# RNNs

The problem: Want to model sequences of arbitrary length.

- Stock prices
- Physiological data
- Genetic code

...

Suppose we have $x_1, x_2, ..., x_{t-1}$ and want to predict $x_t$.

$$X_t \sim P(x_t | x_1, x_2, ..., x_{t-1})$$

<span style="color:purple">Size of this depends on $t$!</span>

So far we have handled via Markov assumption

$$P(x_t | x_1 ... x_{t-1}) = P(x_t | x_{t-1})$$

$$P(x_1, x_2, ..., x_T) = P(x_1) P(x_2|x_1) P(x_3|x_2) ... P(x_T|x_{T-1})$$

<span style="color:purple">First-order Markov</span>

Pros?
- Terms will (likely) be well estimated

- Very simple

Cons?

- Obviously wrong

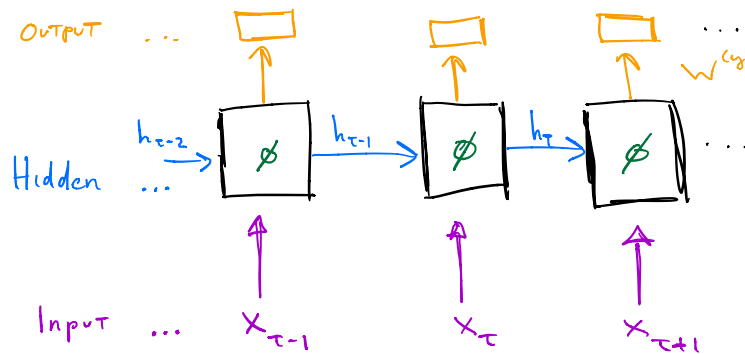Enter recurrent neural networks (RNNs). These pack history into a fixed size vector $h_t$ s.t.

$$P(x_t | x_1, x_2, \ldots, x_{t-1}) = P(x_t | h_t)$$

Context vector, updated at each time step (word)

Activation

$$h_t = \phi \{ x_t W^{(x)} + h_{t-1} W^{(h)} + b^{(h)} \}$$

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

Word at Position $t$

last hidden state

Model parameters

(figure below)

- The model is recurrent in that it consumes its own output from the previous step

- Consequently, it can be passed over inputs of arbitrary length: the number of params is fixed!
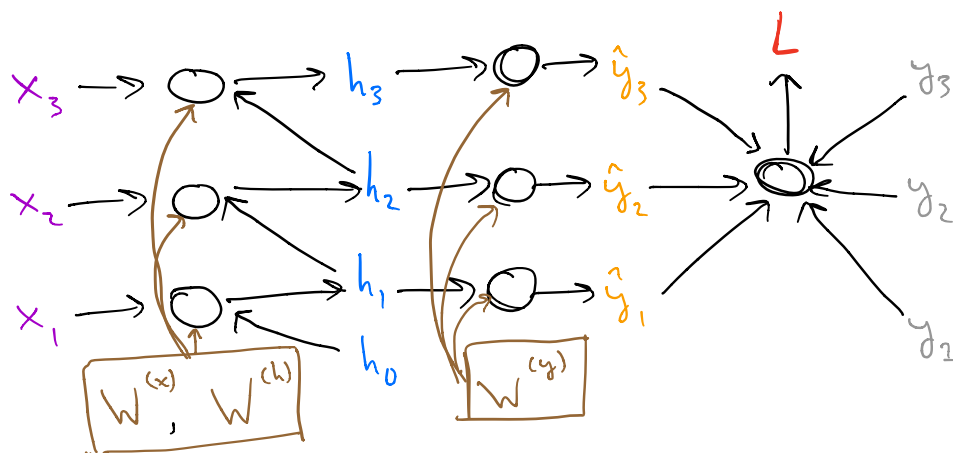
Downside Requires Sequentially processing input!

Can we avoid This?

The output layer will be some function of
h. In language modeling this might just
be a dense + SoftMax over the vocab.

Let's think about backprop in RNNs.

$$L = \frac{1}{T} \sum_{\tau=1}^{T} \mathcal{L}(\hat{y}_\tau, y_\tau)$$



$$\frac{\partial L}{\partial \hat{y}_\tau} = \frac{1}{T} \cdot \frac{\partial \mathcal{L}(\hat{y}_\tau, y_\tau)}{\partial \hat{y}_\tau}$$

$$\frac{\partial L}{\partial W^{(y)}} = \sum_{\tau=1}^{T} \frac{\partial L}{\partial \hat{y}_\tau} \cdot \frac{\partial \hat{y}_\tau}{\partial W^{(y)}}$$

The hidden vectors a bit trickier due to recurrence. Consider the last state $h_3$

$$\frac{\partial L}{\partial h_3} = \frac{\partial L}{\partial \hat{y}_3} \cdot \frac{\partial \hat{y}_3}{\partial h_3}$$

How about $h_2$? Gradient flows in from $\hat{y}_2$ and from $\hat{y}_3$ via $h_3$. So we add these

$$\frac{\partial L}{\partial h_2} = \frac{\partial L}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial h_2} + \frac{\partial L}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2}$$

Finally, similar for $h_1$

$$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial h_1} + \frac{\partial L}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1}$$

These are used to derive gradients for $W^{(h)}$ and $W^{(x)}$:

$$\frac{\partial L}{\partial W^{(h)}} = \sum_{\tau=1}^{T} \frac{\partial L}{\partial h_\tau} \cdot \frac{\partial h_\tau}{\partial W^{(h)}}$$

As sequences grow, gradients must travel "backward" further and further. Consider

$$\left\| \frac{\partial h_\tau}{\partial h_k} \right\| = \left\| \prod_{j=k+1}^{\tau} \frac{\partial h_j}{\partial h_{j-1}} \right\|$$

Jacobians!

Assume upper-bound

$$\left\| \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq \beta$$

$$\longrightarrow \leq \beta^{\tau-k}$$

So if $t-k$ is big we can see how this might be problematic!

Specifically: Vanishing when the error signal attenuates. Alternatively, gradients may explode for the same reason.

Jane walked into the room. John walked in too. It was late in the day. John said hi to ____.
    [ex. from Socher, cs224d]

Next time, we'll intro RNN variants that try to address this.