

STATS 507

Data Analysis in Python

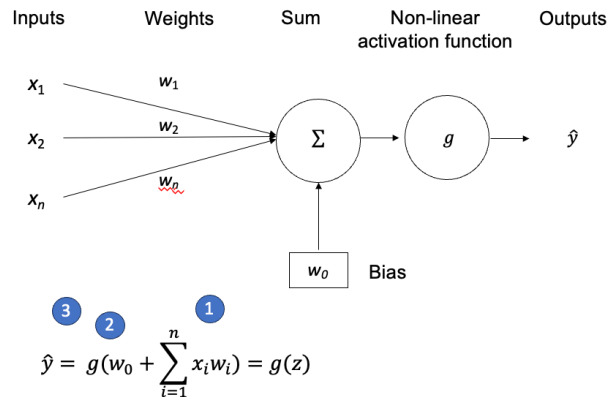
Week13-1: Deep Sequential Modeling (RNN, LSTM)

Dr. Xian Zhang

Adapted from slides by Ava Amini
MIT Introduction to Deep Learning

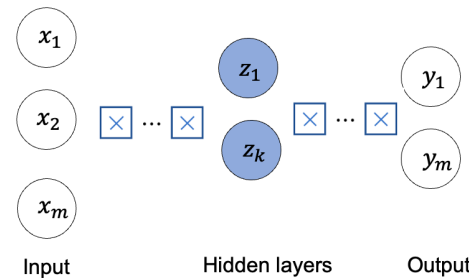
Recap: Core Foundation Review

The perceptron



- Structural building blocks
- Numerical operator
- Nonlinear activation functions

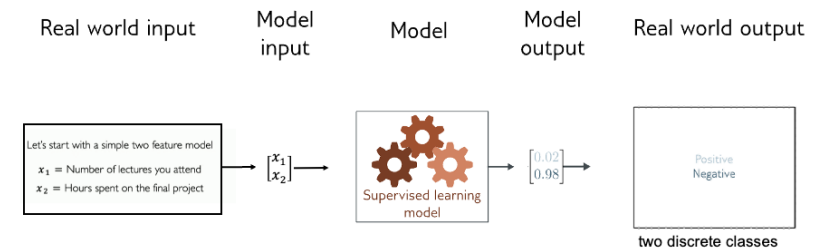
Neural Networks



- Stacking Perceptions to form neural networks (MLP)
- Optimization through backpropagation

Training in Practice

Applying DNN: Will I pass this class?



- Adaptive learning rate
- Batching
- Regularization

Recap: Gradient descent

Full Batch

- Batch size = N
- 1 update per epoch
- Use all the data

$$\frac{\partial L}{\partial \mathbf{W}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial L_i(\mathbf{W})}{\partial \mathbf{W}}$$

- Very slow updates
- Most stable
- Can stuck in local minima

Mini-batch GD

- Batch size = bs
- $N/32$ updates per epoch
- Use sub dataset

$$\frac{\partial L}{\partial \mathbf{W}} = \frac{1}{B} \sum_{i=1}^B \frac{\partial L_i(\mathbf{W})}{\partial \mathbf{W}}$$

- More accurate estimation of gradient
- Smoother convergence
- Allow for larger learning rates

Stochastic gradient descent

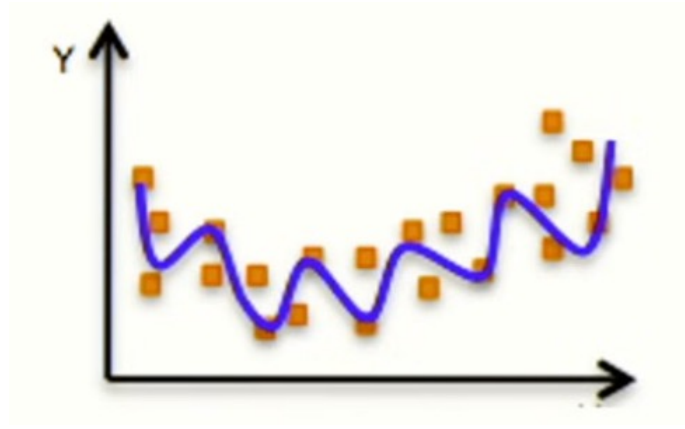
- Batch size = 1
- N updates per epoch
- Use single sample

$$\frac{\partial L_i}{\partial \mathbf{W}}$$

- Most efficient
- Very noisy updates
- Unstable

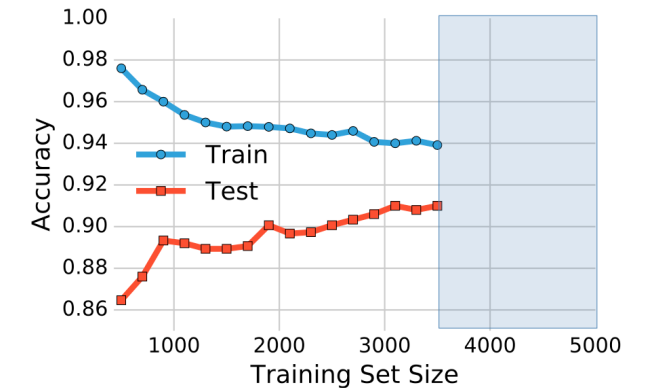
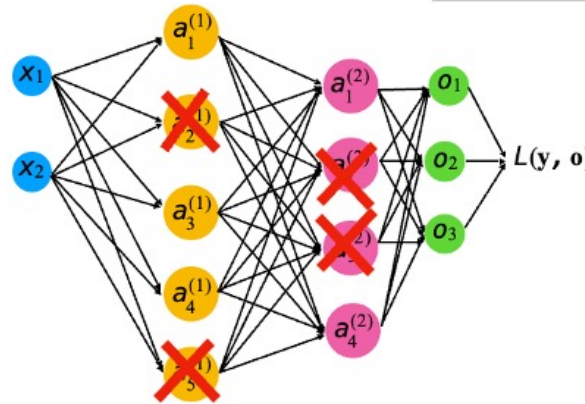
Recap: Training in practice

Model is doing really well on training data, but very badly on test data



Overfitting

Use regulation to discourage complex models



Slide Credit: Alexander Amini

Modified from MIT open course: 6.S191

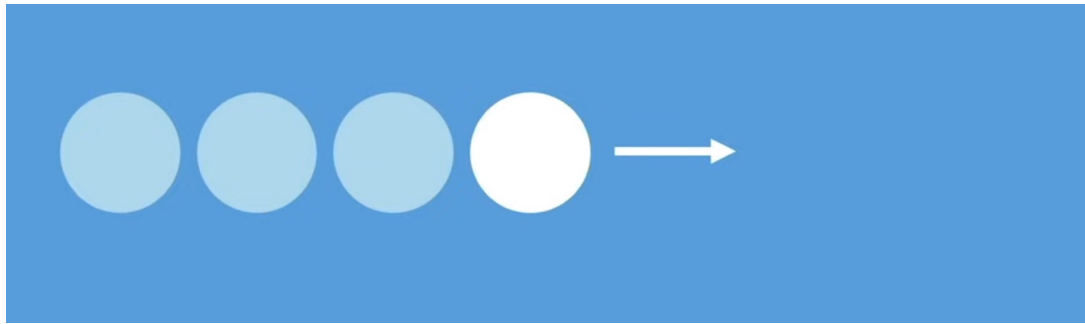
1. What is sequence modeling

2. RNN

3. LSTM

Intro: Sequence Data?

Given an image of a ball, can you predict where it will go next?



Audio:



Stock Market:



Language: This is Stats 507, where we learn data analytics using Python...

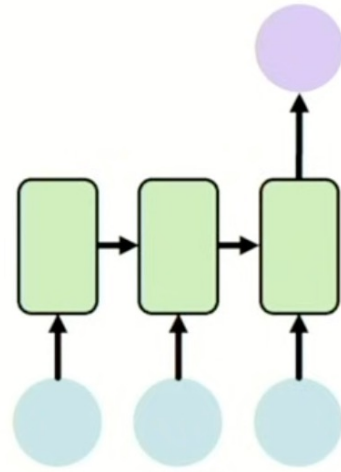
Intro: Sequence modeling applications



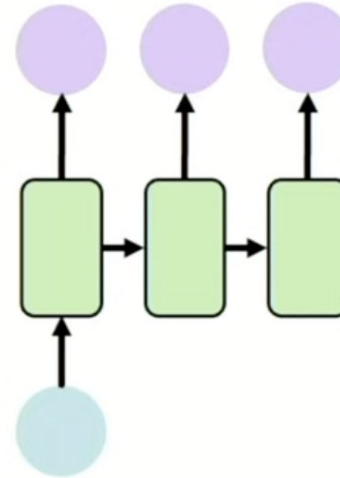
One to One
Binary Classification



"Will I pass this class?"
Student → Pass?



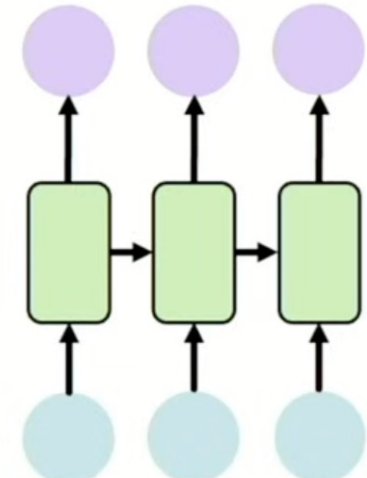
Many to One
Sentiment Classification



One to Many
Image Captioning



"A baseball player throws a ball."

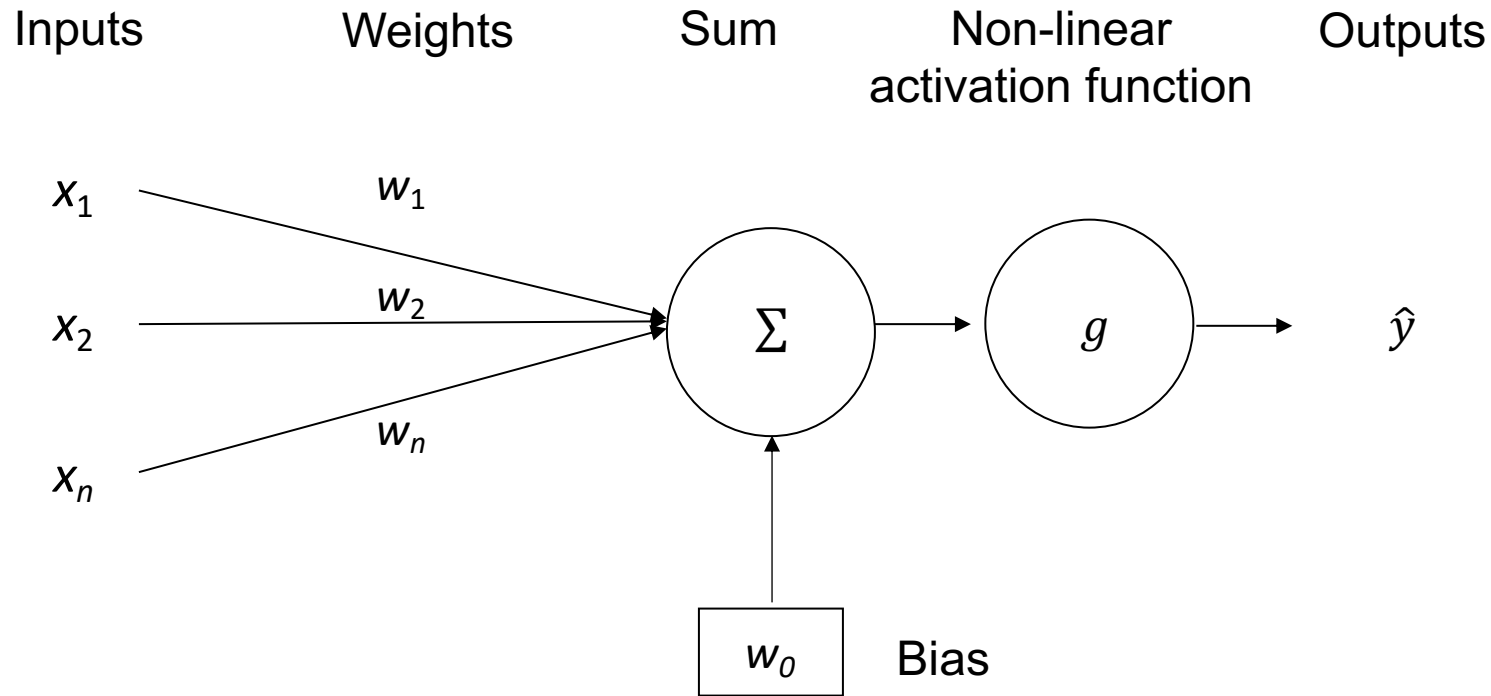


Many to Many
Machine Translation



How to build models for sequence?

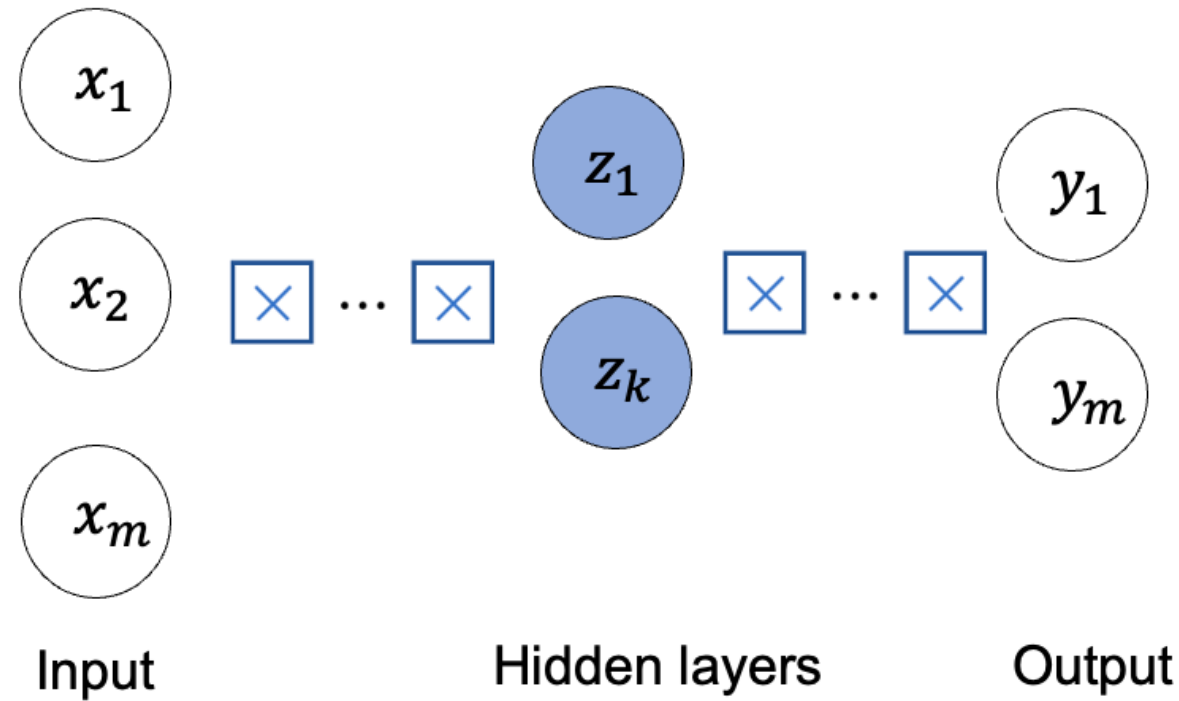
The perceptron (revisited)



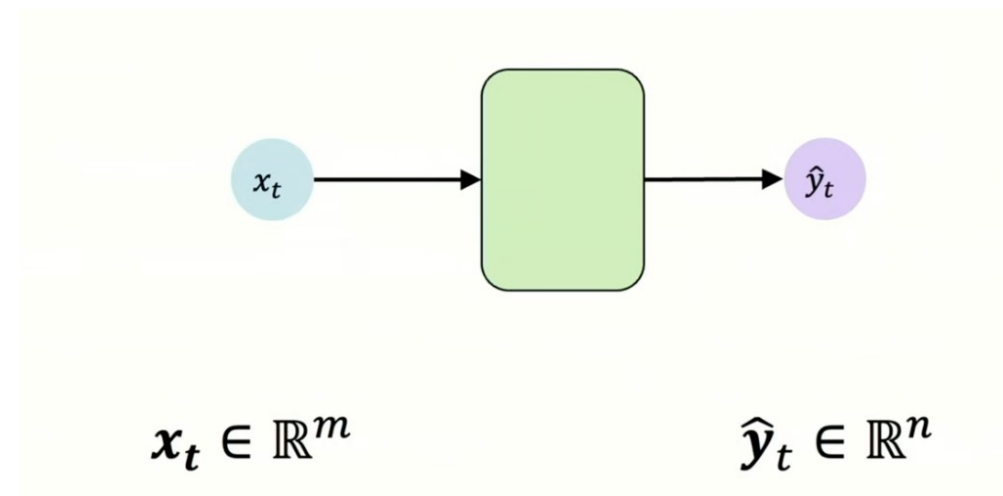
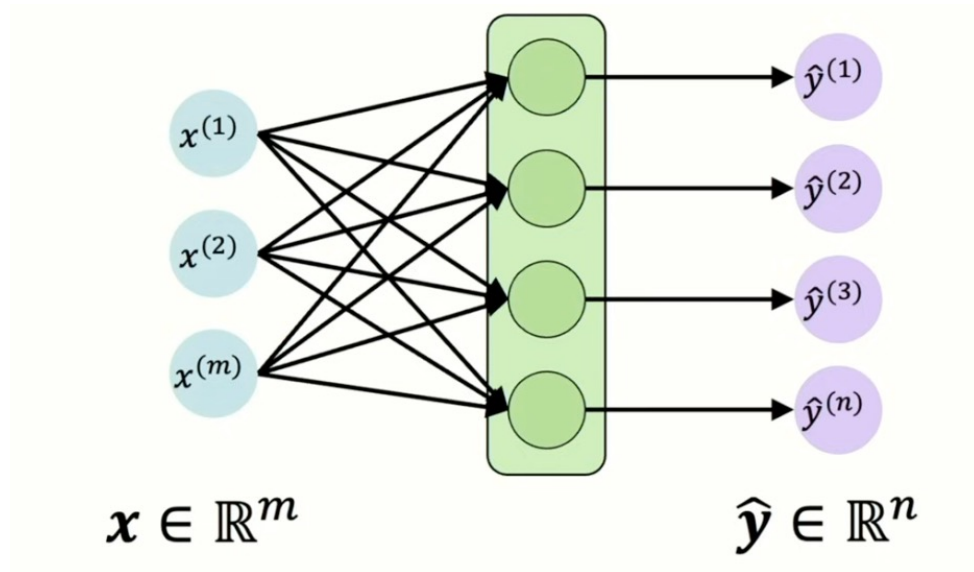
$$\hat{y} = g(w_0 + \sum_{i=1}^n x_i w_i) = g(z)$$

Diagrammatic annotations for the equation above: a blue circle with '3' is above the \hat{y} ; a blue circle with '2' is above the g ; a blue circle with '1' is above the x_i term.

Feed-forward Networks (revisited)

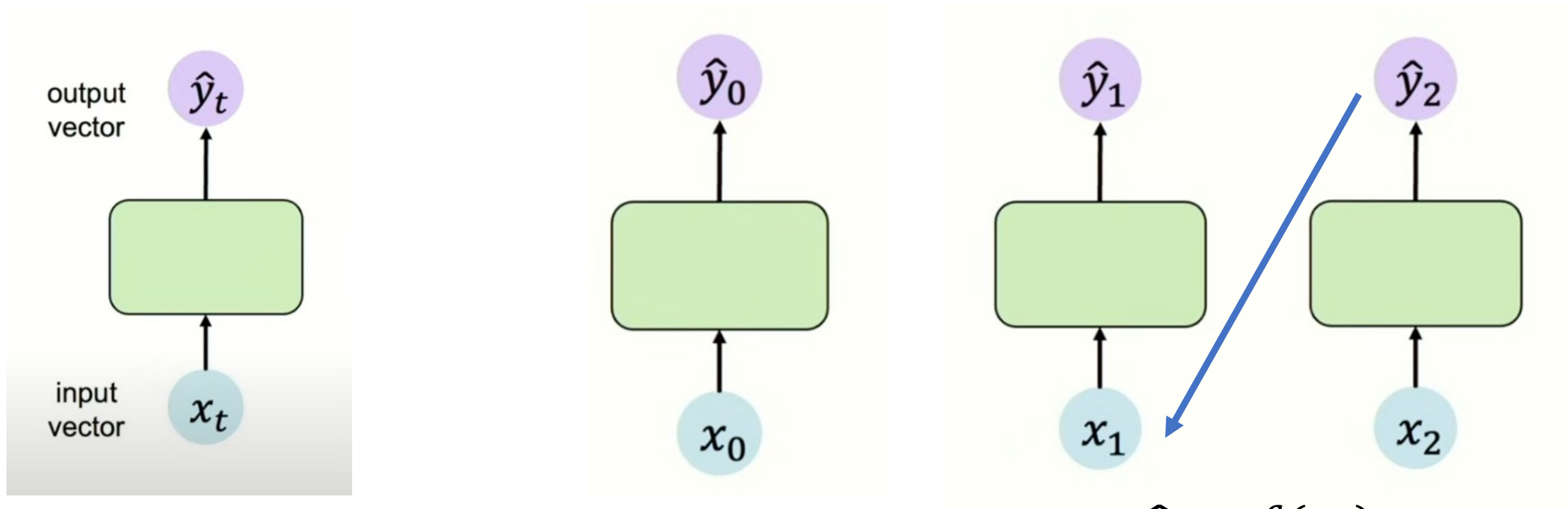


Adding the notion of time



Handling individual time steps

Apply the same model stepwise to **time** slice



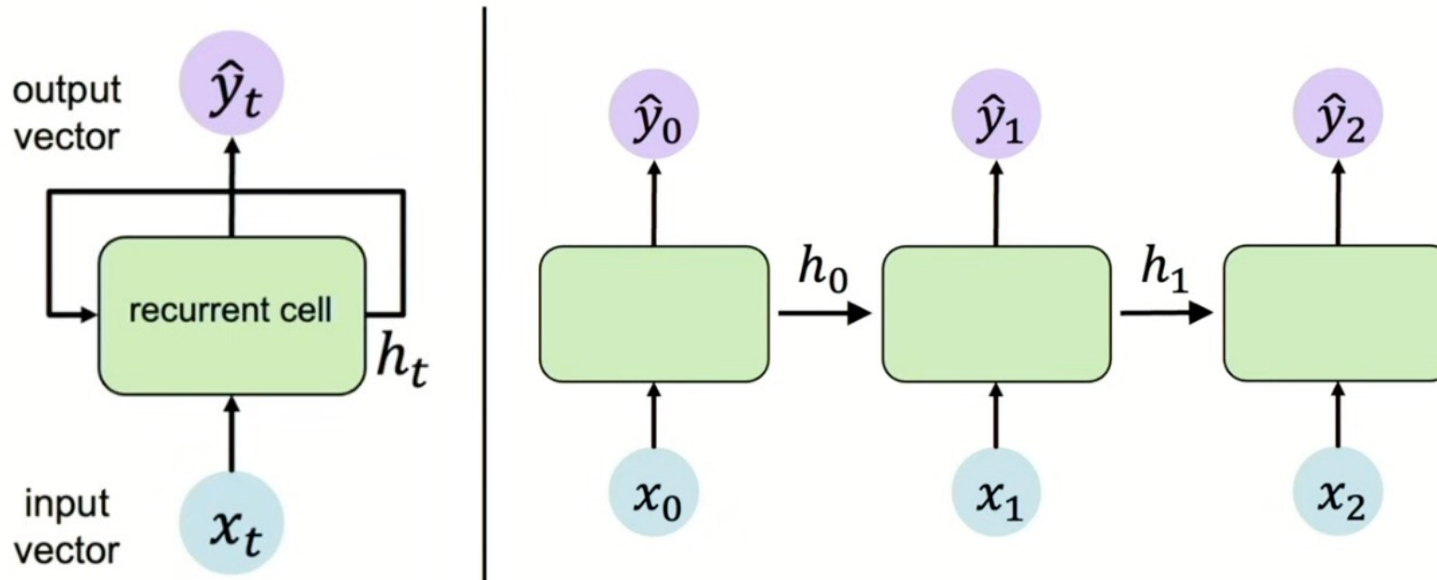
$$\hat{y}_t = f(x_t)$$

What could the issue?

There is no link between time steps.

Neurons with Recurrence

We need to build a neural network that can explicitly model the time step to time step relation: neuron with **recurrence**.



$$\hat{y}_t = f(x_t, h_{t-1})$$

output input **past memory**

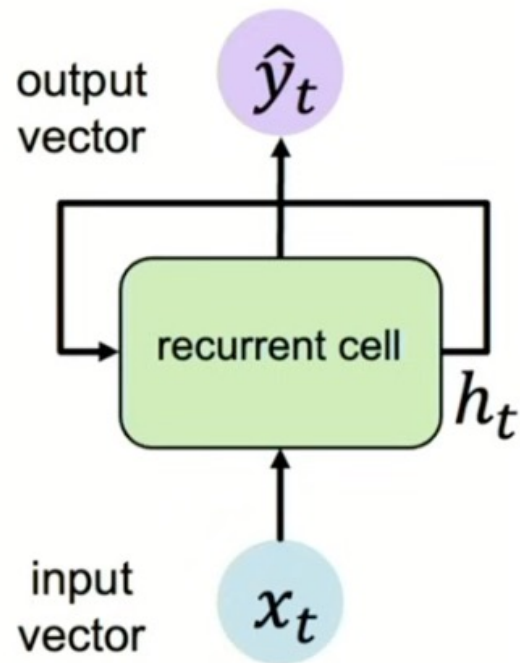
1. What is sequence modeling

2. RNN

3. LSTM

Recurrent Neural Networks (RNNs)

Add a cell state: where we apply a recurrence relation at very time step to process sequence data.



$$h_t = f_w(x_t, h_{t-1})$$

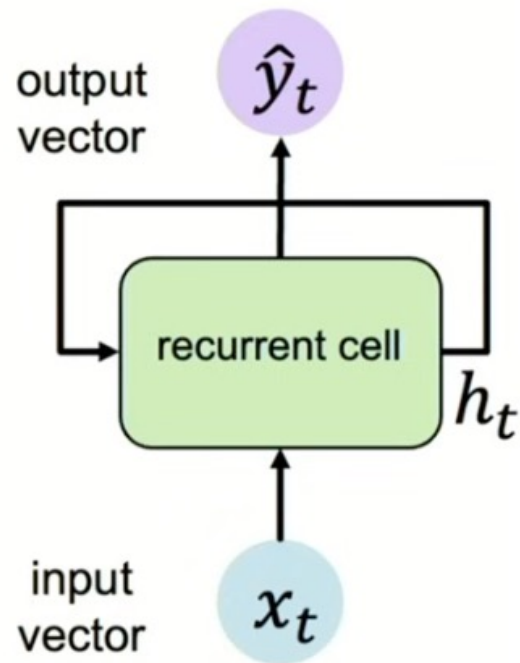
Hidden State input Past memory

Function weights

Note: in RNN, we use the **SAME** function and set of parameters at every time step.

State update and output update

Update hidden state (cell state)



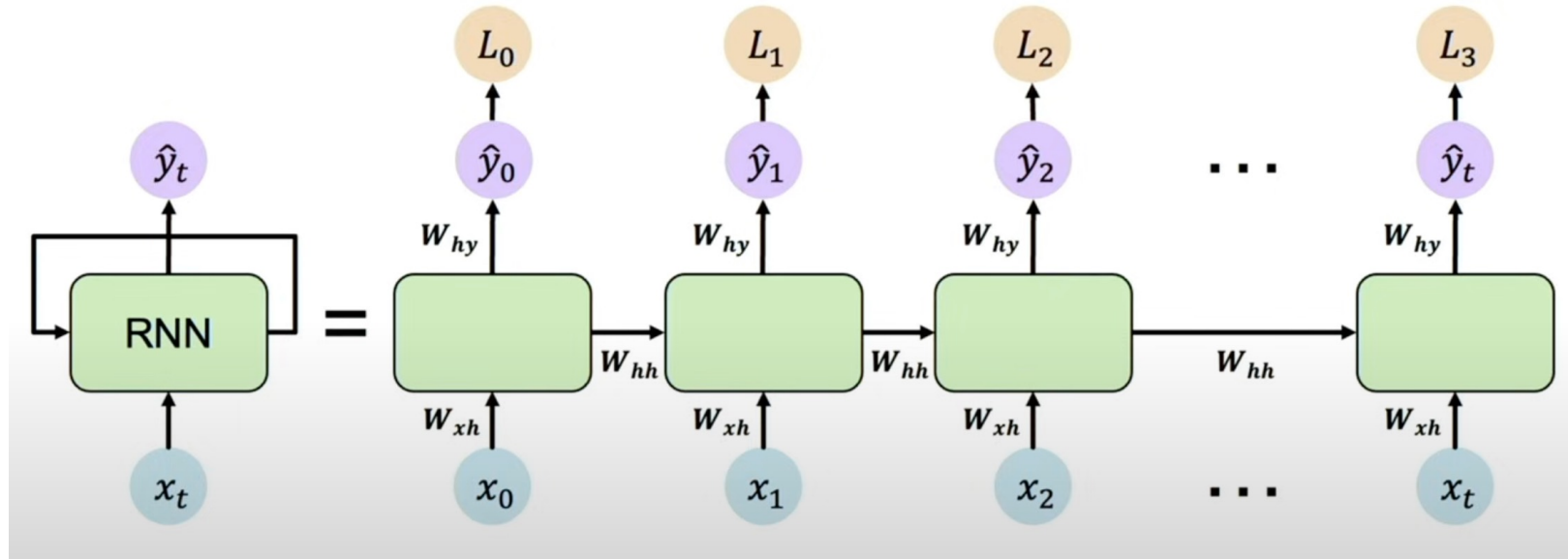
$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$

Hidden StateActivation Previous Input vector
function hidden state

$$\hat{y}_t = W_{hy}^T h_t$$

Note: in RNN, we use the **SAME** function and set of parameters at every time step.

RNNs: computational Graph



→ Forward Pass

$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$

$$\hat{y}_t = W_{hy}^T h_t$$

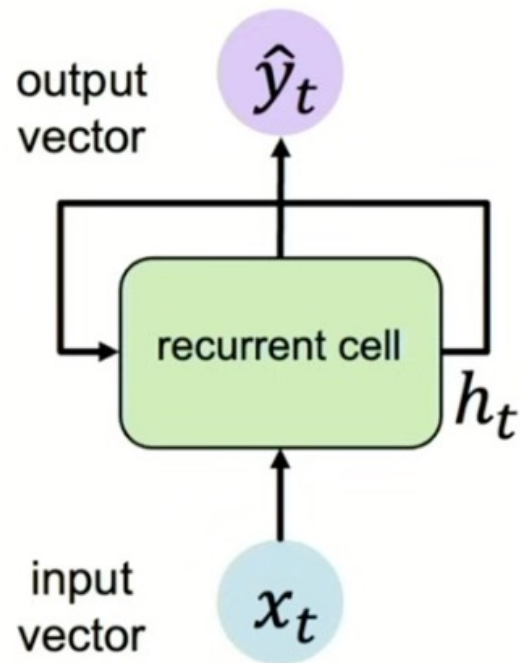


Backpropagation to update the weights

In-class practice

RNNs from Scratch

RNNs from Scratch



```
class SimpleRNNCell(nn.Module):
    def __init__(self, input_size, hidden_size, output_size):
        super().__init__()

        # Define the weights as Parameters so they're tracked by PyTorch
        self.Wxh = nn.Parameter(torch.randn(hidden_size, input_size) * 0.01)
        self.Whh = nn.Parameter(torch.randn(hidden_size, hidden_size) * 0.01)
        self.Why = nn.Parameter(torch.randn(output_size, hidden_size) * 0.01)
        self.hidden_size = hidden_size
        self.h = None

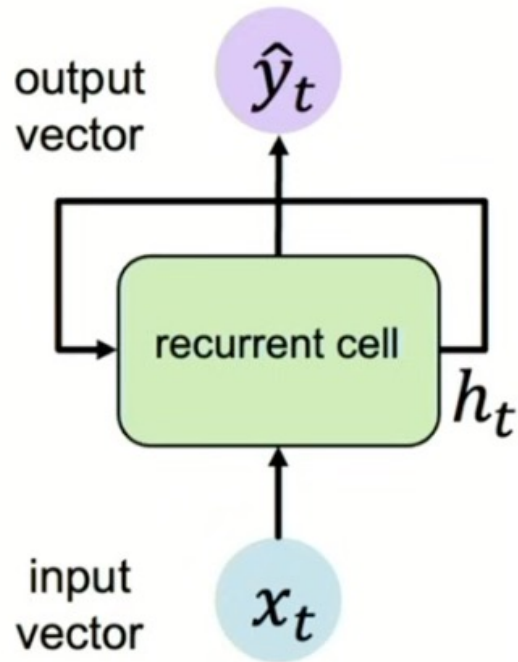
    def forward(self, x):
        if self.h is None:
            self.h = torch.zeros(x.size(0), self.hidden_size, device=x.device)

        self.h = torch.tanh(
            x @ self.Wxh.t() +      # Changed F.linear to matrix multiplication
            self.h @ self.Whh.t()  # Changed F.linear to matrix multiplication
        )
        output = F.linear(self.h, self.Why)
        return output
```

Simple one-line RNN using predefined module

```
self.rnn = nn.RNN(input_size=input_size, hidden_size=hidden_size, num_layers=1)
self.fc = nn.Linear(hidden_size, output_size)
```

RNNs: design criteria



- Handle **variable-length** sequences
- Track long-term dependencies
- Maintain information about order
- Share parameters across the sequence

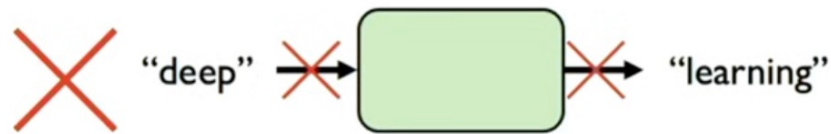
RNN meets those criteria.

Example: predict the next word

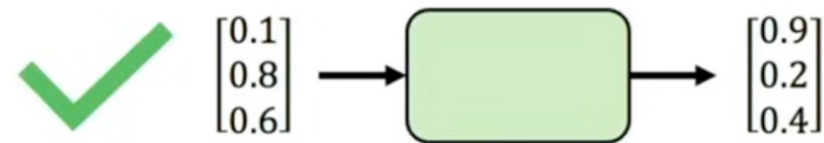
This morning, I took my cat for a __

walk

Representing the word (encoding language)

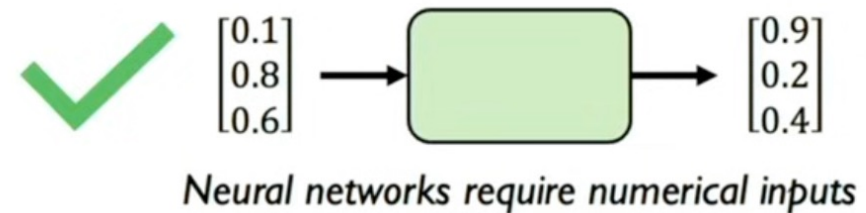
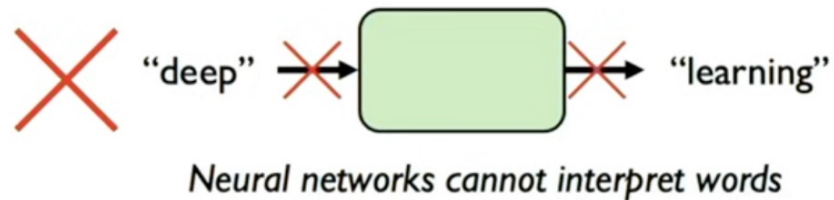


Neural networks cannot interpret words

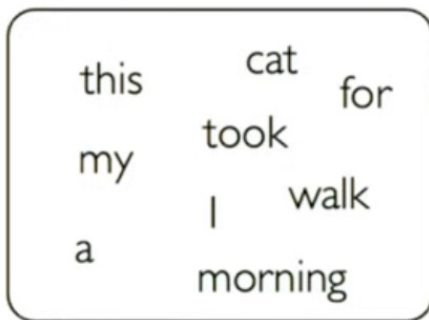


Neural networks require numerical inputs

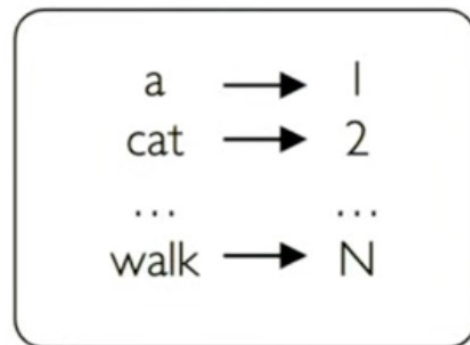
Encoding language for NN



Embedding: transform **indexes** into a vector of fixed size

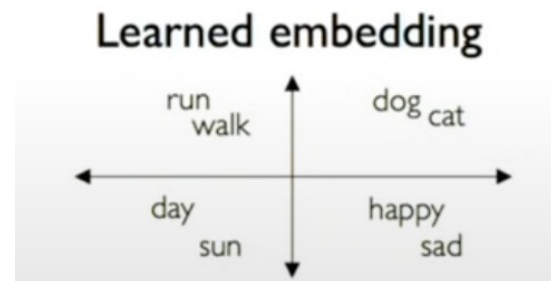


1. Vocabulary:
Corpus of words



2. Indexing:
Word to index

One-hot embedding
"cat" = $[0, 1, 0, 0, 0, 0]$
↑
i-th index



Capture differences in **order**



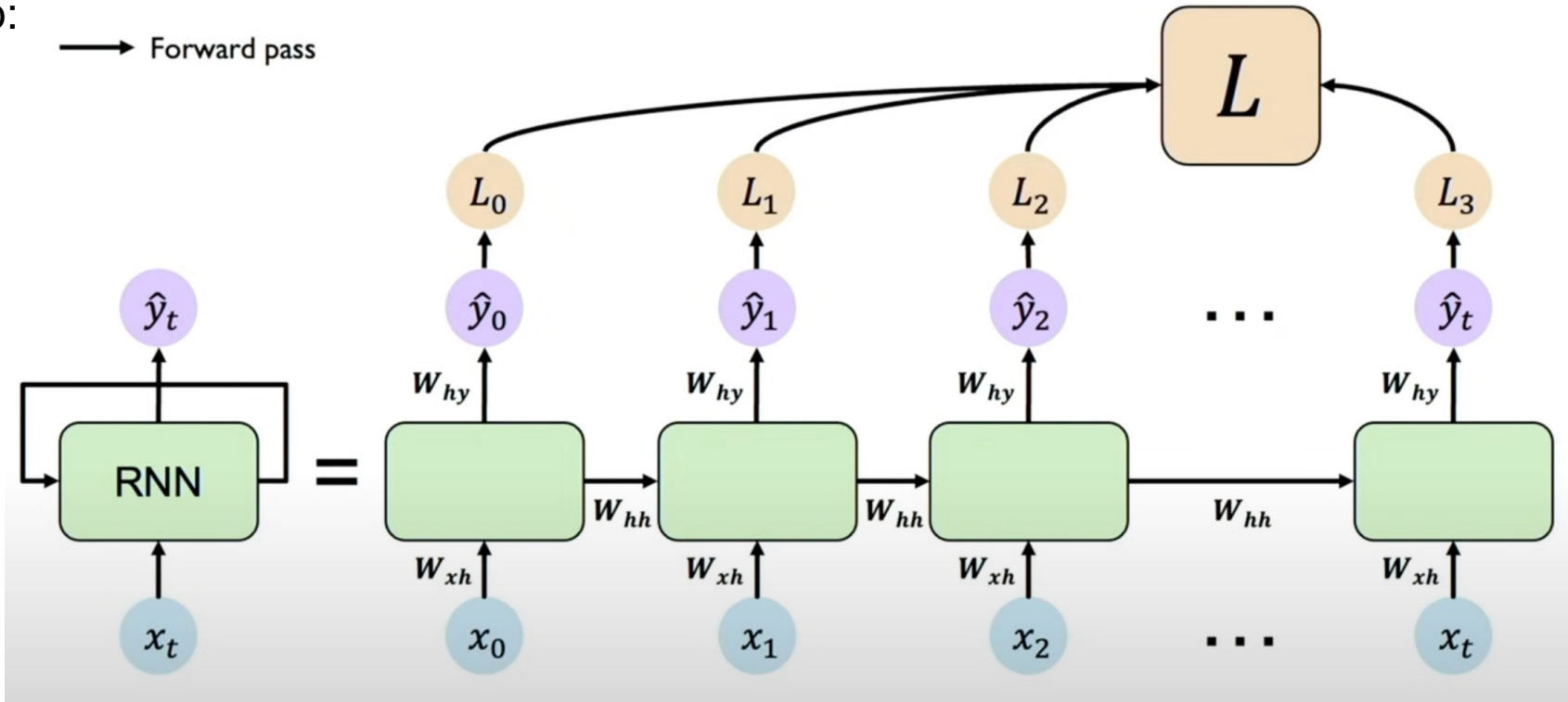
The food was good, not bad at all.

The food was bad, not good at all.



Define the loss

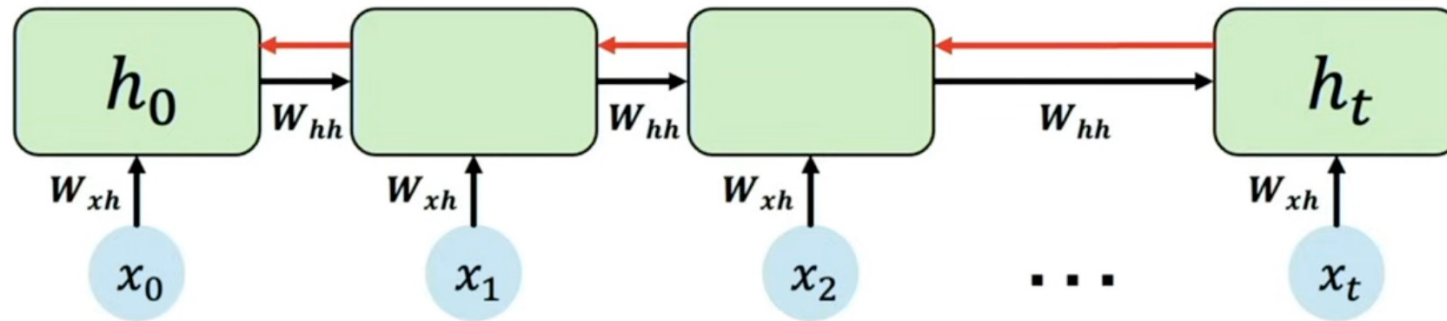
In the case of a RNN, the loss function of all time steps is defined based on the loss at every time step:



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^i; \mathbf{W}), y^i)$$

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathcal{L}(\hat{y}^{<t>}, y^{<t>})$$

Backpropagation through time



Backpropagation is done at each point in time. At timestep T , the derivative of the loss with respect to weight matrix W is expressed as:

$$\frac{\partial \mathcal{L}^{(T)}}{\partial W} = \sum_{t=1}^T \frac{\partial \mathcal{L}^{(T)}}{\partial W} \Big|_{(t)}$$

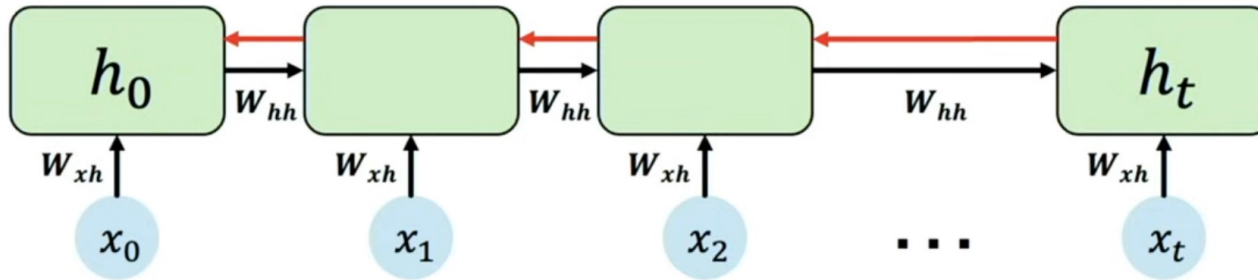
In-class practice

Predict the next word in PyTorch

- Implement the training loop

Problem with gradients

It is difficult to capture long term dependencies because of multiplicative gradient that can be exponentially decreasing/increasing with respect to the number of layers.



Exploding gradient

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Vanishing gradient

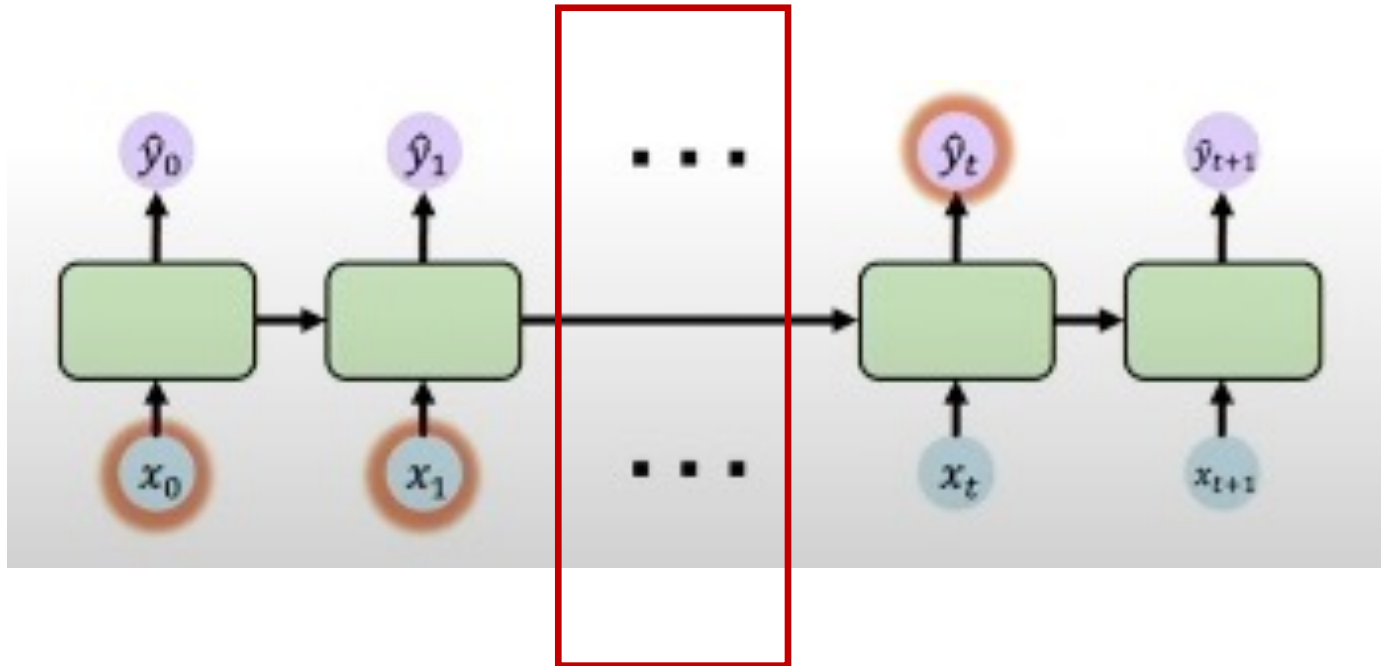
Many values < 1 :
vanishing gradients

1. Activation function
2. Weight initialization
3. Network architecture

The problem of long-term dependencies

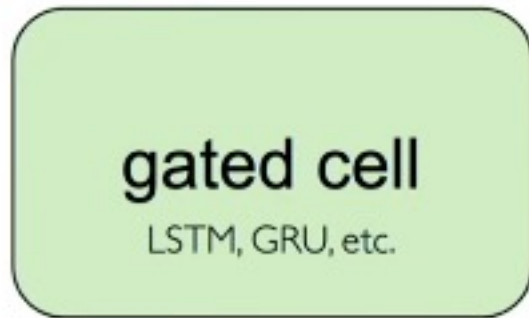
Why is vanishing gradient is a problem?

Vanishing gradient is a problem because if we multiply many small numbers together, we would have smaller and smaller gradient, and the bias parameters are only there to capture short-term dependencies.

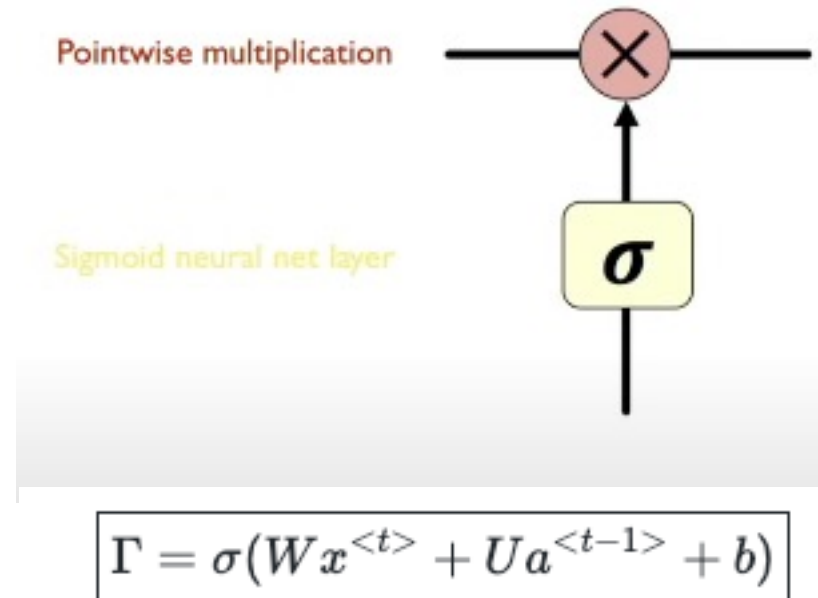


The idea of gated cells

To resolve vanishing gradient, one of the ideas is to use **gates** to selectively add or remove information within each recurrent unit with:



Gates can optionally let information through the cell



Gated Recurrent Unit (GRU) and Long Short-Term Memory units (LSTM) use gated cells

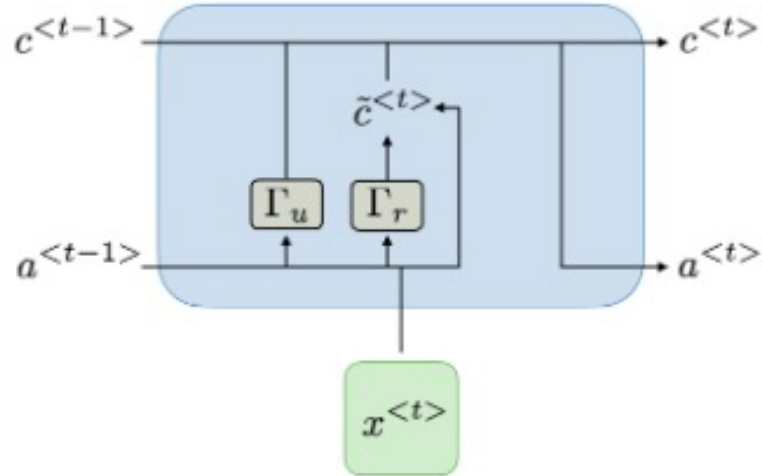
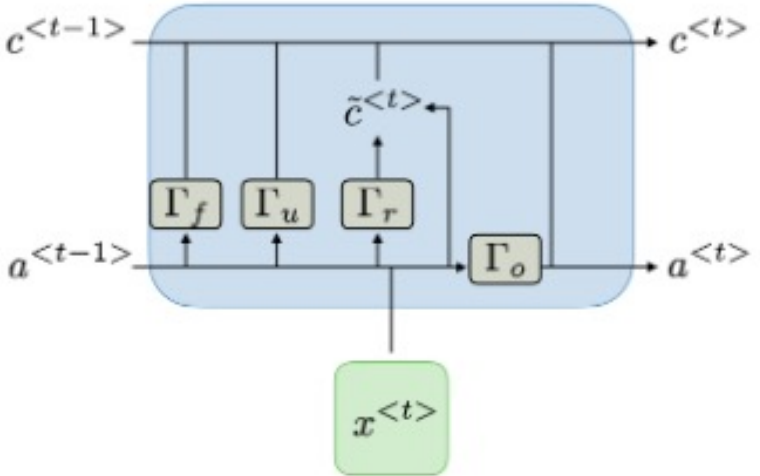
The idea of gated cells

In order to remedy the vanishing gradient problem, specific gates are used and each has a well-defined purpose.

Type of gate	Role	Used in
Update gate Γ_u	How much past should matter now?	GRU, LSTM
Relevance gate Γ_r	Drop previous information?	GRU, LSTM
Forget gate Γ_f	Erase a cell or not?	LSTM
Output gate Γ_o	How much to reveal of a cell?	LSTM

LSTM

Use gates to control the flow of information. (Forget, store, update and output ...)

Characterization	Gated Recurrent Unit (GRU)	Long Short-Term Memory (LSTM)
$\tilde{c}^{<t>}$	$\tanh(W_c[\Gamma_r \star a^{<t-1>}, x^{<t>}] + b_c)$	$\tanh(W_c[\Gamma_r \star a^{<t-1>}, x^{<t>}] + b_c)$
$c^{<t>}$	$\Gamma_u \star \tilde{c}^{<t>} + (1 - \Gamma_u) \star c^{<t-1>}$	$\Gamma_u \star \tilde{c}^{<t>} + \Gamma_f \star c^{<t-1>}$
$a^{<t>}$	$c^{<t>}$	$\Gamma_o \star c^{<t>}$
Dependencies	 <p>The diagram shows the dependencies for the GRU. It features a blue rounded rectangle representing the unit. An input $x^{<t>}$ (in a green box) is fed into the unit. The previous hidden state $c^{<t-1>}$ and previous activation $a^{<t-1>}$ are also inputs. Inside the unit, $a^{<t-1>}$ and $x^{<t>}$ are combined with weights Γ_u and Γ_r to produce $\tilde{c}^{<t>}$. Then, $\tilde{c}^{<t>}$ and $c^{<t-1>}$ are combined with weight Γ_u to produce the new hidden state $c^{<t>}$. Finally, $c^{<t>}$ is used to produce the new activation $a^{<t>}$.</p>	 <p>The diagram shows the dependencies for the LSTM. It features a blue rounded rectangle representing the unit. An input $x^{<t>}$ (in a green box) is fed into the unit. The previous hidden state $c^{<t-1>}$ and previous activation $a^{<t-1>}$ are also inputs. Inside the unit, $a^{<t-1>}$ and $x^{<t>}$ are combined with weights Γ_f, Γ_u, and Γ_r to produce $\tilde{c}^{<t>}$. Then, $\tilde{c}^{<t>}$ and $c^{<t-1>}$ are combined with weight Γ_u to produce the new hidden state $c^{<t>}$. Finally, $c^{<t>}$ is combined with weight Γ_o to produce the new activation $a^{<t>}$.</p>

Other things

HW8 out.

Final project guideline out (**start early**)

Coming next:

Deep Sequential model (RNN, LSTM, Transformer...)