# DS 4440 Recurrent Neural Networks
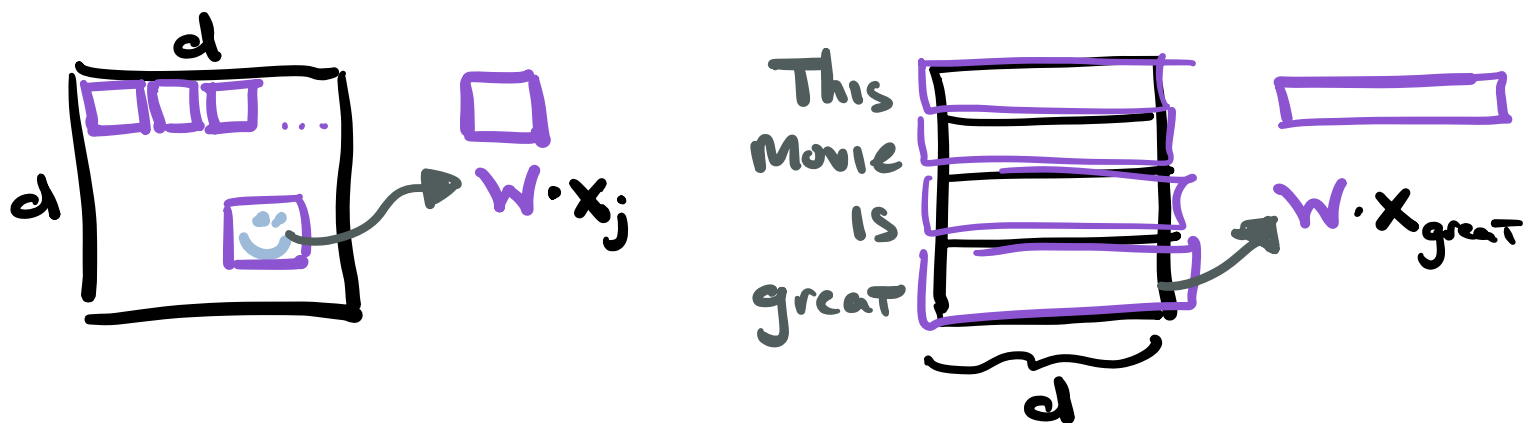
Last few lectures... ConvNets



These **filters** are position invariant.

But data often is **Sequentially Structured** — e.g.

- Language
- Physiological data
- Genetic Code
- Stock prices

$$X_\tau \sim P(X_\tau | X_{\tau-1} \ldots X_1)$$
$$Y_\tau = f(X_\tau \ldots X_1)$$

How can we model an arbitrary length Sequence with a finite Set of parameters?

## Markov Assumption

$$P(x_\tau | x_{\tau-1} \cdots x_1) = P(x_\tau | x_{\tau-1} \cdots x_{t-k})$$

**RNNs** learn to induce a fixed length representation of $x_1 \cdots x_T$.
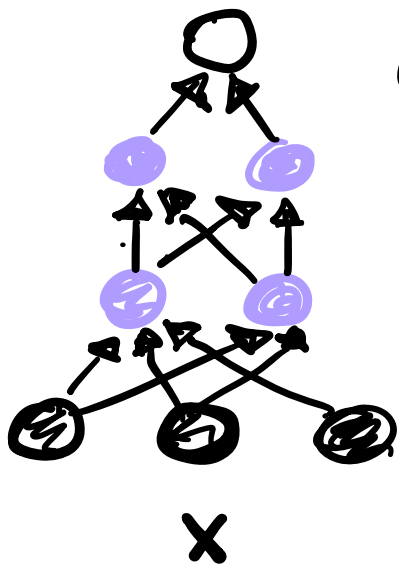
## Example Learning to Count in Variable length Sequences

$$
\begin{bmatrix}
[a\ b\ b] \\
[b\ b\ a\ a\ a\ b] \\
[a\ a\ b\ a\ b\ b\ b\ b]
\end{bmatrix}
\begin{bmatrix}
1 \\
0 \\
0
\end{bmatrix}
\begin{cases}
1 \text{ if Count(b)} \\
\quad \text{is even} \\
0 \text{ otherwise}
\end{cases}
$$

$X$ $\qquad\qquad$ $Y$

Unclear how we could model
this with MLPs or ConvNets.
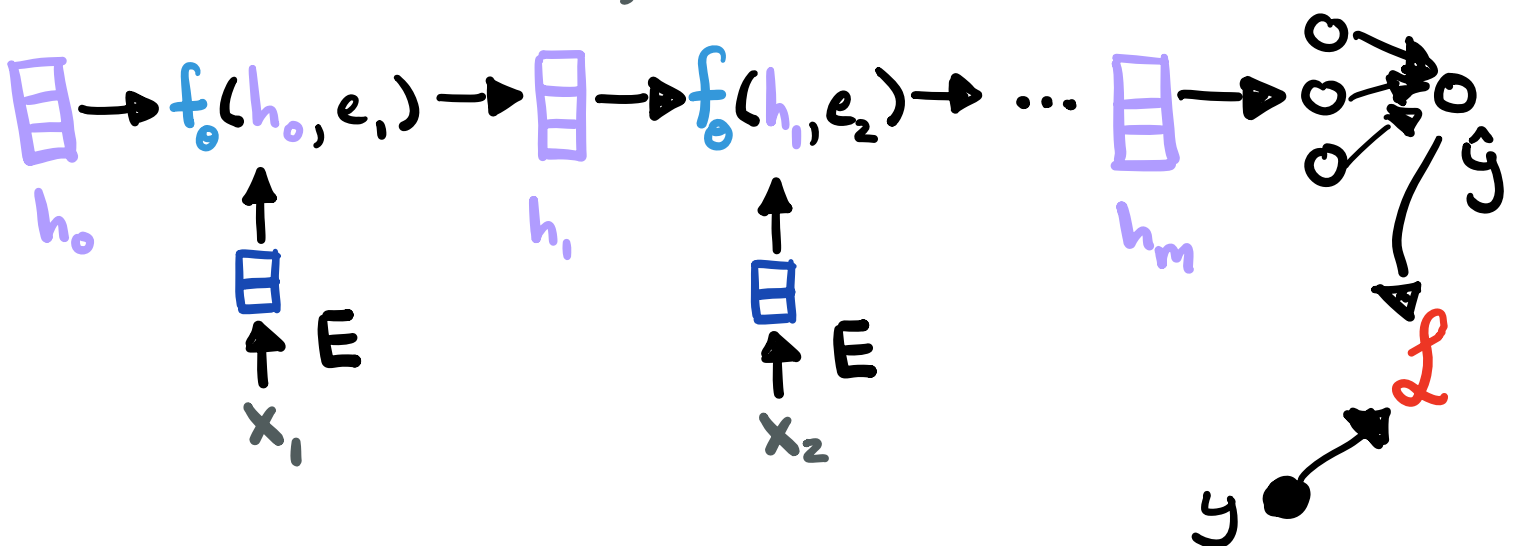We need State.

MLPs have hidden units



X

Q. Why doesn't this suffice?

A. Need to process X in
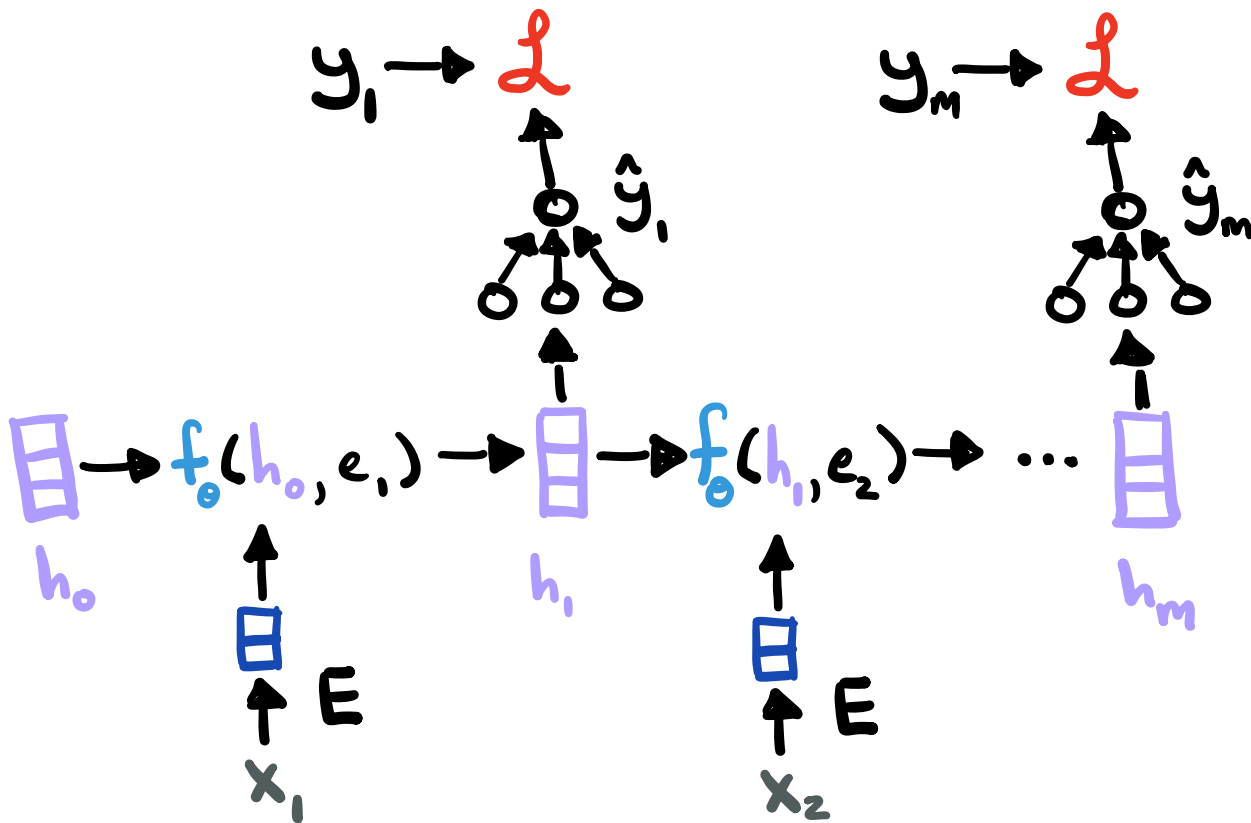   order, update $h$ at
   each step.

## RNNs

$X = \{x_1, x_2, \dots x_m\}, \quad y \in \mathbb{R}$

$f_\theta$ is shared.

Sequence Tagging: $\hat{y}$ at each $\tau$.

$x = \{x_1, x_2, \dots x_m\}, \quad y = \{y_1, y_2 \dots y_m\}$



Formally (batched)

$x_\tau \in \mathbb{R}^{b \times d}$    $h_\tau \in \mathbb{R}^{b \times h}$

batch size   input dims      hidden size

$$h_{\tau+1} \leftarrow \sigma(x_\tau W_x + h_\tau W_h + b_h)$$

$(b \times h)$      $(b \times d)(d \times h)$   $(b \times h)(h \times h)$    $(1 \times h)$

(See CoLab Notebook)

# Backprop through time

$h_0$

$X_1 \rightarrow \boxed{W^h} \rightarrow h_1$

$X_2 \rightarrow \boxed{W^h} \rightarrow h_2$

$X_3 \rightarrow \boxed{W^h} \rightarrow h_3 \rightarrow \boxed{W^o} \rightarrow \hat{y} \rightarrow \mathcal{L} \leftarrow y$

Assume Simple Model

$$\hat{y} = W^o \cdot h_3$$

$$\nabla_{W^o} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \nabla_{W^o} \hat{y}$$

What about $\nabla_{W^h}$?

$$\nabla_{W^h} = \sum_{\tau=1}^{T} \nabla_{W^h} h_\tau \cdot \nabla_{h_\tau} \mathcal{L}$$

$$\nabla_{h_3} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \nabla_{h_3} \hat{y}$$

$$= W^o \in \mathbb{R}^d$$

But $\nabla_{h_2}$ More Complicated

$$\nabla_{h_2}\mathcal{L} = (\nabla_{h_2} h_3)\cdot(\nabla_{h_3}\mathcal{L})$$

$$d\begin{bmatrix} \frac{\partial h_{3_1}}{\partial h_{2_1}} & \cdots & \frac{\partial h_{3_1}}{\partial h_{2_d}} \\ & \cdots & \\ \frac{\partial h_{3_d}}{\partial h_{2_1}} & & \frac{\partial h_{3_d}}{\partial h_{2_d}} \end{bmatrix} \text{Jacobian}$$

$$d$$

## Problem long sequences require repeat multiplies. Consider

$$\nabla_{h_1}\mathcal{L} = (\nabla_{h_1} h_2)\cdot(\nabla_{h_2} h_3)\cdot(\nabla_{h_3}\mathcal{L})$$

... And this is a Toy example.

Gradients $\nabla$ Tend to {Explode} or VANISH.

Can address Explosions via gradient Clipping:

If $\text{Norm}(\nabla_\theta) \geq \tau$

$$\nabla_\theta \leftarrow \nabla_\theta / 2$$