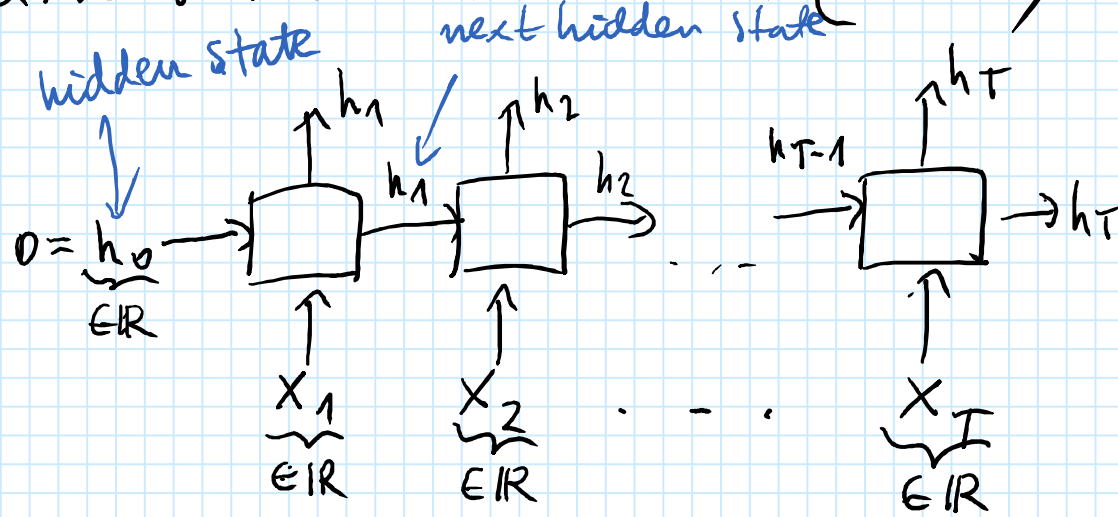


Machine Learning 9.12.2022

Recurrent Neural Network (RNN)



$$h_t = \sigma(Wx_t + Uh_{t-1} + b), \quad \sigma - \text{sigmoid}$$

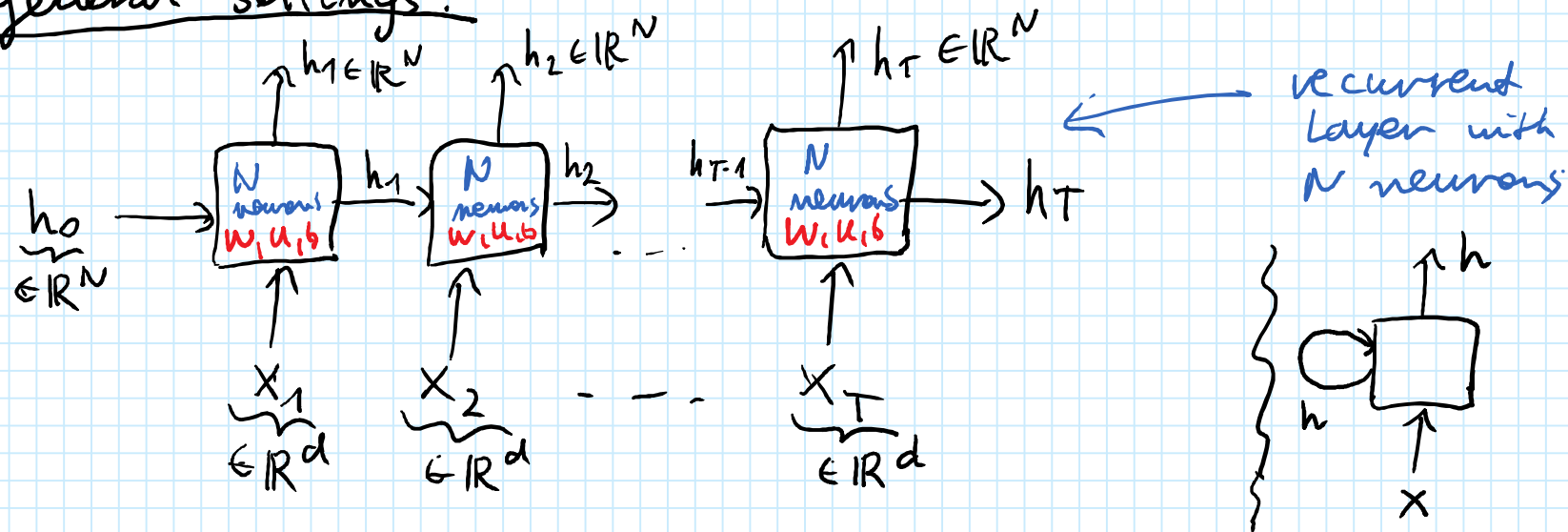
Applications

- Modelling Time Series
- Natural Language Processing (NLP)

Examples

- identifying sarcastic comments (1s) on Reddit
- identifying number of stars in movie review
- predicting closing price of stocks

More general settings:



$$\mathbb{R}^N \ni h_t = \sigma \left(\underbrace{W}_{\in \mathbb{R}^{N \times d}} \underbrace{x_t}_{\in \mathbb{R}^d} + \underbrace{U}_{\in \mathbb{R}^{N \times N}} \underbrace{h_{t-1}}_{\in \mathbb{R}^N} + \underbrace{b}_{\in \mathbb{R}^N} \right)$$

- σ (sigmoid) is acting on every coordinate
- W, U, b are the same at every time-step
- Memory is shared between all neurons

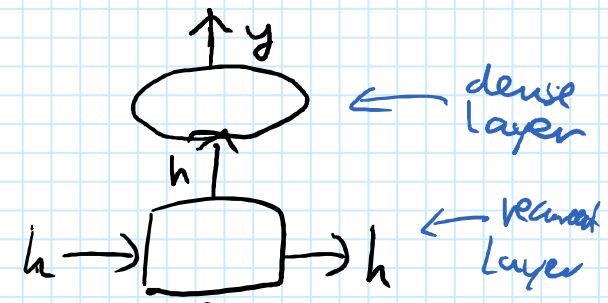
$$\text{Number of weights} = Nd + N^2 + N$$

- Final output is h_T or (h_1, h_2, \dots, h_T)

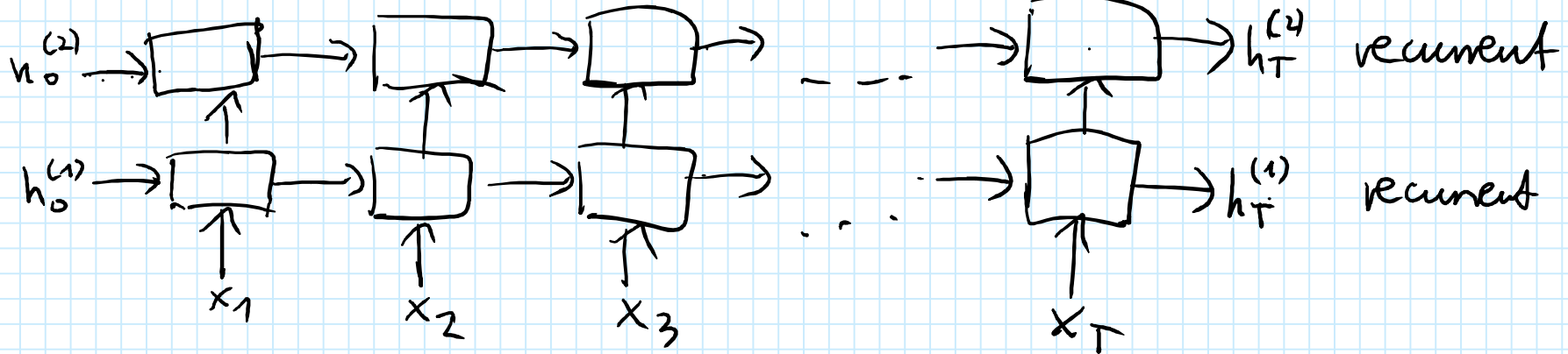
Examples

1) Elman Network

$$h_t = \sigma (w_h x_t + u_h h_{t-1} + b_h)$$
$$y_t = \sigma_y (w_h h_t + b_y)$$



2)



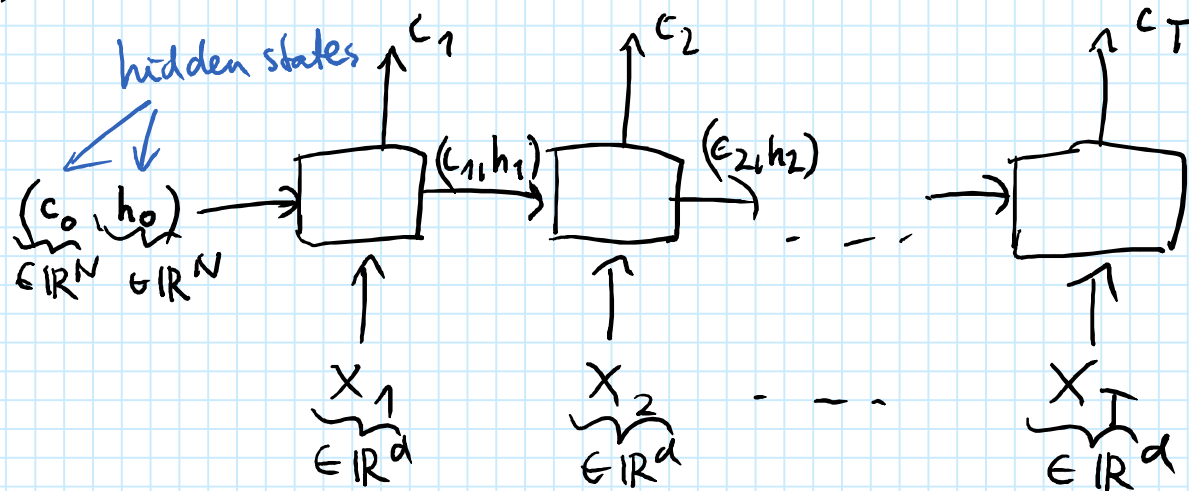
PROBLEM: Vanishing Gradients!

If T is large, then distance from x_1 to x_T is very long and the network is very deep in fact!

SOLUTION LSTM!

Long Short-Term Memory (LSTM)

Reference: Hochreiter, Schmidhuber 1997



$$f_t = \sigma_g \left(\underbrace{W_f}_{N \times d} x_t + \underbrace{U_f}_{N \times N} h_{t-1} + b_f \right)$$
$$i_t = \sigma_g \left(\underbrace{W_i}_{N \times d} x_t + \underbrace{U_i}_{N \times N} h_{t-1} + b_i \right)$$
$$o_t = \sigma_g \left(\underbrace{W_o}_{N \times d} x_t + \underbrace{U_o}_{N \times N} h_{t-1} + b_o \right)$$
$$\tilde{c}_t = \sigma_c \left(\underbrace{W_c}_{N \times d} x_t + \underbrace{U_c}_{N \times N} h_{t-1} + b_c \right)$$

σ_g - sigmoid, σ_c - tanh

f_t - forget gate

i_t - input gate

o_t - output gate

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \sigma_c(c_t)$$

Number of weights = $4(Nd + N^2 + N)$

GRU (Gated Recurrent Unit)

- less parameters than in LSTM
- very similar results

$$\text{Number of weight} = 3(Nd + N^2 + N)$$

$$z_t = \sigma_g (W_z x_t + U_z h_{t-1} + b_z)$$

update gate

$$r_t = \sigma_g (W_r x_t + U_r h_{t-1} + b_r)$$

reset gate

$$\hat{h}_t = \sigma_h (W_h x_t + U_h h_{t-1} + b_h)$$

candidate activation

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t$$

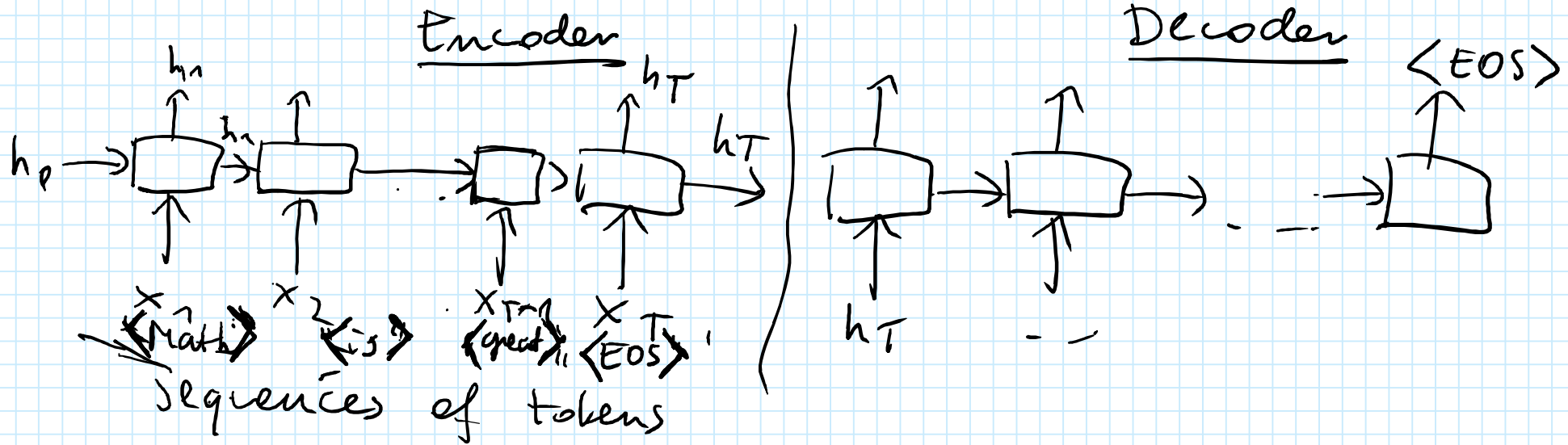
σ_g - sigmoid

σ_h - tanh

$z_t \approx 1$ on some coordinate, then $h_t \approx \hat{h}_t$, so we choose new candidate (\hat{h}) as h

$z_t \approx 0$ on some coordinate, then $h_t \approx h_{t-1}$
so we choose "old" h

Encoder-Decoder



[
 < >
 < is >
 < great >
 < Math >
 < UNK >
 < EOS >
]

<UNK> - unknown word
 <EOS> - end of sentence

Keras

- RepeatVector
- return_sequences = TRUE