

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

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

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

Sequences

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

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

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)

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?)

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?)

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?)
- Task: regression (how favorable is this product review?)

Causal models

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

Causality assumption

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

Prediction/forecasting/generative models

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

-

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

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

Autoregressive model

- For $t \geq s$

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

Autoregressive model

- For $t \geq s$

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

- Loss

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

Autoregressive model

- For $t \geq s$

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

- Loss

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

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

Autoregressive model

- For $t \geq s$

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

- Loss

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

- Solve for w_0, \dots, w_{s-1}, b
- Next: `seq2col`

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

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

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.] )
```


Another way to look at autoregressive models

- Drop the bias term for simplicity

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

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 $t - 1$)

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

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 $t - 1$)

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

- Next: directly model the hidden state

Recurrent unit (with linear activation)

- Recurrent unit

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

Recurrent unit (with linear activation)

- Recurrent unit

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

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

Recurrent unit (with linear activation)

- Recurrent unit

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

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

Recurrent unit (with linear activation)

- Recurrent unit

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

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

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

Recurrent unit (with linear activation)

- Recurrent unit

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

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

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

- $\mathbf{W}_{\text{rec}} \in \mathbb{R}^{d_h \times d_h}$, $\mathbf{w}_{\text{in}} \in \mathbb{R}^{d_h}$ and $\mathbf{w}_{\text{out}} \in \mathbb{R}^{d_h}$ 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?

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

Recurrent unit (with linear activation)

- Recurrent unit

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

- \mathbf{W}_{rec} and \mathbf{W}_{in} are parameters
- $\mathbf{h}^{(t)}$ (hidden) states at time t

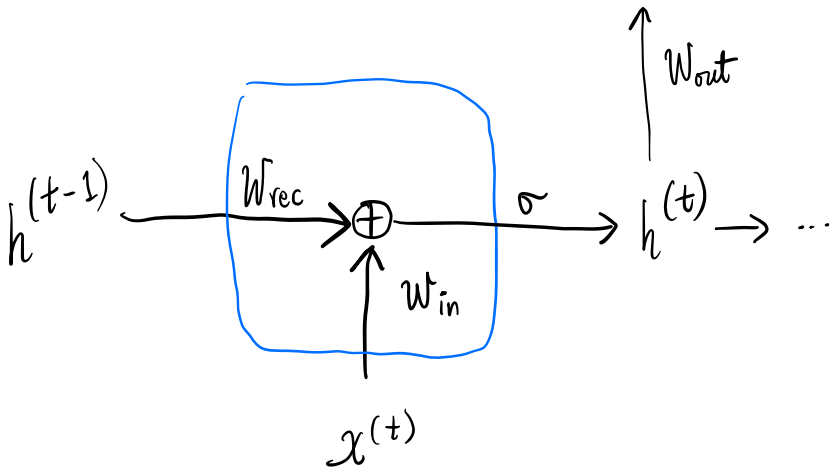
Recurrent unit (with linear activation)

- Recurrent unit

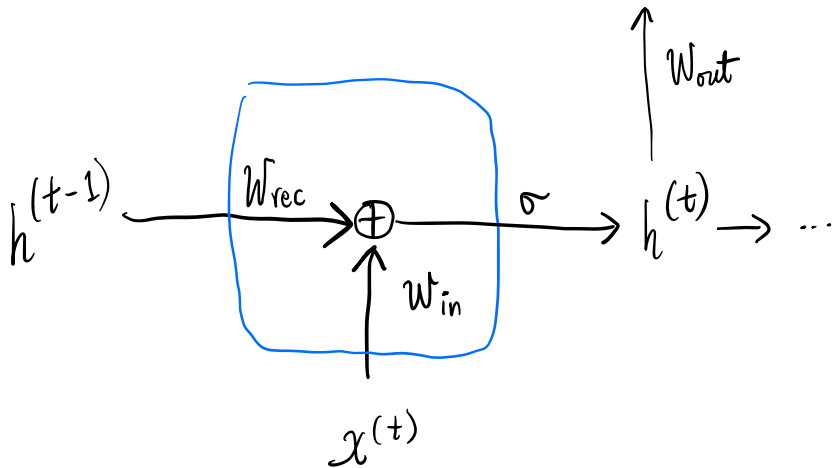
$$\mathbf{h}^{(t)} = f(\mathbf{x}^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \sigma(\mathbf{W}_{\text{rec}}\mathbf{h}^{(t-1)} + \mathbf{W}_{\text{in}}\mathbf{x}^{(t)})$$

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

Recurrent unit



Recurrent unit



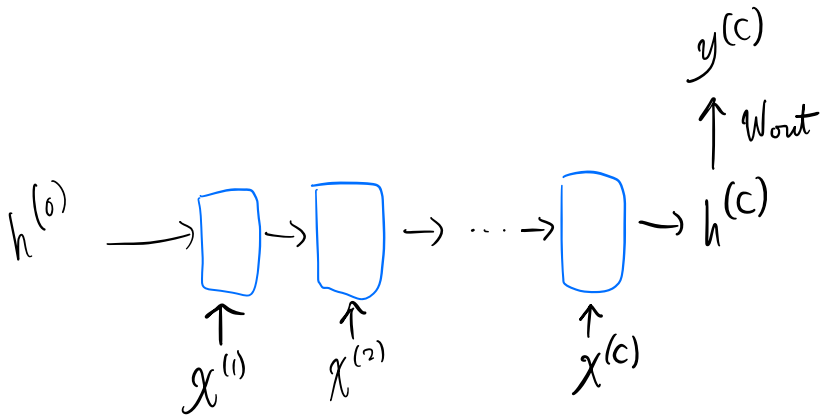
Recurrent unit (with linear activation)

- Recurrent unit

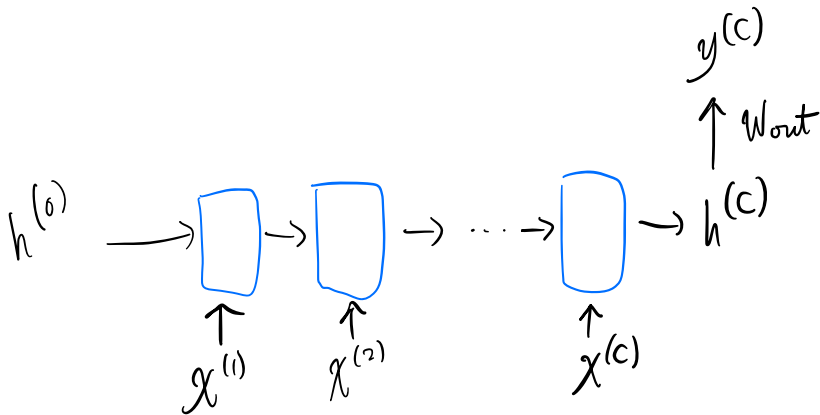
$$\mathbf{h}^{(t)} = f(\mathbf{x}^{(t)}, \mathbf{h}^{(t-1)}; \boldsymbol{\theta}) = \sigma(\mathbf{W}_{\text{rec}}\mathbf{h}^{(t-1)} + \mathbf{W}_{\text{in}}\mathbf{x}^{(t)})$$

```
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



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()
```

In class exercise 2

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

```
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
```

Vectorize the minibatch?

```
1 minibatch_size = 2**8
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()
```

In class exercise 3

```
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
```

An annoying loop

```
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     # [...]
```

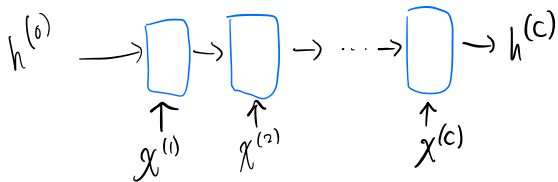
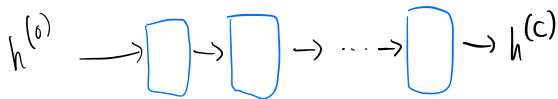
- No, at least not yet (active research area)

An annoying loop

```
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)

Recurrent Neural Network with multiple layers



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

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

Seq-2-“one” Causal models

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

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

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

Causality assumption

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

Sequences

$$\mathbf{X} = [\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \dots \quad \mathbf{x}^{(C)}]$$

where

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

Self-attention

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

- parameters

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

- (“ Q , K , V ” stands for “query”, “key”, “value”, respectively)
where

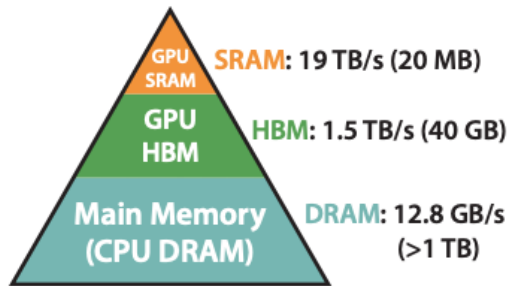
$$\mathbf{W}^{(Q)} \text{ and } \mathbf{W}^{(K)} \text{ and } \mathbf{W}^{(V)} \in \mathbb{R}^{d \times d}.$$

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

Is this causal?

In class exercise 4

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 $\mathbf{S} = \mathbf{QK}^\top$, 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 $\mathbf{O} = \mathbf{PV}$, write \mathbf{O} to HBM.
 - 4: Return \mathbf{O} .
-

From Dao et al 2022 (Flash Attention)

References I
