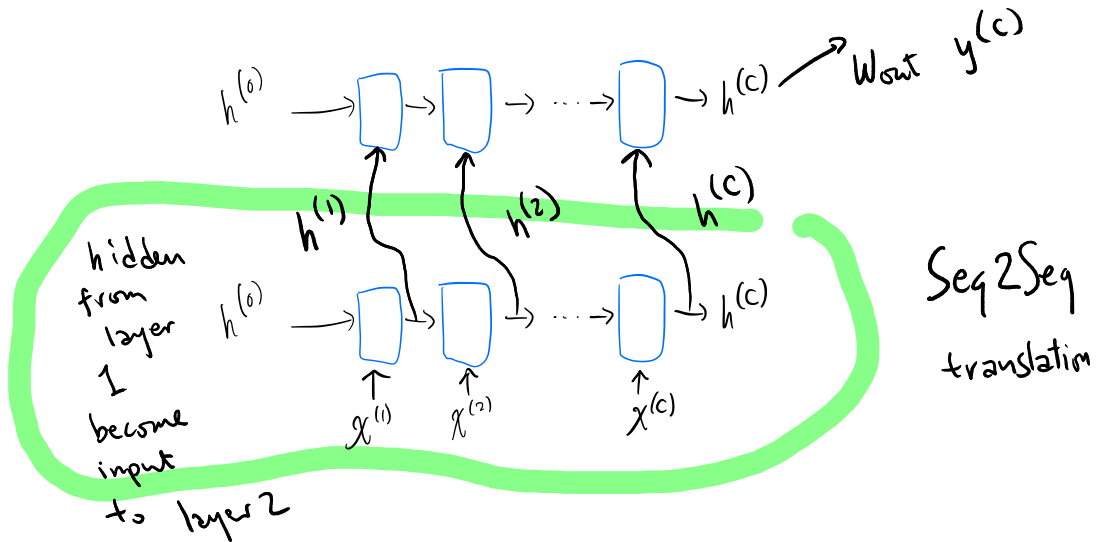$i, j \in \{\text{all params}\}$

```
1  def forward(seqs, Wrec, Win, Wout):
2      """
3      INPUT
4      seqs - a (minibatch_size, seq_length, dim_input) tensor
5
6      RETURN
7      yhat - a (minibatch_size, ) tensor
8      """
9      h = torch.zeros(seqs.shape[0], dim_hidden)
10
11     for x in seqs.transpose(0, 1): # CAN WE GET RID OF THIS?
12         h = recurrent_unit(x, h, Wrec, Win)
13     # [...]
```

- No, at least not yet (active research area)
- "Parallelizing non-linear sequential models over the sequence length" (Lim et al., 2024 ICLR)

↳ Neural ODE apply DL to physics

# Recurrent Neural Network with multiple layers



$h^{(0)} \rightarrow \square \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow h^{(c)} \nearrow W_{out} \; y^{(c)}$

$h^{(1)}$   $h^{(2)}$   $h^{(c)}$

hidden from layer 1 become input to layer 2

$h^{(0)} \rightarrow \square \rightarrow \square \rightarrow \cdots \rightarrow \square \rightarrow h^{(c)}$

$\uparrow x^{(1)}$   $\uparrow x^{(2)}$   $\uparrow x^{(c)}$

Seq 2 Seq translation

# Implementation

```
1  class recurrent_cell(nn.Module):
2      def __init__(self, input_dim, hidden_dim):
3          # [...]
4          self.hidden_dim = hidden_dim
5          self.input_to_hidden = nn.Linear(input_dim, hidden_dim)
6          self.hidden_to_hidden = nn.Linear(hidden_dim, hidden_dim)
7
8      def forward(self, x, h_prev):
9          h_next = torch.relu(self.input_to_hidden(x) + self.
    hidden_to_hidden(h_prev))
10         return h_next
```

# Implementation

```python
class RNN_layer(nn.Module):
    def __init__(self, embed_dim, hidden_dim):
        # [...]
        self.rnn_layer = recurrent_cell(embed_dim, hidden_dim)

    def forward(self, x_seq, h):
        outputs = []
        for t in range(x_seq.size(1)):
            h = self.rnn_layer(x_seq[:, t, :], h)
            outputs.append(h)

        # Stack outputs
        x_transformed = torch.stack(outputs, dim=1)
        return x_transformed
```

# Implementation

```
1  class RNN(nn.Module):
2      def __init__(self, vocab_size, embed_dim, hidden_dim, num_layers):
3          # [...]
4          self.embedding = nn.Embedding(vocab_size, embed_dim)
5          self.layers = nn.ModuleList([RNN_layer(embed_dim, hidden_dim)
   for _ in range(num_layers)])
6          self.classification_head = nn.Linear(embed_dim, vocab_size)
7
8      def forward(self, x):
9          x = self.embedding(x)
10         h = torch.zeros(x.size(0), self.layers[0].rnn_layer.hidden_dim,
   device=x.device)  # Initial hidden state
11
12         for rnn_layer in self.layers:
13             x = rnn_layer(x, h)
14
15         logits = self.classification_head(x)
16         return logits
```

# Seq-2-"one" Causal models

$$f(\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(t)}; \theta) = \mathbf{y}^{(t)}$$

### Causality assumption

The output $\mathbf{y}^{(t)}$ does not depend on future values of $\mathbf{x}^{(t')}$ where $t' > t$

# Seq-2-seq causal models

$$h^{(1)} \dots h^{(C)}$$

$$f(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}; \theta) = \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(C)} \quad \text{translation}$$

### Causality assumption

The output $\mathbf{y}^{(t)}$ does not depend on future values of $\mathbf{x}^{(t')}$ where $C \geq t' > t \geq 1$

transformer $\rightarrow$ Self-attention    is a Seq2seq model

$\quad\quad\quad\quad \hookrightarrow$ multilayer

Does not suffer from non-parallelization in time / seq length

## Sequences

*linear algebra-ification*

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \cdots & \mathbf{x}^{(C)} \end{bmatrix} \in \mathbb{R}^{d \times C}$$

where

$X^{(i)}$

- $\mathbf{x}^{(t)} \in \mathbb{R}^d$
- $t \in \{1, \ldots, C\}$,
- we should really write $\mathbf{X}^{(i)}$ where $i \in \{1, \ldots, N\}$ but hide from notation for convenience

# Self-attention

$$X \quad Y \quad Z \qquad \text{self} \quad X = Y = Z$$
can't plug in different things

- $$\text{self-attention}(\mathbf{X}; \theta) := \underbrace{\mathbf{W}^{(V)\top}} \mathbf{X} \, \text{softmax}\left( \underbrace{\mathbf{X}^\top \mathbf{W}^{(K)\top}}_{K} \underbrace{\mathbf{W}^{(Q)} \mathbf{X}}_{Q} \right) \in \mathbb{R}^{d \times C}.$$

- parameters

$$\theta^{(\text{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

$(C, d) \qquad (d, C)$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively) where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

$(C, C)$

- "**seq-2-seq**": maps the sequence $\mathbf{X}$ to another sequence $\text{attention}(\mathbf{X}; \theta)$:

- **Problem**: is this causal?

$$\frac{\text{softmax}(K^\top Q)}{(C, C)} - \text{column col sum to 1} \quad \text{stochastic}$$

# Self-attention

$(d,d)$    $(d,C)$      $[C,C)$

- 
$$\texttt{attention}(\mathbf{X};\theta) := \mathbf{W}^{(V)\top}\mathbf{X}^{(i)}\mathrm{softmax}\left(\mathbf{X}^\top\mathbf{W}^{(K)\top}\mathbf{W}^{(Q)}\mathbf{X}\right) \in \mathbb{R}^{d\times C}.$$

- parameters
$$\theta^{(\texttt{att})} = [\mathbf{W}^{(Q)}, \mathbf{W}^{(K)}, \mathbf{W}^{(V)}]$$

$(d,C)$

- ("$Q$, $K$, $V$" stands for "query", "key", "value", respectively) where

$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d\times d}.$$

- "**seq-2-seq**": maps the sequence $\mathbf{X}$ to another sequence $\texttt{attention}(\mathbf{X};\theta)$:
- **Problem**: is this causal?

input.shape = output.shape

# Is this causal?

$$\boxed{\exp K Q} \leftarrow \text{softmax}(K^T Q)$$

start here
my rec

$C = 4$

In class exercise 4

$$\text{attention}(X) =$$

$x_0 \; x_1 \; x_2 \; x_3$

$$W_v \left[ \; | \; | \; | \; | \; \right]$$

$\underbrace{\qquad}_{X}$

force it zero

$\Sigma$ to 1

$$= \left( | \; | \; | \; | \right)$$

$\underbrace{\text{att}(X)[0,2]}_{\substack{\uparrow \\ \text{2nd} \\ \text{context}}} = W_v \left( p_0 x_0 + p_1 x_1 + p_2 x_2 + p_3 x_3 \right) .$

$\begin{bmatrix} p_0 \\ p_1 \\ p_2 \\ \boxed{p_3} \end{bmatrix} \neq 0$

## Is this causal?

In class exercise 4

$$W_V \times \left( \text{softmax}(KQ) * M \right)$$

Causal mask
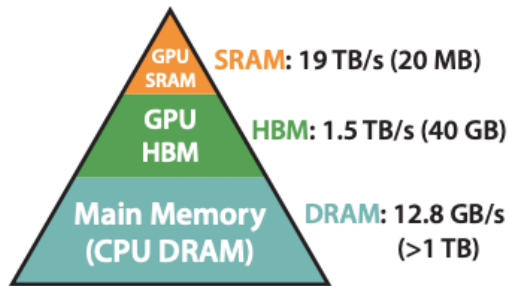
not quite
right

but on
right
track

# Flash attention



**Memory Hierarchy with Bandwidth & Memory Size**

From Dao et al 2022 (Flash Attention)

# Flash attention

**Algorithm 0** Standard Attention Implementation

**Require:** Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.
1: Load $\mathbf{Q}, \mathbf{K}$ by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^{\top}$, write $\mathbf{S}$ to HBM.
2: Read $\mathbf{S}$ from HBM, compute $\mathbf{P} = \text{softmax}(\mathbf{S})$, write $\mathbf{P}$ to HBM.
3: Load $\mathbf{P}$ and $\mathbf{V}$ by blocks from HBM, compute $\mathbf{O} = \mathbf{P}\mathbf{V}$, write $\mathbf{O}$ to HBM.
4: Return $\mathbf{O}$.

From Dao et al 2022 (Flash Attention)

# References I