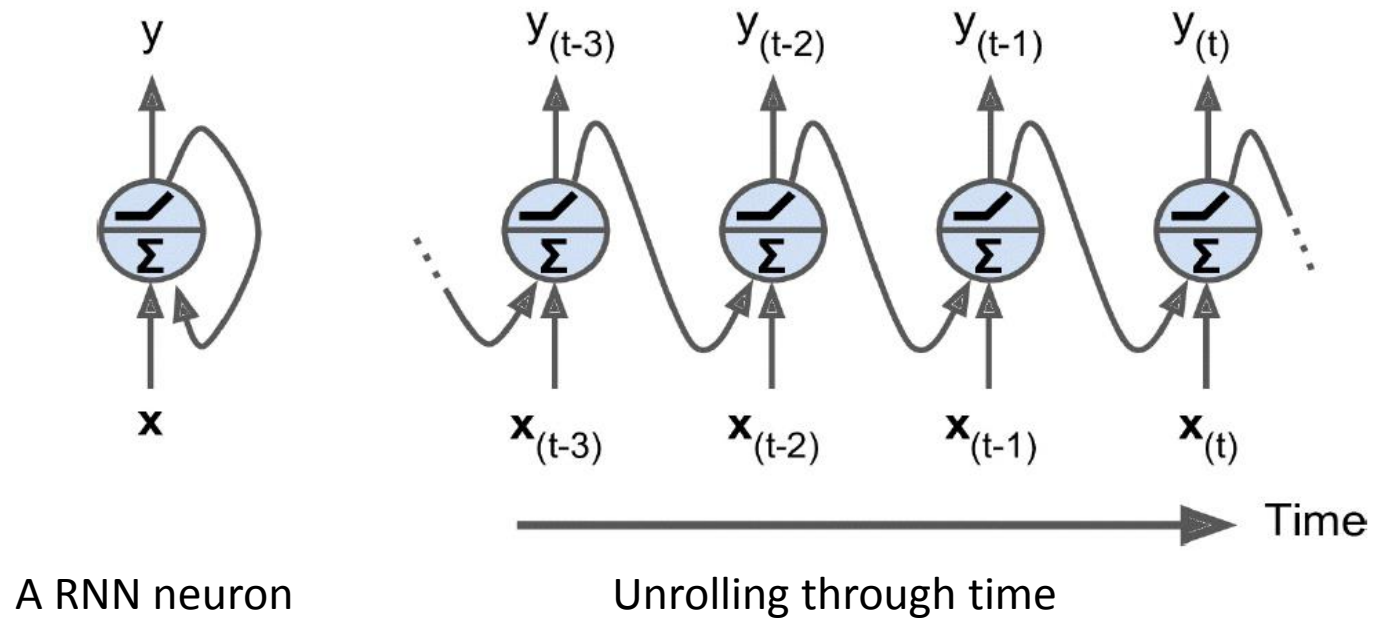


Recurrent Neural Networks

CMPUT328

Nilanjan Ray

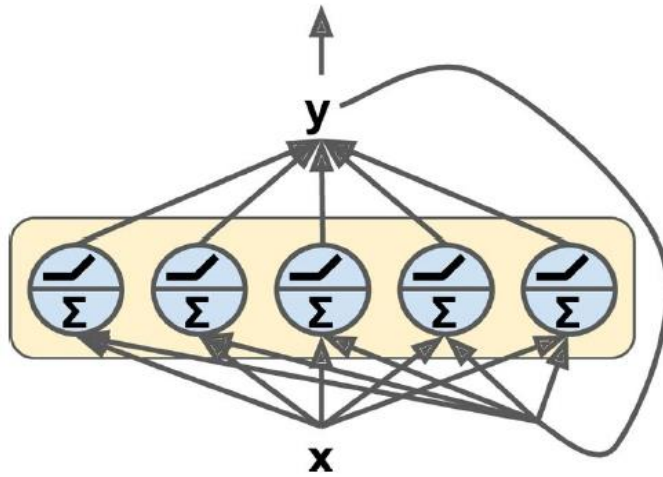
RNN: Unrolling through time



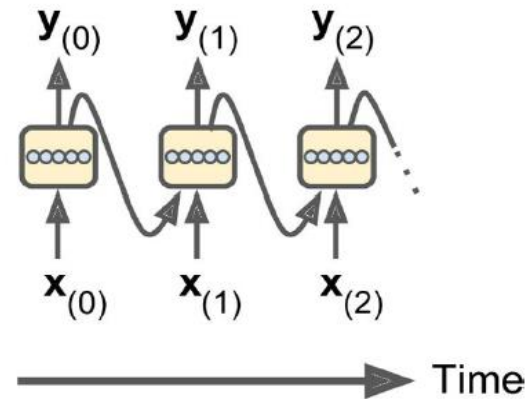
Note: The unrolled RNN is a DAG

Why is the DAG structure important?

RNN: A Layer of recurrent neurons



RNN

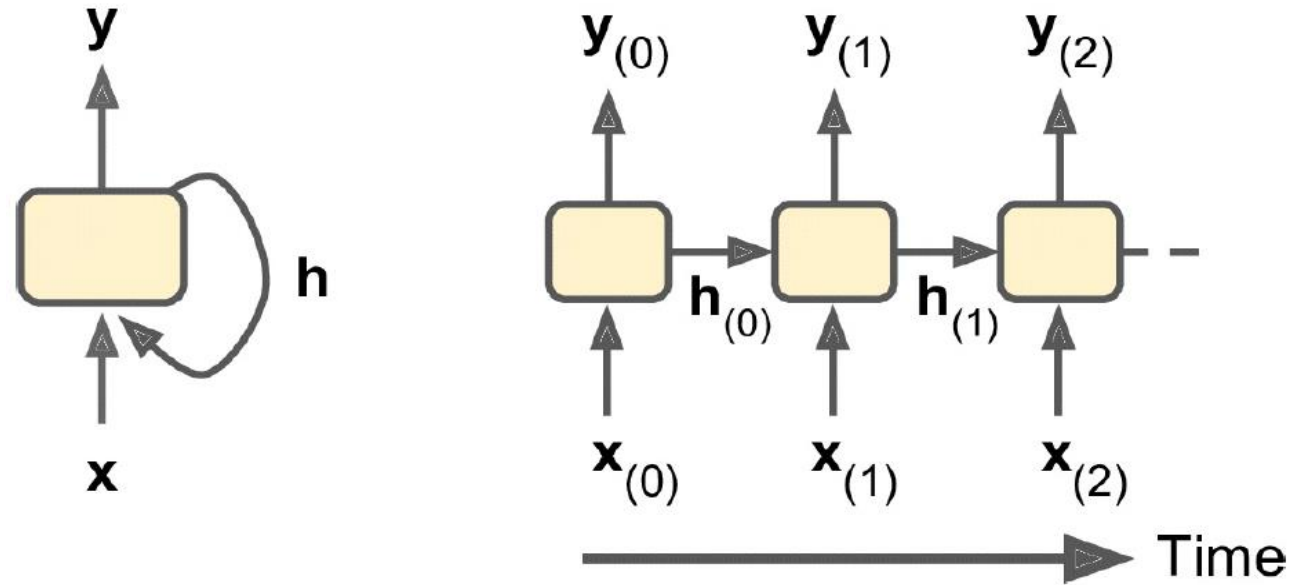


Unrolled through time

$$Y(t) = \varphi(X(t)W_x + Y(t-1)W_y + b) \quad \text{Note that parameters } W_x \text{ and } W_y \text{ are shared.}$$

φ is a non-linear activation function, such as ReLU

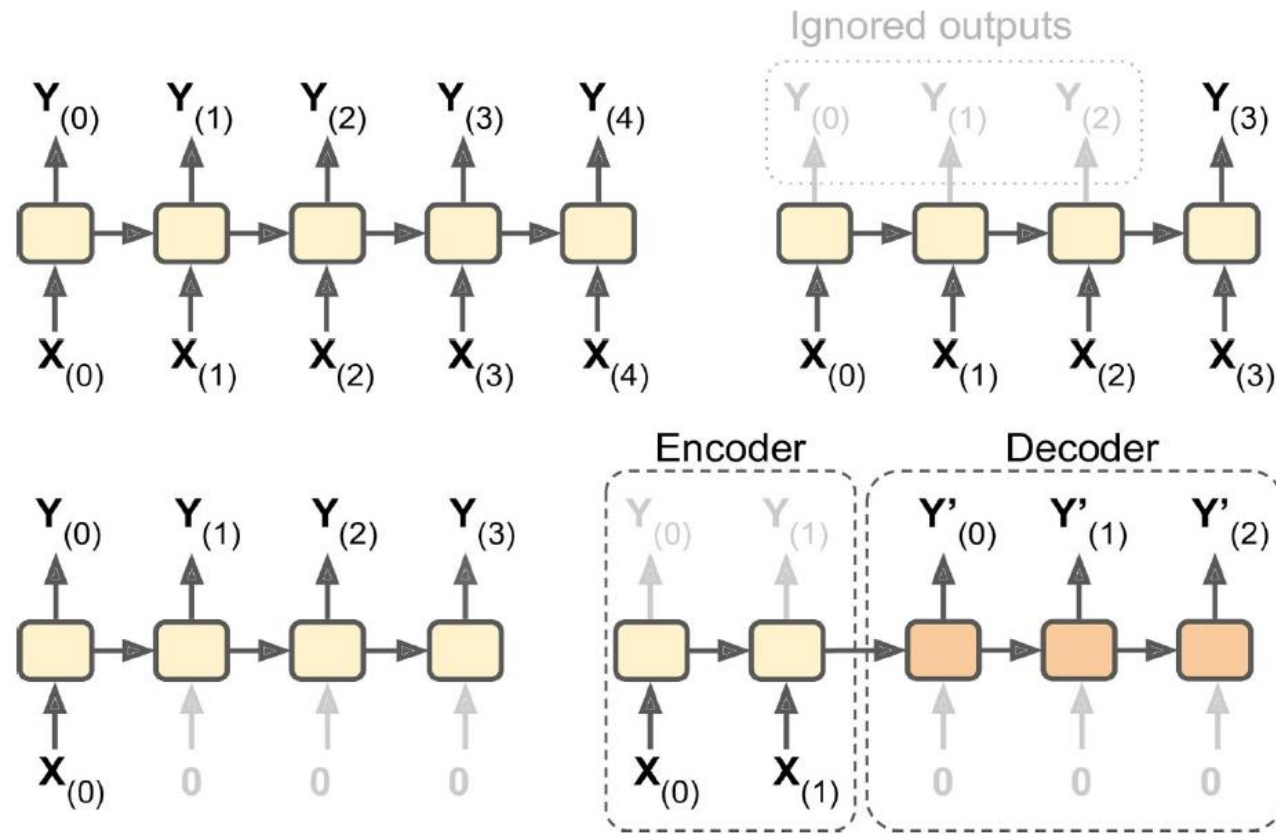
RNN: Memory cell



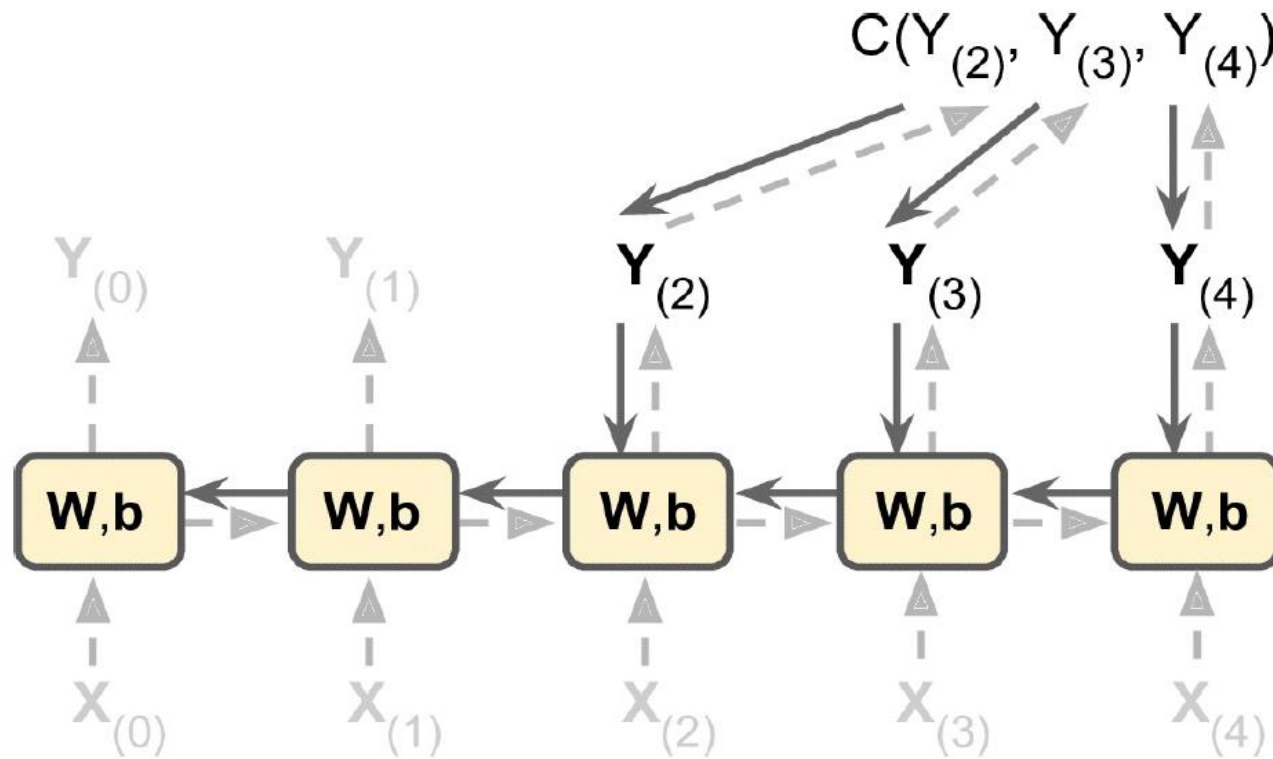
$$\mathbf{h}(t) = f(\mathbf{h}(t-1), \mathbf{x}(t))$$

$\mathbf{h}(t)$ acts as a memory cell holding “memories” from past until time point t .

Types of inputs and outputs in RNN



Backpropagation through time

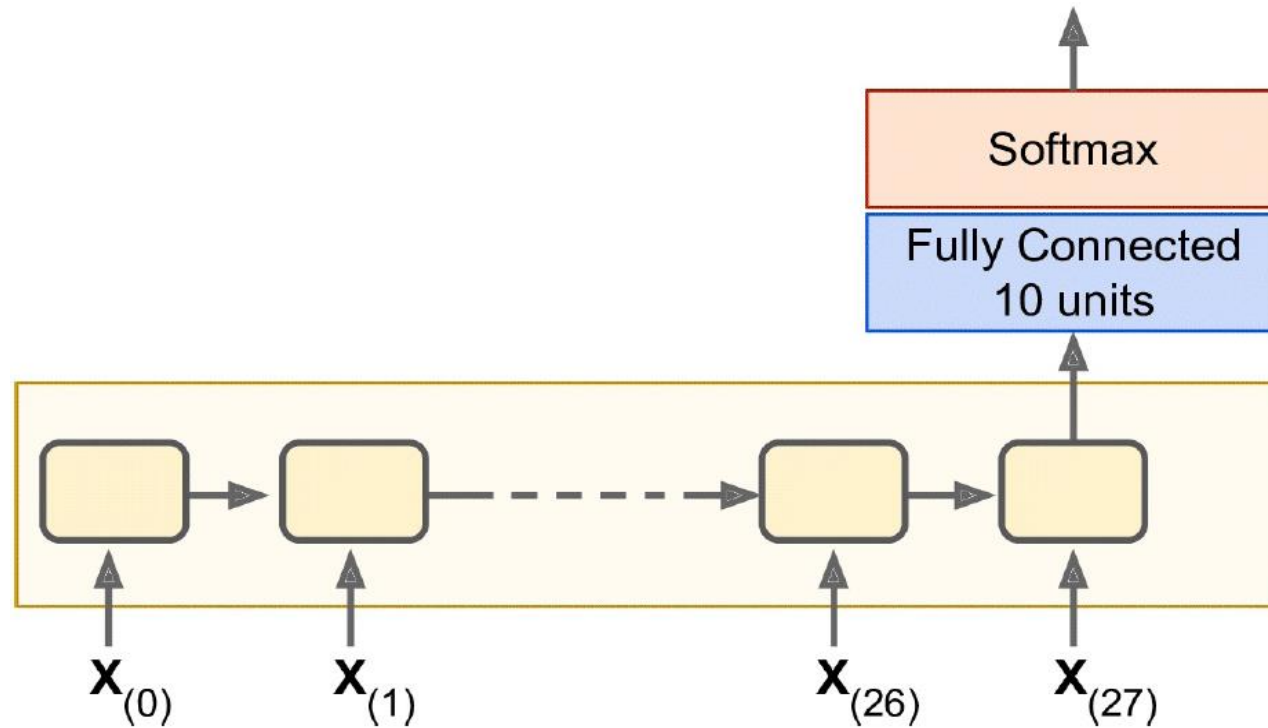


C: cost function
Dotted lines: forward pass
Solid lines: backward pass

Note: parameters are shared!

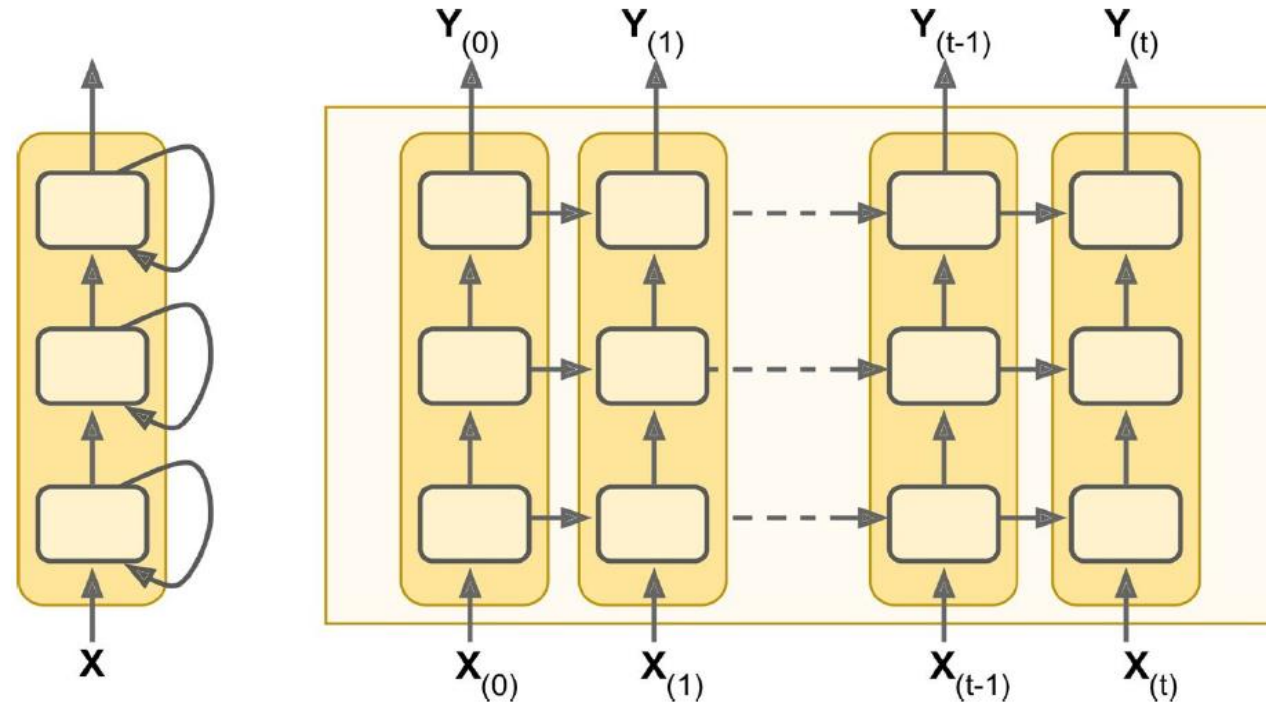
Backpropagation through time: normal backpropagation through unrolled RNN

MNIST classification: Using sequence!



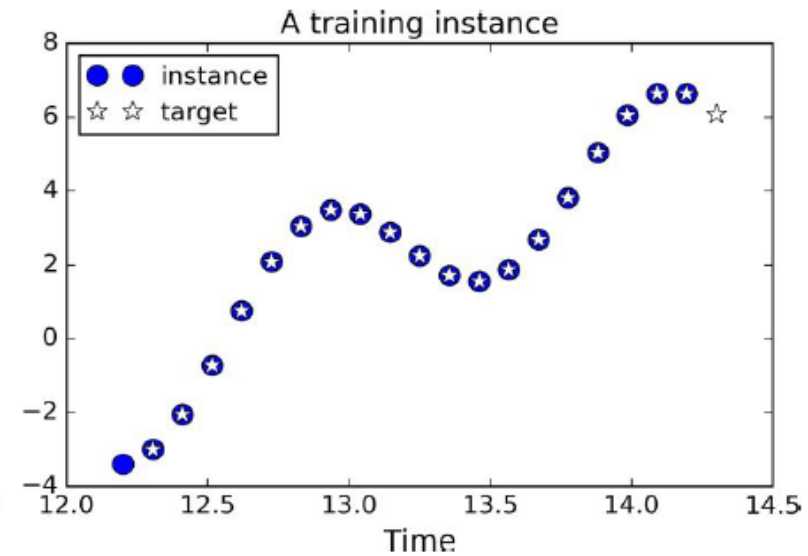
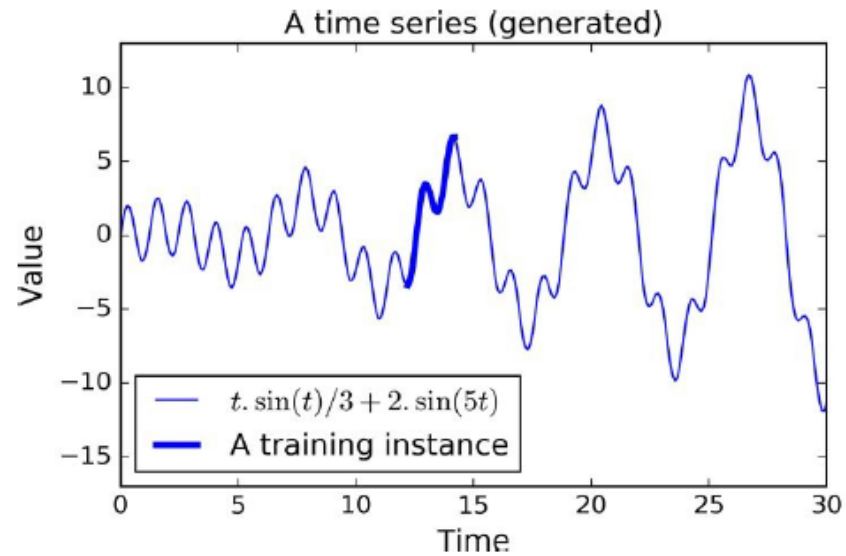
Treat each as a sequence of 28 rows

Deep (Multi-layer) RNN

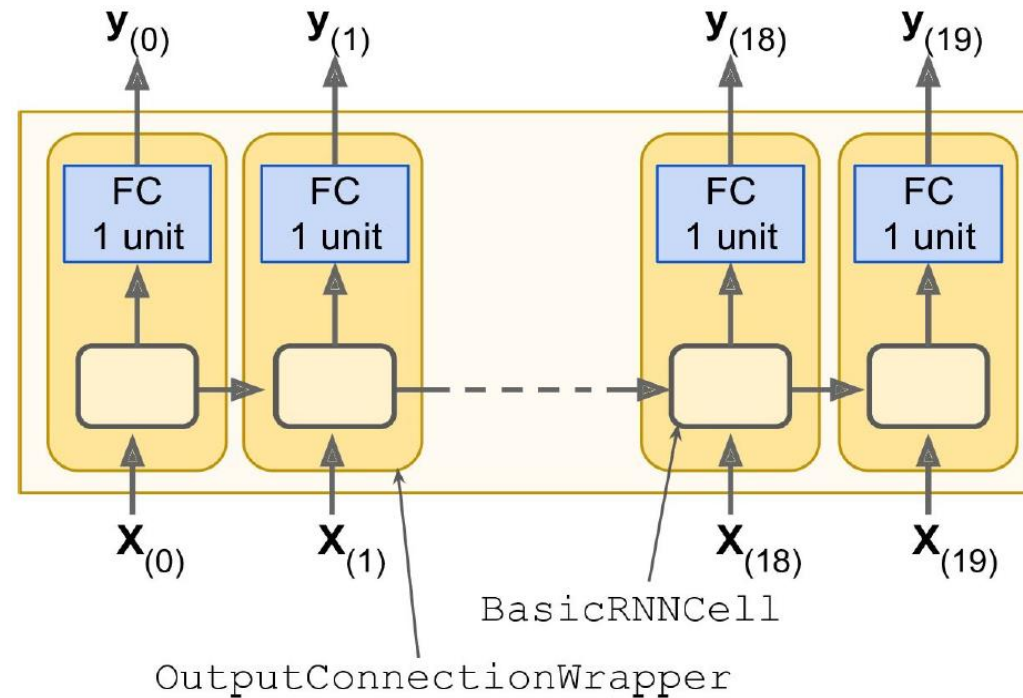


Example implementation for MNIST classification

Predicting a time series

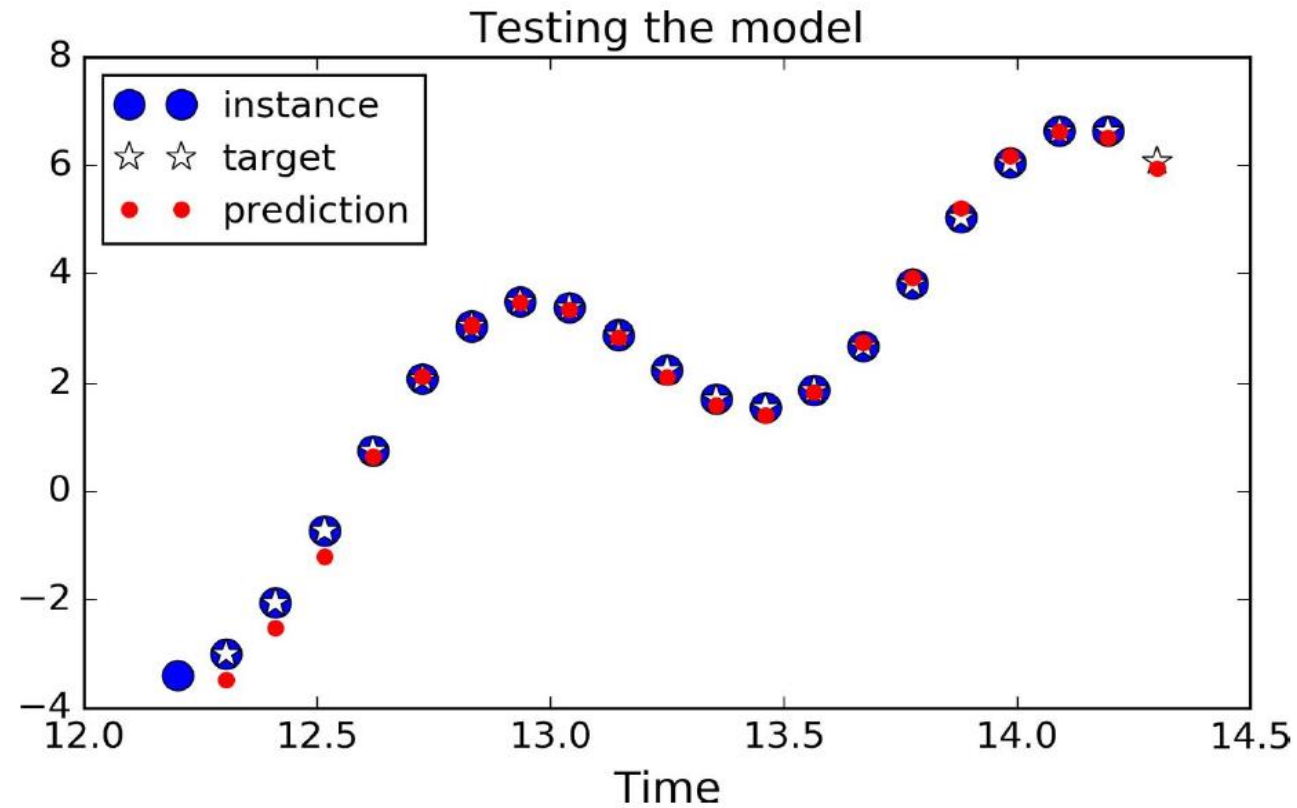


Predicting a time series...

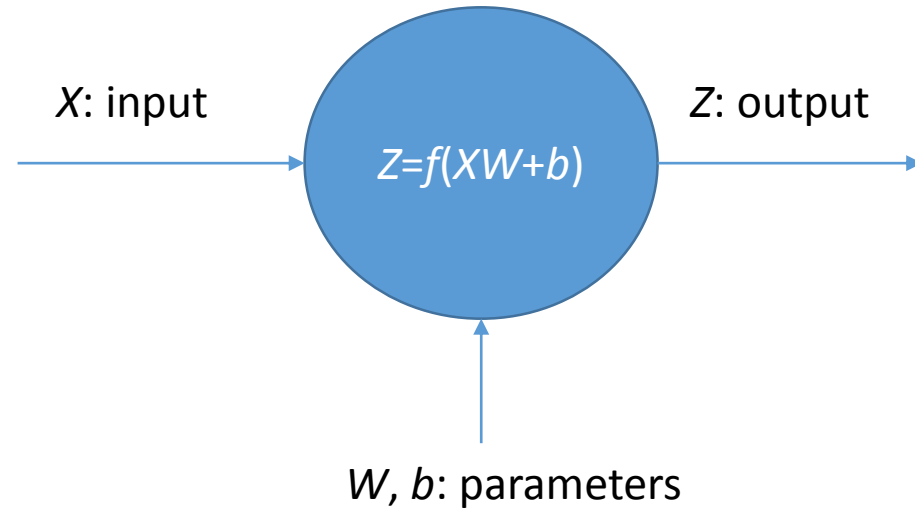


Wrapper function in TensorFlow for dimensionality reduction

Predicting a time series...



Backpropagation in RNN



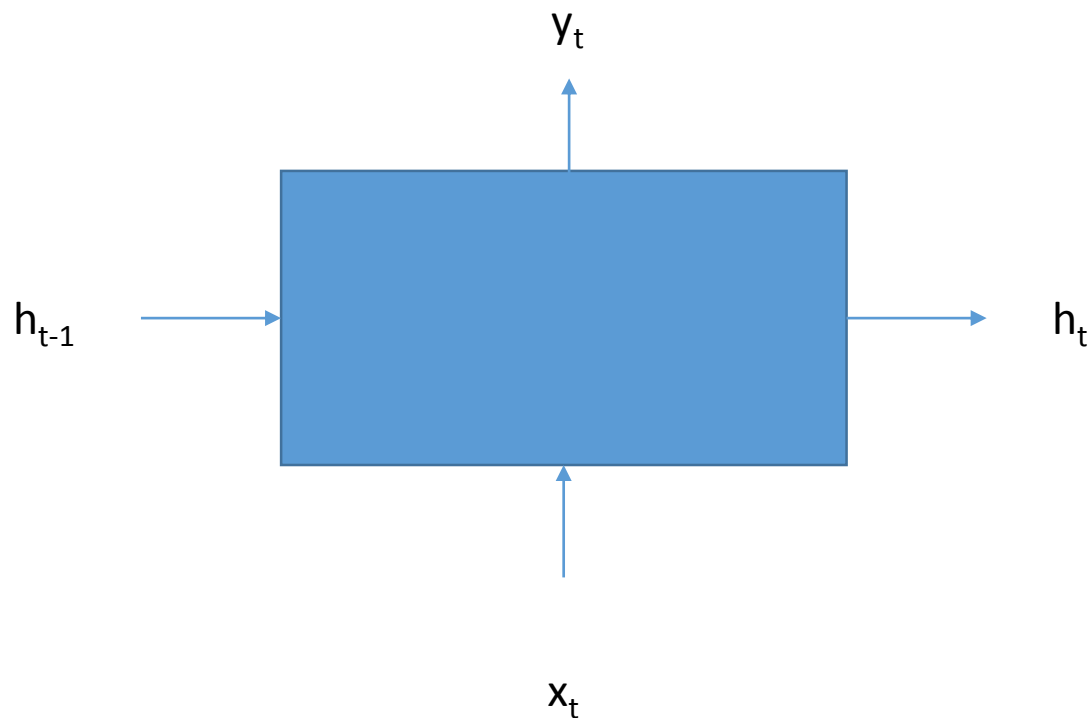
Let's recall rule of BP
for a computation node

$$\delta X = [f'(XW + b) \cdot \delta Z] * W^T$$

$$\delta W = X^T * [f'(XW + b) \cdot \delta Z]$$

$$\delta b = \sum_k [f'(XW + b) \cdot \delta Z]_{k,:}$$

Backpropagation in RNN...



How do we apply BP here?

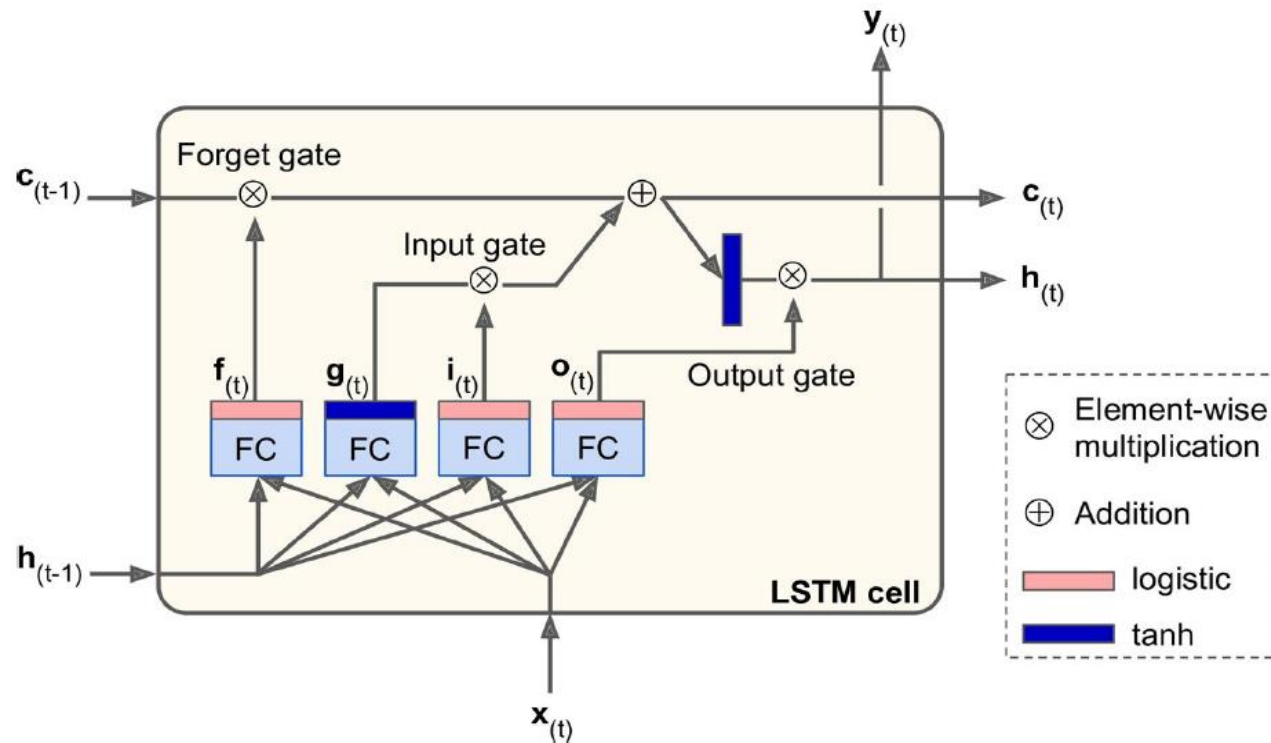
$$h_t = f(x_t W_x^h + h_{t-1} W_h^h + b^h)$$

$$y_t = g(x_t W_x^y + h_{t-1} W_h^y + b^y)$$

Vanishing / Exploding gradient problem

- Use BP formula to understand this issue

Long Short-Term Memory (LSTM)



LSTM computations:

$$i_{(t)} = \sigma(W_{xi}^T \cdot x_{(t)} + W_{hi}^T \cdot h_{(t-1)} + b_i)$$

$$f_{(t)} = \sigma(W_{xf}^T \cdot x_{(t)} + W_{hf}^T \cdot h_{(t-1)} + b_f)$$

$$o_{(t)} = \sigma(W_{xo}^T \cdot x_{(t)} + W_{ho}^T \cdot h_{(t-1)} + b_o)$$

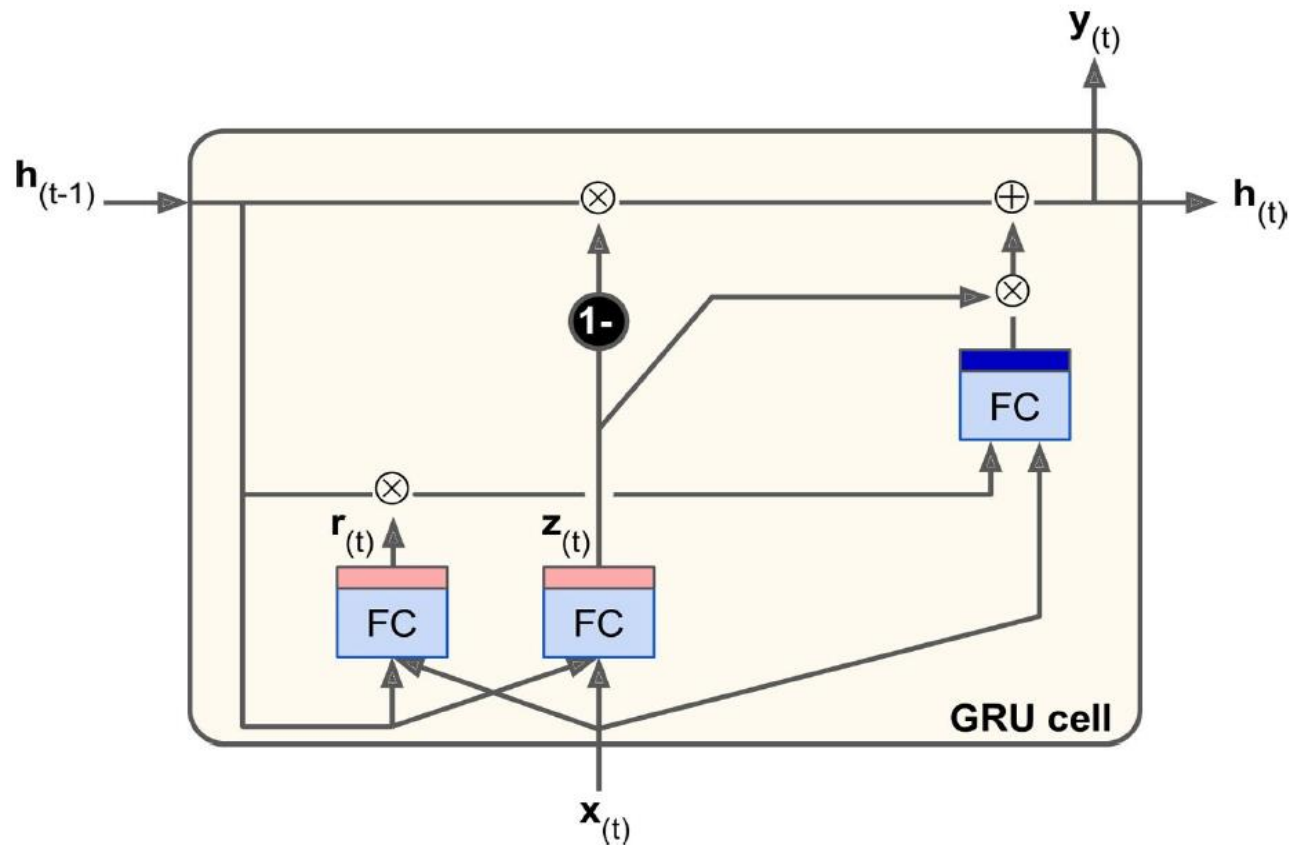
$$g_{(t)} = \tanh(W_{xg}^T \cdot x_{(t)} + W_{hg}^T \cdot h_{(t-1)} + b_g)$$

$$c_{(t)} = f_{(t)} \otimes c_{(t-1)} + i_{(t)} \otimes g_{(t)}$$

$$y_{(t)} = h_{(t)} = o_{(t)} \otimes \tanh(c_{(t)})$$

This architecture helps to mitigate vanishing/exploding gradient problem. Why?
Let's use LSTM in our time series prediction.

Gated Recurrent Unit (GRU)



GRU computations;

$$\begin{aligned} \mathbf{z}_{(t)} &= \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)}) \\ \mathbf{r}_{(t)} &= \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)}) \\ \mathbf{g}_{(t)} &= \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)})) \\ \mathbf{h}_{(t)} &= (1 - \mathbf{z}_{(t)}) \otimes \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{h}_{(t-1)}) + \mathbf{z}_{(t)} \otimes \mathbf{g}_{(t)} \end{aligned}$$

GRU is a much more simplified recurrent unit; but it is almost as good as LSTM.