

# Recurrent Neural Networks

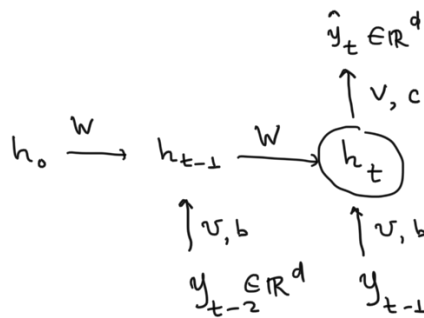
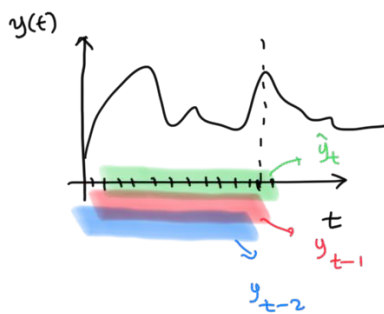
Setup: We have sequence data-set  $\{y_t : t=1, 2, \dots, T\}$ ,  $y_t \in \mathbb{R}^d$

Goal: Predict the next sequence  $\hat{y}_t$  as a function of previous values (lags), i.e.  $\hat{y}_t = f_\theta(y_{t-1}, y_{t-2}, \dots, y_{t-L})$

$L$ : # of lags

$f_\theta$ : recurrent neural net.

Example: Prediction using 2 lags, i.e.  $\hat{y}_t = f_\theta(y_{t-1}, y_{t-2})$



Params:

$$\theta := \{U, b, W, V, c\}$$

⊕ Notice how  $U, b, W$  are shared.

$$\begin{cases} \hat{y}_t = h_t V + c \\ h_t = \tanh(h_{t-1} W + y_{t-1} U + b) \\ h_{t-1} = \tanh(h_0 W + y_{t-2} U + b) \\ h_0 = 0 \end{cases}$$

User defined hyper- $\theta$

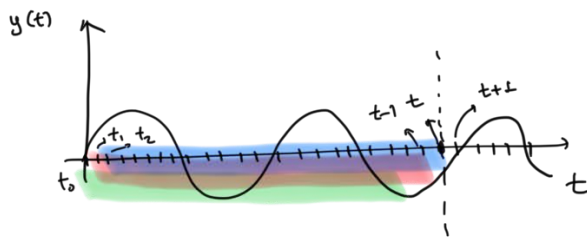
- 1.) # of lags
- 2.) dimension of the hidden state  $q$

Training: Minimize an MSE loss:

$$L(\theta) := \frac{1}{T-2} \sum_{t=3}^T (y_t - \hat{y}_t)^2$$

$y_t \in \mathbb{R}^d$

Example: Predict the dynamics of a sine wave, i.e.  $y(t) = \sin(t)$



$$\hat{y}_t = f_\theta(y_{t-1}, y_{t-2})$$

$y_t \in \mathbb{R}^{N \times D}$  (here  $D=1$  since we have a single state/time series)

$$y_t := \underbrace{[y(2), y(3), \dots, y(t)]}_N \in \mathbb{R}^N$$

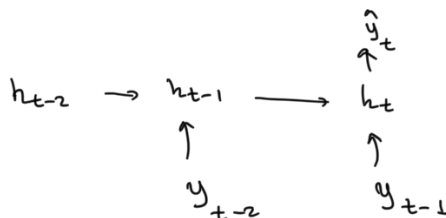
Input:  $(y_{t-1}, y_{t-2})$

$$y_{t-1} := \underbrace{[y(1), y(2), \dots, y(t-1)]}_N \in \mathbb{R}^N$$

$X: L \times N \times D$   
 # of lags, # of measurements, dim of state

$$y_{t-2} := [y(0), y(1), \dots, y(t-2)]$$

$Y = y_t, N \times D$



$$\begin{cases} h_{t-2} = 0 \\ h_{t-1} = \tanh(h_{t-2}W + X[0, :, :]V + c) \\ h_t = \tanh(h_{t-1}W + X[1, :, :]V + c) \\ \hat{y}_t = h_t V + c \end{cases}$$

## Long Short-Term Memory network:

In an LSTM one replaces the hidden units:

Standard RNN,  $h_t = \tanh(h_{t-1}W + y_{t-1}V + b)$   $\theta := \{W, V, b\}$   
 Cell update rule

Replace it with:

$$\begin{cases} h_t = o_t \odot \tanh(s_t) & : \text{output vector} \\ o_t := \sigma(h_{t-1}W_o + y_{t-1}V_o + b_o) & : \text{output gate} \\ s_t := f_t \odot s_{t-1} + i_t \odot \tilde{s}_t & : \text{cell state} \end{cases}$$

↑  
element-wise multiplication

Input:  $h_{t-1}, y_{t-1}$

Output:  $h_t$

$$\left[ \begin{array}{l} \tilde{s}_t := \tanh(h_{t-1} W_s + y_{t-1} V_s + b_s) \\ i_t := \sigma(h_{t-1} W_i + y_{t-1} V_i + b_i) : \text{external inp gate} \\ f_t := \sigma(h_{t-1} W_f + y_{t-1} V_f + b_f) : \text{forget gate} \end{array} \right.$$

Now we have more parameters:

$$\Theta := \{ W_o, V_o, b_o, W_s, V_s, b_s, W_i, V_i, b_i, W_f, V_f, b_f, \underbrace{V, c}_{\text{final out layer param}} \}$$