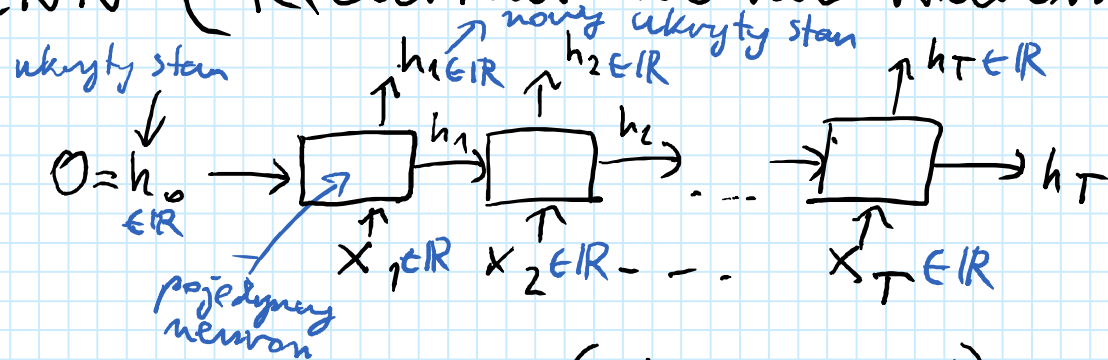


Uczenie maszynowe 7.12.2022

RNN (Recurrent Neural Network)



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

W, U, b - wagi

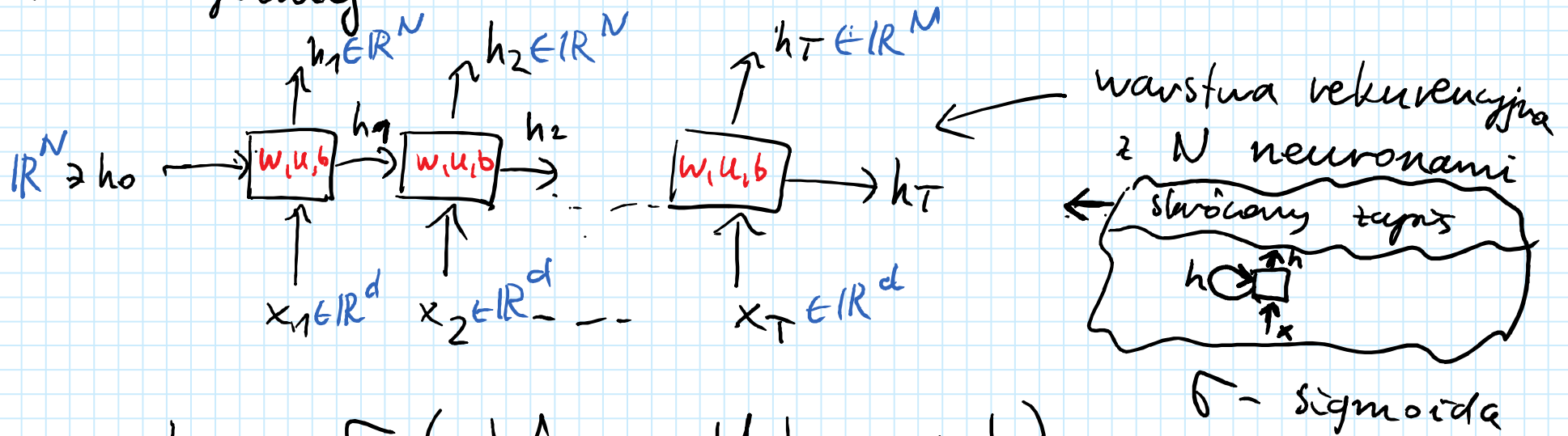
Zastosowania

- modelowanie szeregów czasowych
- przetwarzanie języka naturalnego

Przykłady

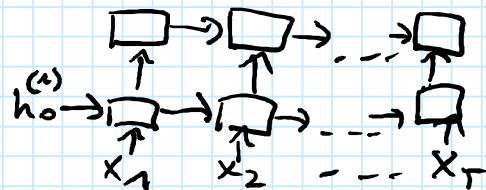
- prognozowanie cen zamknięcia akcji na podstawie danych historycznych
- klasyfikacja recenzji ("gwiazdki")
- wykrywanie sarkazmu (1s na Reddicie) itd.

RNN - ogólniej



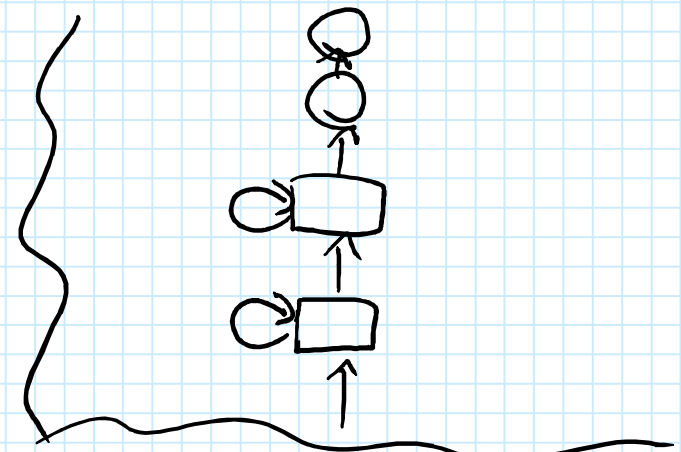
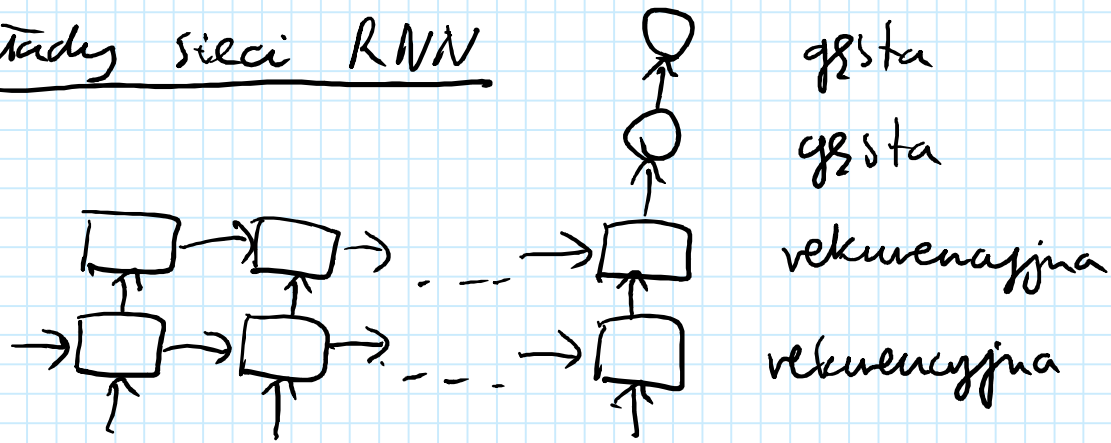
$$h_t \in \mathbb{R}^N = \sigma \left(\underbrace{W}_{\mathbb{R}^{N \times d}} x_t \in \mathbb{R}^d + \underbrace{U}_{\mathbb{R}^{N \times N}} h_{t-1} \in \mathbb{R}^N + \underbrace{b}_{\mathbb{R}^N} \right)$$

- σ nakładamy po współrzędnych
- W, U, b są takie same w każdym kroku czasowym
- pamięć jest współdzielona pomiędzy wszystkie neurony
- liczba wag = $Nd + N^2 + N$
- ostateczne wyjście z sieci to h_T albo (h_1, h_2, \dots, h_T)



Przykłady sieci RNN

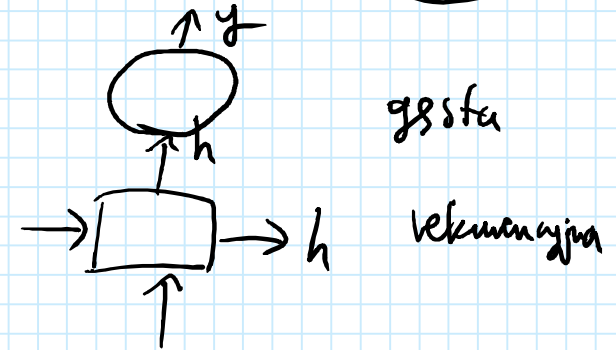
1)



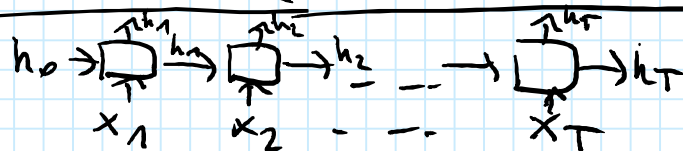
2) Elman Network

$$h_t = \sigma (W_h \cdot x_t + U_h h_{t-1} + b_h)$$

$$y_t = \sigma_y (W_y x_t + b_y)$$



PROBLEM ZNIKAJĄCYCH GRADIENTÓW

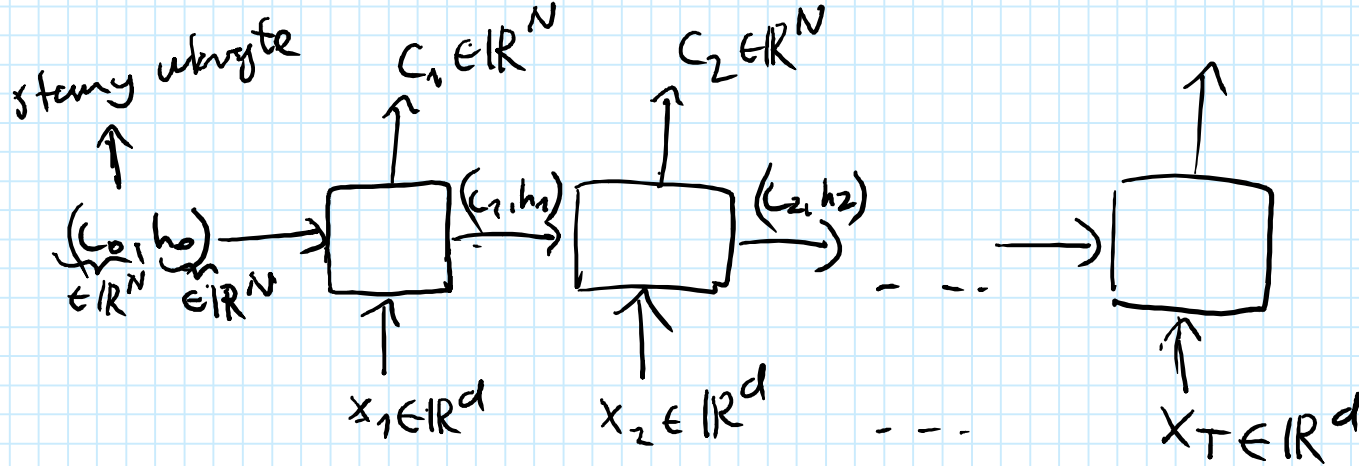


- Ostatni x_T przechodzi tylko przez jedną grupę neuronów ale x_1 przechodzi przez znacznie więcej, więc sieć jest efektywnie bardzo głęboka
- Jeśli T jest duże, to x_1 jest bardzo daleko od x_T i trudno mu „dojść do h_T ”.

Rozwiązanie: LSTM!

LSTM (Long Short-Term memory)

Hochreiter, Schmidhuber 1997



Liczba wag
 $4(Nd + N^2 + N)$

$$f_t = \sigma_g \left(\underbrace{W_f}_{\in \mathbb{R}^{N \times d}} x_t + \underbrace{U_f}_{\in \mathbb{R}^{N \times N}} h_{t-1} + \underbrace{b_f}_{\in \mathbb{R}^N} \right) \in \mathbb{R}^N$$

forget gate

$$i_t = \sigma_g \left(W_i x_t + U_i h_{t-1} + b_i \right)$$

input gate

$$o_t = \sigma_g \left(W_o x_t + U_o h_{t-1} + b_o \right)$$

output gate

$$\tilde{c}_t = \sigma_c \left(W_c x_t + U_c h_{t-1} + b_c \right)$$

σ_g - sigmoid

σ_c, σ_h - tanh

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

\odot - mnożenie po współrzędnych (Hadamard product)

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

- $f_t^{(k)} \leftarrow k\text{-ta współrzędna } f_t$
 $f_t^{(k)} \approx 1$
 $i_t^{(k)} \approx 0$
 \Rightarrow propagujemy tę k -tą współrzędną dalej
- $f_t^{(k)} \approx 0$
 $i_t^{(k)} \approx 1$
 \Rightarrow zapominamy tę k -tą współrzędną i zastępujemy ją nowym kandydatem \tilde{c}

GRU (Gated Recurrent Unit)

Kyunghyun Cho et al. "Learning Phrase Representation..." 2014

- mniej wag niż dla LSTM
- podobne działanie do LSTM
- brak output gate

$$\text{Liczba wag} = 3(Nd + N^2 + N)$$

$$z_t = \sigma_g \left(W_z x_t + U_z h_{t-1} + b_z \right) \in \mathbb{R}^N \quad \text{update gate} \quad \sigma_g - \text{sigmoid}$$

$$r_t = \sigma_g \left(W_r x_t + U_r h_{t-1} + b_r \right) \quad \text{reset gate} \quad \sigma_h - \text{tanh}$$

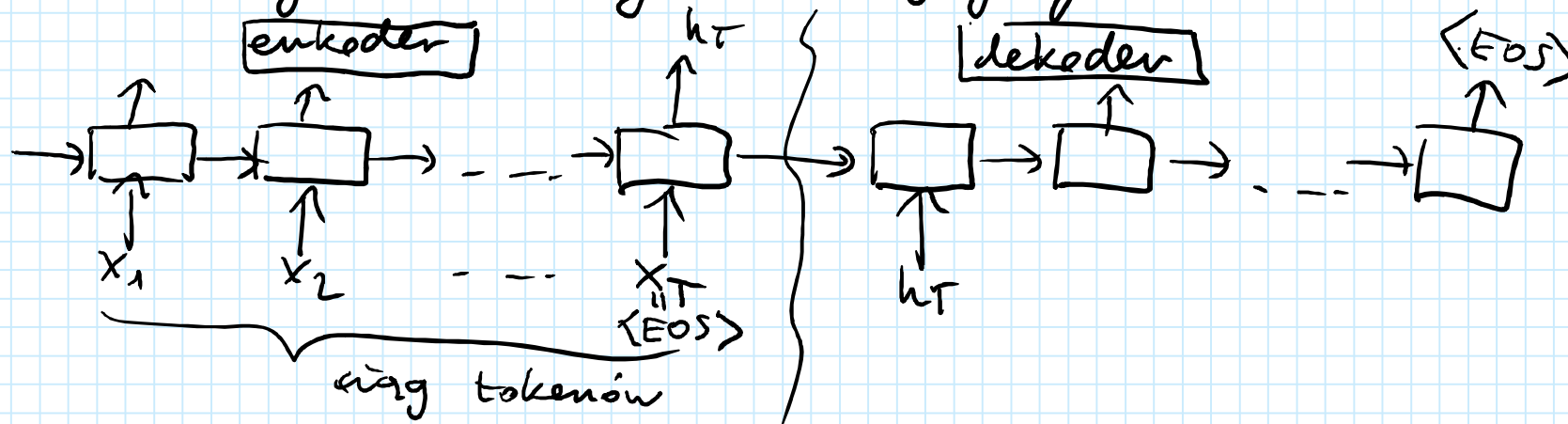
$$\hat{h}_t = \sigma_h \left(W_h x_t + U_h (r_t \circ h_{t-1}) + b_h \right) \quad \text{candidate activation}$$

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

- Jeśli $z_t \approx 0$ na pewnej współrzędnej, to na tej współrzędnej zostawiamy h_t bez zmian
- Jeśli $z_t \approx 1$ na pewnej współrzędnej, to na tej współrzędnej zapominamy przeszłość i zastępujemy ją nowym \hat{h}_t

Architektura encoder - decoder

- są 2 części rekurencyjne: encoder oraz decoder
- encoder koduje zdanie zapisane w jednym języku, a decoder je odekoduje na inny język.



$\left[\begin{array}{c} \vdots \\ \langle \text{UNK} \rangle \\ \langle \rangle \end{array} \right]$ – specjalne słowa np. $\langle \text{UNK} \rangle$ – nieznanym wyraz, którego nie ma w słowniku; $\langle \text{EOS} \rangle$ – end of sentence (koniec zdania)

- Repeat Vector
- return_sequences = TRUE