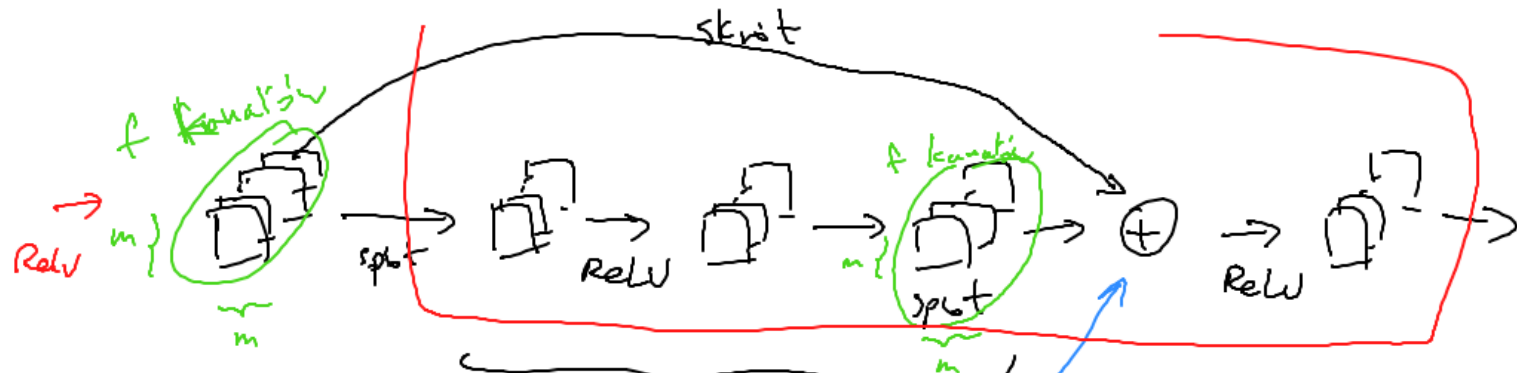


ResNet (residual neural network)



2 splity przed i po warstwą aktywacyjną Relu (czasami bierze się 3)

zły kodowanie po wyznaczeniu było wykonane bez dodatkowych zabiegów, to wymiar warstw zwracanych na siebie muszą być identyczne (Można zrobić to więcej różnic gdy siebie warstwy nie są identyczne, ale wtedy trzeba coś dodać).

Pomysł polega na tym, że gdyby warstwy w splitach w części czerwonej były takie, to część czerwona jest (prawie) identyczna.

Łączenie składowych rozwiązań problem znikających gradientów.

Sequential

Optymalizacja:

• SGD (stochastyczny gradient descent) — to jest to, co było
jak w regresji liniowej

metoda opisana dla sieci neuronowych (głębokich)
(on-line lub batch)

• przy użyciu się SGD, to zwykle stopniowo zmniejsza się stopień uczenia

• często używa się dodatkowo metody „bezwładności” (momentum):

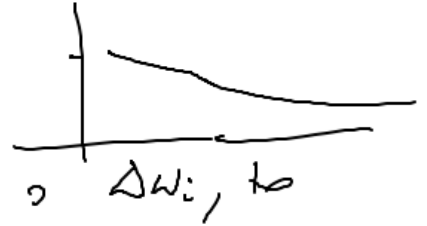
jeśli w i -tym kroku aktualniemy pewną wagę w o Δw_i , to
w kroku $i+1$:

$$w := w - \eta \cdot \frac{\partial L}{\partial w} + \alpha \cdot \Delta w_i$$

↑
stopień uczenia

↑
stopień bezwładności $\in [0, 1)$
często $\alpha = 0.9$

~~...~~

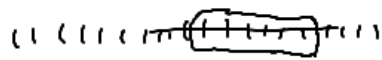


• adeptacyjne : duża pojemność reaktorów od konkretnych wag

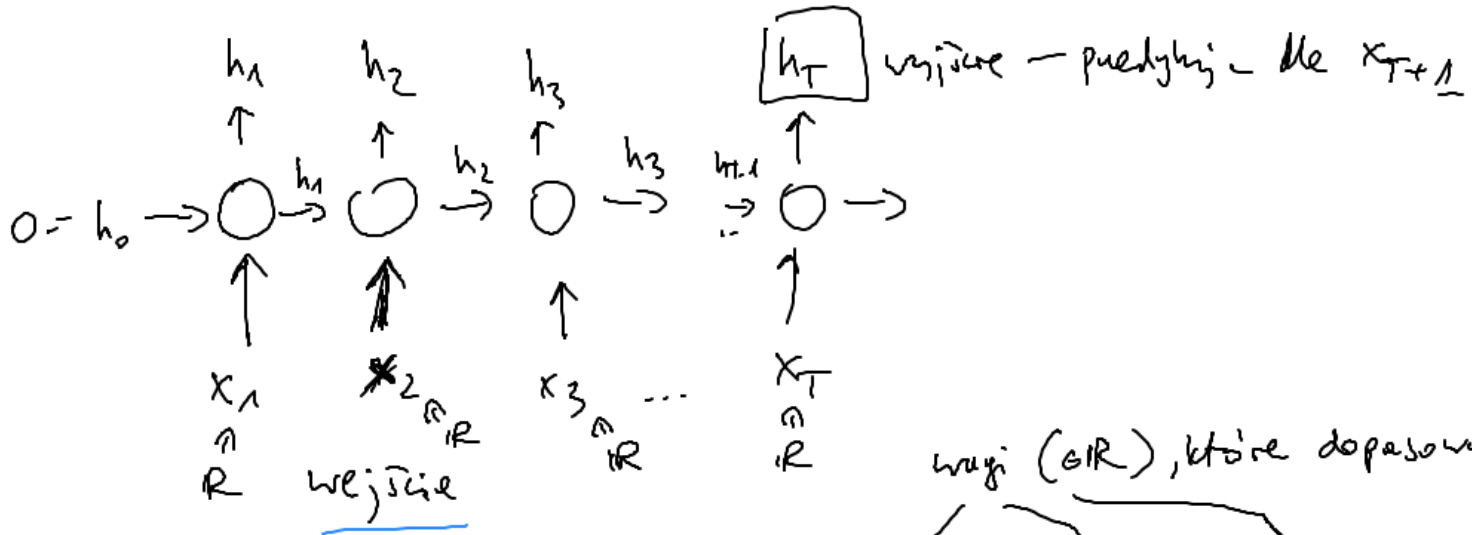
AdaGrad, Ada, RMSprop, Adam

wiele zmian w konfigacji mikro architektury sieci (tzw. framework)

Rekurencyjne sieci neuronowe (RNN)



zaczynamy od prostego przypadku pojedynczego wejścia, warstwy, komórki



wagi ($\in \mathbb{R}$), które dopasujemy

prosta komórka RNN:
$$h_t = \sigma(W \cdot x_t + U \cdot h_{t-1} + b) \quad , t = 1, 2, 3, \dots$$

↑
wejście

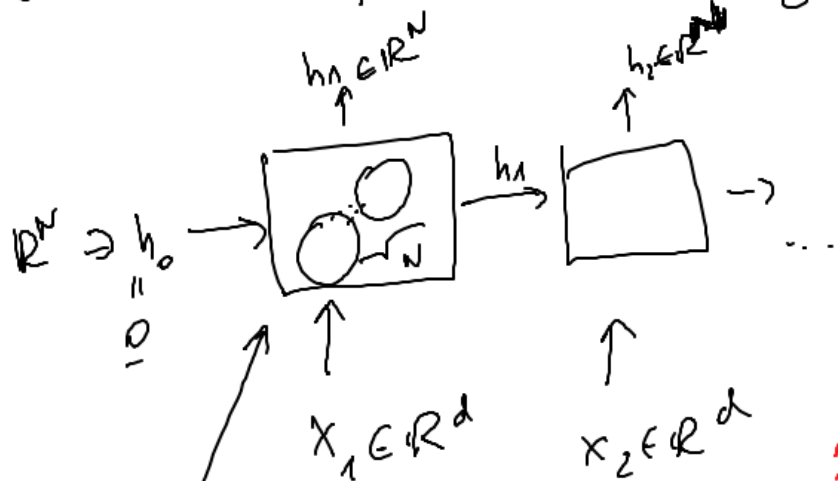
↑
poprzednie wyjście ("historia")

da

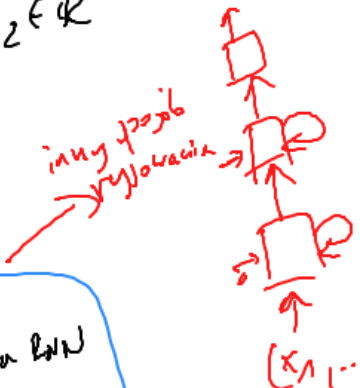
σ - funkcja aktywacji, np.
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Uogólnienie na przepływe wyliczeniowy:

Linia permutiv: $Nd + N^2 + N$



warstwa skrocona z wielk komorek (N)

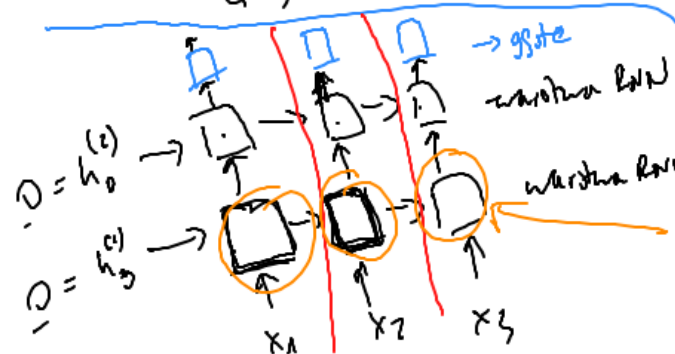


Wagi, które są dopasowywane w procesie uczenia

(działanie współrzędnych)

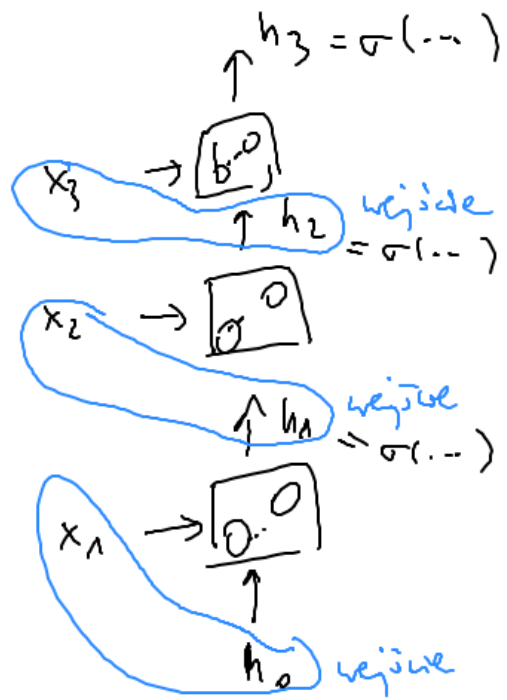
$$h_t = \sigma(W \cdot x_t + U \cdot h_{t-1} + b) = \begin{bmatrix} \sigma(\dots) \\ \sigma(\dots) \\ \vdots \\ \sigma(\dots) \end{bmatrix} \in \mathbb{R}^N$$

$N \times d$ $d \times 1$ $N \times N$ \mathbb{R}^N \mathbb{R}^N



są to same wagi!
jest wyliczana operacja tego samego typu, ale dla innych wejść

pojedyncza warstwa RNN dla $T=3$



LSTM (long-short term memory)

95-97-99 Hochreiter, Schmidhuber

warstora 2N kwadrata

parametrow jest $4(Nd + N^2 + N)$

$$\mathbb{R}^N \rightarrow f_t = \sigma_g \left(\overbrace{W_f}^{N \times d} \cdot \overbrace{x_t}^{\in \mathbb{R}^d} + \overbrace{U_f}^{N \times N} \cdot \overbrace{h_{t-1}}^{\in \mathbb{R}^N} + b_f \right) \in [0,1]^N$$

forget gate

$$\sigma_g(x) = \frac{1}{1 + e^{-x}}$$

$$\mathbb{R}^N \rightarrow i_t = \sigma_g \left(\overbrace{W_i}^{N \times d} \cdot \overbrace{x_t}^d + \overbrace{U_i}^{N \times N} \cdot \overbrace{h_{t-1}}^N + b_i \right) \in [0,1]^N$$

input gate

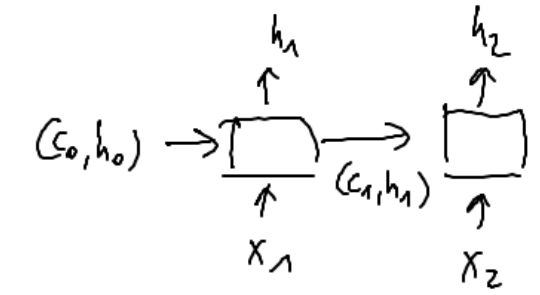
$$\sigma_c(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\mathbb{R}^N \rightarrow o_t = \sigma_g \left(\overbrace{W_o}^{N \times d} \cdot \overbrace{x_t}^d + \overbrace{U_o}^{N \times N} \cdot \overbrace{h_{t-1}}^N + b_o \right) \in [0,1]^N$$

output gate

$$\mathbb{R}^N \rightarrow \tilde{c}_t = \sigma_c \left(\overbrace{W_c}^{N \times d} \cdot \overbrace{x_t}^d + \overbrace{U_c}^{N \times N} \cdot \overbrace{h_{t-1}}^N + b_c \right)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \in \mathbb{R}^N \leftarrow \text{memory cell}$$



$$\mathbb{R}^N \rightarrow h_t = o_t \odot \sigma_h(c_t)$$

← wyjście i bieżący stan

↑ \mathbb{R}^N ↑ \mathbb{R}^N
 ↑ po wyprodukowaniu
 ↑ wyjście i bieżący stan

$$(a_1, \dots, a_n) \odot (b_1, \dots, b_n) = (a_1 b_1, a_2 b_2, \dots, a_n b_n)$$

Dalsza modyfikacja: GRU (gated recurrent unit)

prze wszystkim jedna brama mniej niż w LSTM

ale zwykle podobne działanie