



MSI

12. Sieci neuronowe głębokie (DNN)

Włodzimierz Kasprzak

Treść

1. Głębokie sieci neuronowe
2. Splotowe sieci neuronowe CNN
3. Rekurencyjna sieć neuronowa RNN
4. Warianty sieci RNN

1. Głębokie sieci neuronowe

Typowe architektury głębokich sieci neuronowych (ang. deep neural networks, DNN):

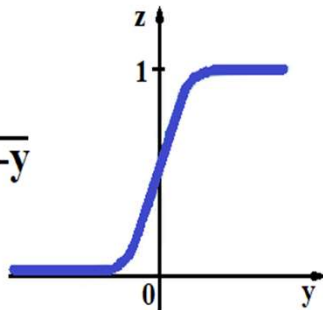
1. Klasyfikacja wzorców poprzez sieci wielowarstwowe („w pełni połączone”) z warstwą wyjściową typu „soft-max”;
2. Splotowe sieci neuronowe (ang. *convolutive neural networks*, CNN) do wyznaczania cech i klasyfikacji obrazów 2D/3D;
3. Sieci splotowe R-CNN i jej modyfikacje;
4. Rekurencyjne sieci neuronowe – RNN, LSTM - do modelowania sekwencji czasowych;
5. *Głębokie auto-enkodery.*
6. *Sieć antagonistyczna;*
7. *Sieci korelacyjne i syjamskie;*
8. *Grafowe sieci (grafowe sieci splotowe).*

Funkcje aktywacji

Typowe funkcje aktywacji (funkcja generująca wyjście neuronu na podstawie pobudzenia od wejść i wag połączeń)

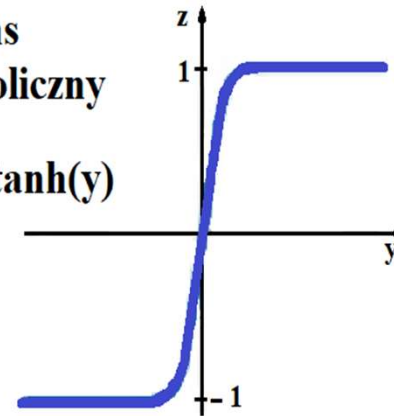
Funkcja
sigmoidalna

$$z(y) = \frac{1}{1 + e^{-y}}$$



Tangens
hiperboliczny

$$z(y) = \tanh(y)$$

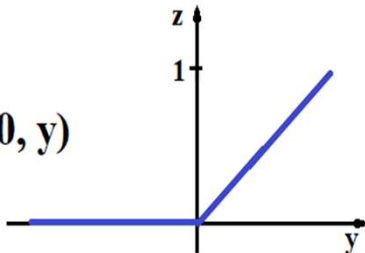


ELU

$$z(y) = \begin{cases} y, & \text{dla } y \geq 0 \\ \alpha(e^y - 1), & \text{dla } y < 0 \end{cases}$$

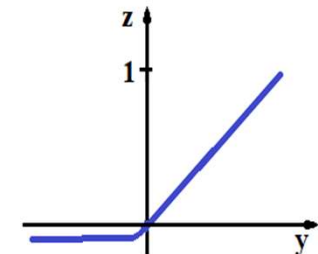
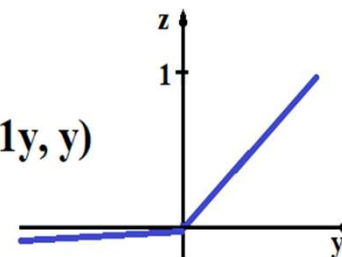
ReLU

$$z(y) = \max(0, y)$$



Leaky ReLU

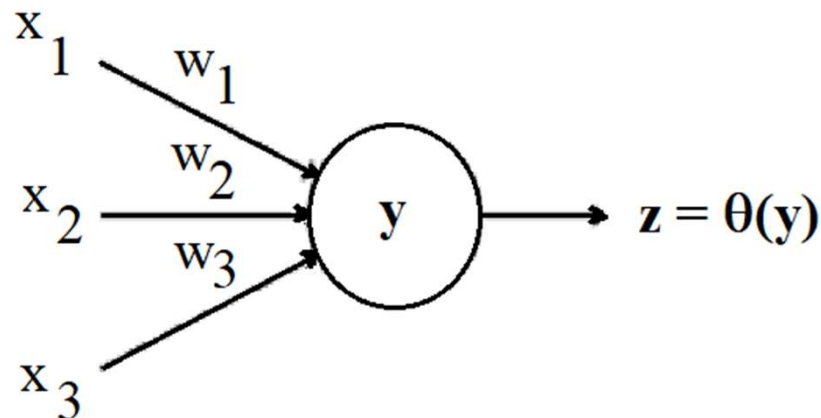
$$z(y) = \max(0.1y, y)$$



Sieci głębokie „w pełni połączone”

- Najprostsze sieci głębokie posiadają warstwy wyjściowe ”w pełni połączone” (ang. *Fully Connected Layer*, FC). Np. wielowarstwowy perceptron jest już przykładem sieci o w pełni połączonych warstwach.
- Wyjście każdego neuronu to wartość funkcji aktywacji dla sumy iloczynów wszystkich neuronów wejściowych i ich wag.

Np.



2. Splotowe sieci neuronowe CNN

Splotowe sieci neuronowe (ang. *convolutive neural networks*, CNN) to sieci jednokierunkowe, które posiadają specyficzną organizację neuronów w warstwy, dostosowaną do danych 2-wymiarowych (np. obrazów cyfrowych) posiadających topologię macierzy 2D (kraty).

Warstwa splotowa realizuje jednoczesną filtrację neuronów wejściowych za pomocą zbioru filtrów splotowych charakteryzowanych 2-wymiarowymi macierzami liczb, tzw. „jądrami” filtra. Jest to część liniowa (funkcja pobudzenia). Typowo stosowaną funkcją aktywacji jest **ReLU**.

Warstwa łącząca („pooling”) zazwyczaj następuje po każdej warstwie splotowej. Odpowiednia funkcja próbkuje lokalne kraty wyników w warstwie splotowej zmniejszając rozdzielczość kraty wyników. Typowe funkcje próbkujące to: „*max pooling*”, która zwraca maksymalną wartość lokalnej kraty, lub „*Lp-pooling*”, która realizuje wygładzanie wyników lokalnej kraty.

Splotowe sieci neuronowe CNN

Część wyjściowa pełni rolę klasyfikatora i zazwyczaj ma postać **dwuwarstwowego perceptronu** (nazywanego siecią „w pełni połączoną” (*fully connected*) lub o „gęstych” połączeniach (*dense*), przy czym warstwa wyjściowa jest zwykle typu „**softmax**” – określa rozkład prawdopodobieństwa *a posteriori* klas.

2-wymiarowa operacja splotu

Dane: mapa wejściowa X,
maska filtra f.

Wynik: mapa cech H

Krok przesunięcia maski
„stride” = (S_w, S_h)

Np. dla filtra o rozmiarze
 $(f_w, f_h) = (3, 3)$:

$$h_{13} = \sum_{i=1}^3 \sum_{j=1}^3 w_{i,j} x_{i,j+2} + b$$

x_{11}	x_{12}	x_{13}				x_{17}
x_{21}						
x_{71}	x_{72}					x_{77}

Mapa wejściowa X
o rozmiarze
 $(X_{\text{width}}, X_{\text{height}})$

w_{11}	w_{12}	w_{13}
w_{21}	w_{22}	w_{23}
w_{31}	w_{32}	w_{33}

Maska
filtra f
o rozmiarze
 (f_w, f_h)

h_{11}	h_{12}	h_{13}		h_{15}
h_{21}				
h_{51}				h_{55}

Mapa cech H
o rozmiarze
 $(H_{\text{width}}, H_{\text{height}})$

x_{11}	x_{12}	x_{13}				x_{17}
x_{21}						
x_{71}	x_{72}					x_{77}

Mapa X

h_{11}	h_{12}	h_{13}		h_{15}
h_{21}				
h_{51}				h_{55}

Mapa cech H

Warstwa splotowa

Kanały wyjściowe: liczba f_{num} masek splotowych (filtrów) stosowanych w warstwie sieci o X_{num} mapach wejściowych (tzw. kanałów) prowadzi do wielokrotnej liczby map wyjściowych H (tzw. kanałów) dla jednej mapy wejściowej:

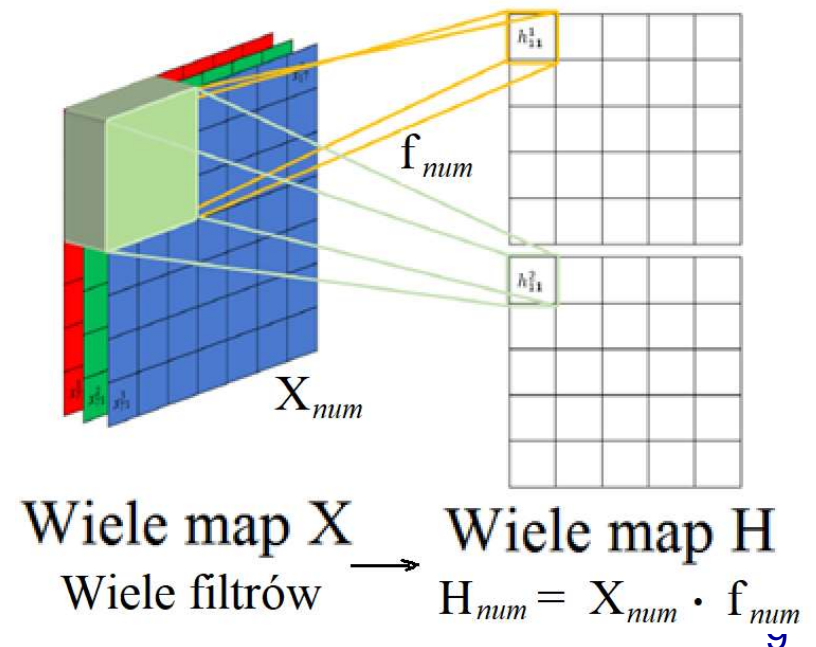
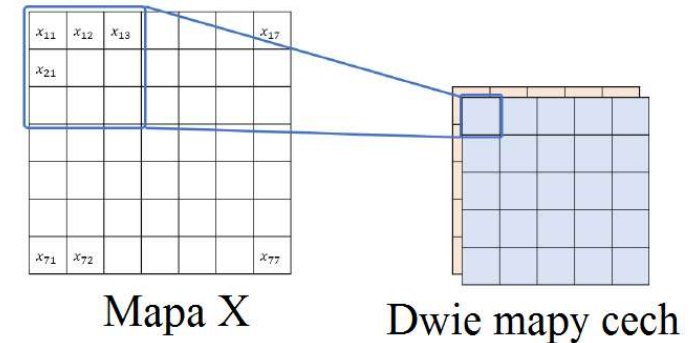
$$H_{num_{wyjścia}} = X_{num_{wejścia}} \cdot f_{num}$$

Np. dwa filtry w warstwie generują dwie mapy (kanały) wyjściowe dla każdej mapy wejściowej.

Wiele map wejściowych
i wielokrotne kanały wyjściowe:

$$H_{num_{wyjścia}} = X_{num_{wejścia}} \cdot f_{num}$$

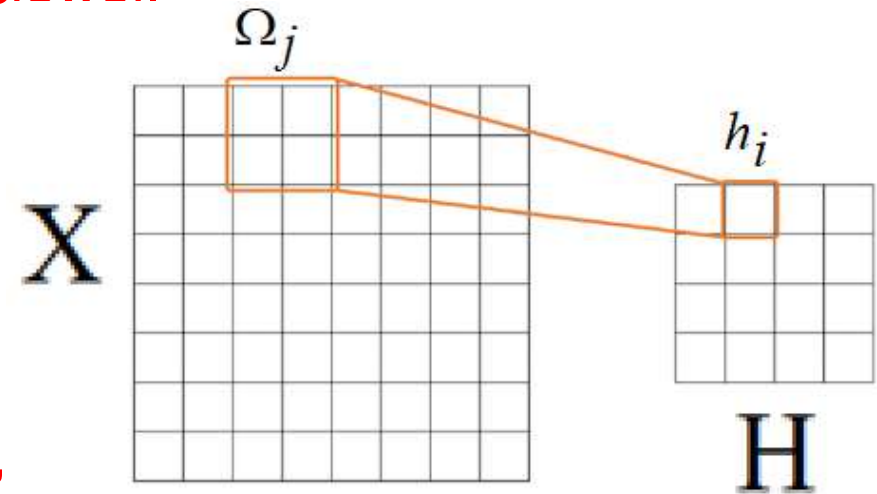
Padding (P_h, P_w) oznacza rozszerzenie mapy wejściowej o P_h wierszy i P_w kolumn dla pełnego wykorzystania brzegowych danych mapy.



Próbkowanie („pooling”)

„Pooling”: redukcja rozmiaru mapy wejściowej.

Np. $pool = 2 \times 2$: każdy podobszar Ω_j o rozmiarze 2×2 wejściowej mapy cech X odwzorowywany jest na jeden element mapy wyjściowej H .



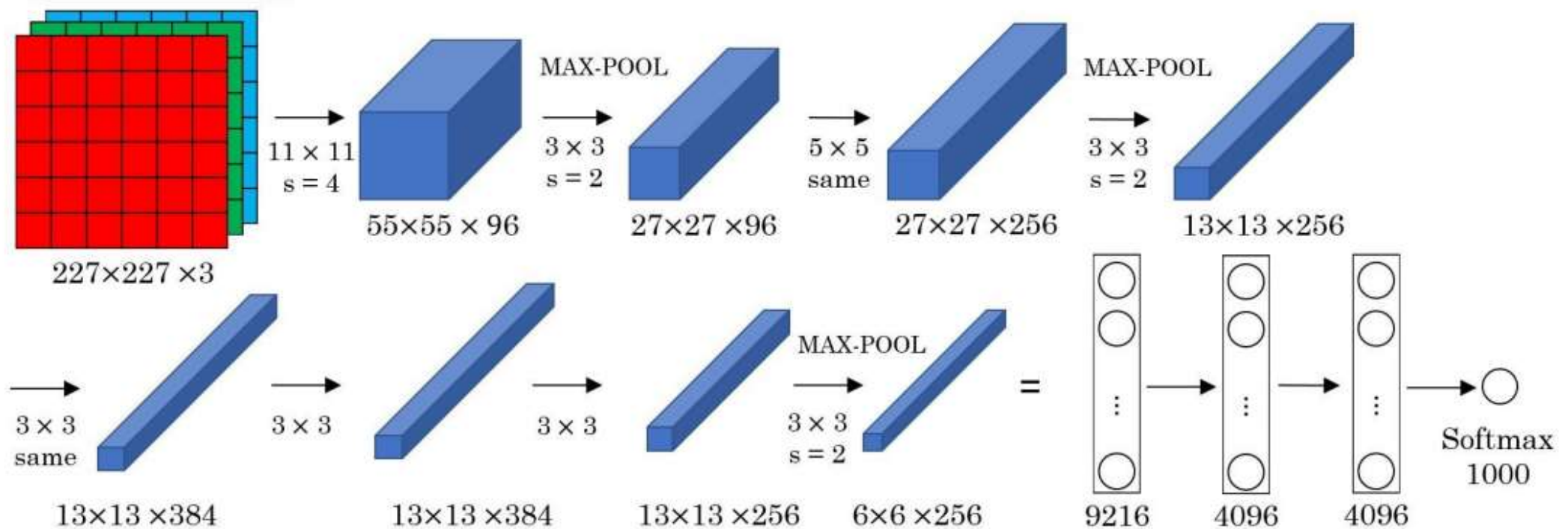
Funkcje stosowane w warstwie „pooling”.

- „Max pooling”: określa maksymalną wartość w każdym podobszarze Ω_j ;
- „Avg pooling”: średnia (lub suma) - określa wartość średnią lub sumę elementów podobszaru Ω_j ;
- „ L_p -pooling”: dla każdego obszaru Ω_j określa normę rzędu p (L_p), czyli

$$h_j = \left(\sum_{i \in \Omega_j} x_i^p \right)^{1/p}$$

Przykład sieci CNN

Sieć „AlexNet”



<https://harangdev.github.io/deep-learning/convolutional-neural-networks/25/>

- Funkcja aktywacji ReLU.
- Trenowana na ImageNet
- Posiada ok. 60 mln parametrów

3. Rekurencyjna sieć neuronowa (RNN)

Dana jest sekwencja obserwacji w czasie: $\mathbf{x}_{1:t} = \{\mathbf{x}_1, \dots, \mathbf{x}_t\}$.

Warstwa rekurencyjna sieci RNN

$$\mathbf{h}_t = \mathbf{f}_h (\mathbf{V}_h \mathbf{x}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h)$$

$$\mathbf{y}_t = \mathbf{f}_y (\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y)$$

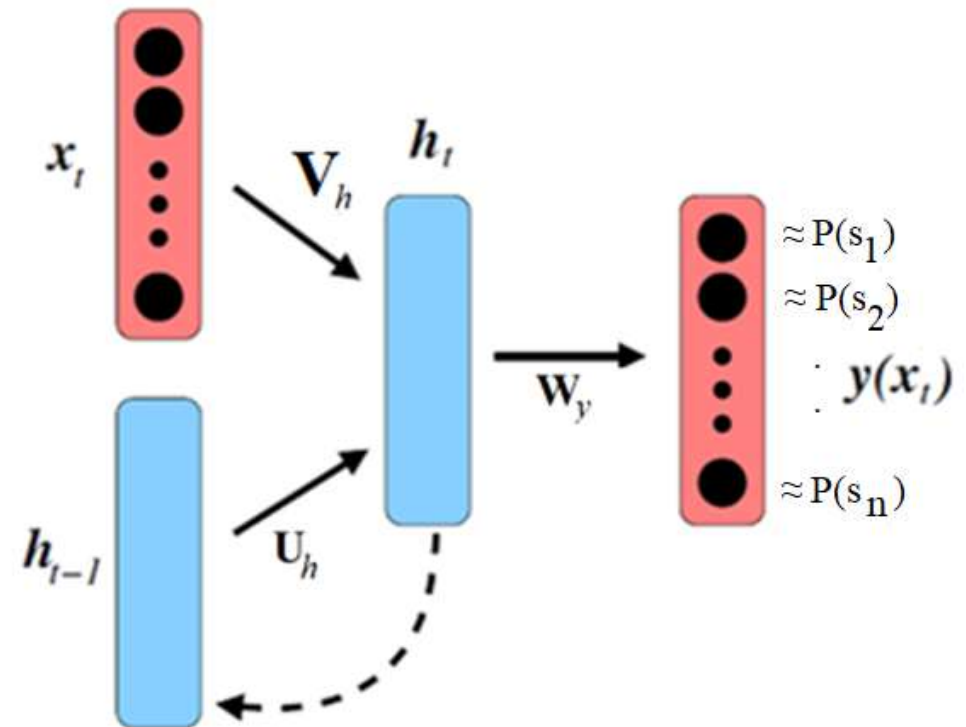
RNN realizuje aproksymację funkcji, $\mathbf{y}(\mathbf{x}_{1:t}) \cong \mathbf{F}(\mathbf{x}_{1:t})$, przy założeniu:

$$\mathbf{F}(\mathbf{x}_{1:t}) = \{ \mathbf{F}(\mathbf{x}_i, \mathbf{x}_{1:i-1}) \mid i=1, \dots, t \}$$

$$\cong \{ \mathbf{F}(\mathbf{x}_i, \mathbf{h}_{i-1}) \mid \dots \}$$

$$\cong \{ \mathbf{F}(\mathbf{h}_i) \mid \dots \} = \{ \mathbf{y}_i \mid i=1, \dots, t \} = \mathbf{y}_{1:t}$$

Efektom jest sekwencja wyjściowa o identycznej długości co sekwencja wejściowa: $\mathbf{x}_{1:t} \rightarrow \mathbf{y}_{1:t}$



Sieć głęboka RNN

Zauważmy, że warstwa rekurencyjna może zostać wielokrotnie powielona, podobnie jak możliwych jest wiele warstw ukrytych w sieci jednokierunkowej MLP.

Wprowadźmy indeks górny (i) reprezentujący numer kolejnej warstwy ukrytej rekurencyjnej, $i = 1, 2, \dots, L$

$$\mathbf{h}_t^{(i)} = \mathbf{f}_h (\mathbf{W}_h^{(i)} \mathbf{h}_t^{(i-1)} + \mathbf{U}_h^{(i)} \mathbf{h}_{t-1}^{(i)} + \mathbf{b}_h^{(i)}) , \text{ gdzie } \mathbf{h}_t^{(1)} = \mathbf{x}_t$$

Wyjście sieci:

$$\mathbf{y}_t = \mathbf{f}_y (\mathbf{W}_y \mathbf{h}_t^{(L)} + \mathbf{b}_y)$$

Daje to możliwość modelowania bardzo złożonych zależności pomiędzy danymi a także prowadzi do znaczącego oddzielenia poziomu zmiennych „stanu” od zmiennych obserwacji.

4. Warianty sieci RNN

4.1 Sekwencja czasowa sieci RNN

- jednokierunkowa
- dwukierunkowa

4.2 Modyfikacja warstwy RNN (dodanie „bramek”):

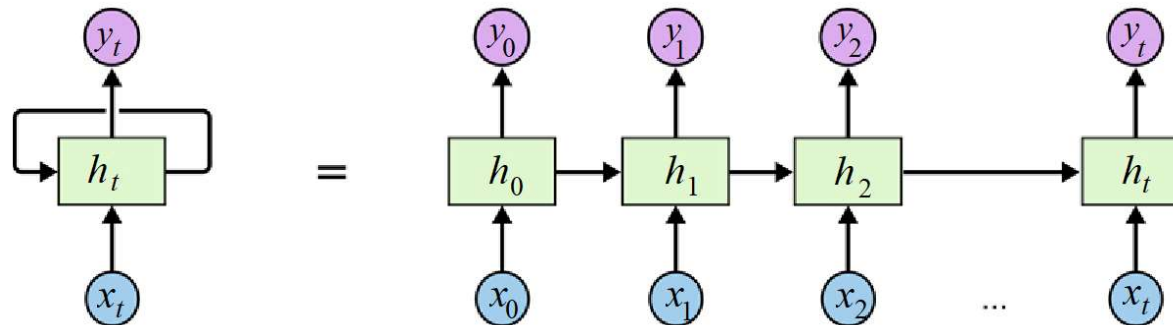
- bramkowane neuronów (GRU)

4.3 Bramkowanie sekwencji czasowej

- sieć „long-short term memory” (LSTM)

4.1 Sekwencja „czasowa” sieci RNN

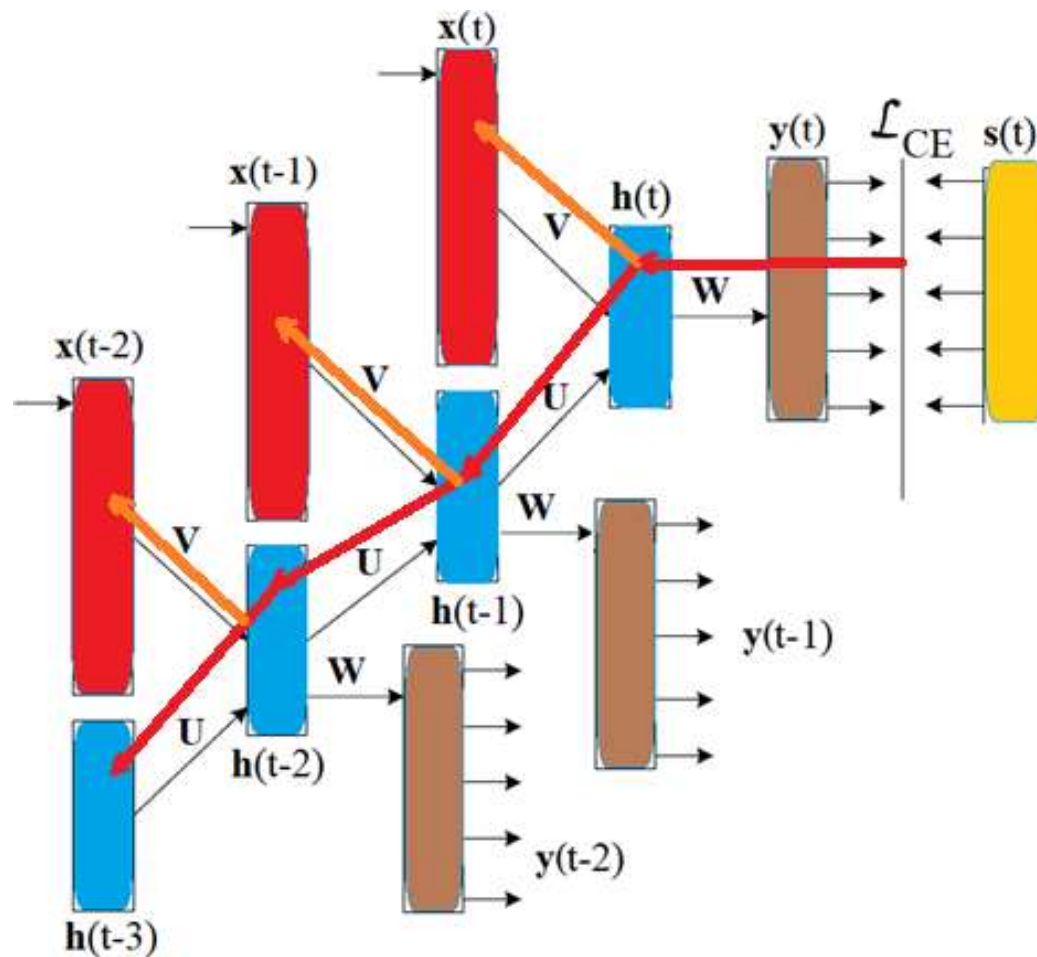
Rozwinięciem sieci RNN jest architektura **sekwencji czasowej takich sieci**, modelująca zależności czasowe dłuższe niż 2 kolejne chwile czasowe. Neurony z połączeniami rekurencyjnymi działają jak pamięć - potencjalnie są one w stanie modelować relacje występujące **w dowolnie długim zakresie** czasowym.



Umożliwia to lepsze modelowanie zależności czasowych a proces uczenia wag sieci oparty jest o dłuższe sekwencje czasowe danych niż o długości 2. Podstawowym algorytmem uczącym takiej sieci jest **wsteczna propagacja względem czasu** BPTT („*BackProgration Through Time*”), a w praktyce algorytm „przyciętego” BPTT („truncated BPTT”).

Uczenie sekwencji sieci RNN (1)

Proces uczenia wag można przedstawić w postaci uczenia równoważnej sieci jednokierunkowej, rozwiniętej z sieci podstawowej RNN w sekwencję „czasową”, o pewnej liczbie K powtórzeń sieci RNN „indeksowanych czasem”, w których wagi każdej z sieci, V , U i W oraz wagi wejść hamujących (*bias*), są wspólne dla wszystkich warstw w „czasie”.



BPTT - uczenie sekwencji sieci RNN

Dane uczące dla sekwencji K sieci RNN to N par danych:

$\langle \mathbf{x}_1, \mathbf{s}_1 \rangle, \langle \mathbf{x}_2, \mathbf{s}_2 \rangle, \dots, \langle \mathbf{x}_N, \mathbf{s}_N \rangle$, gdzie \mathbf{x}_i to dane wejściowe a \mathbf{s}_i odpowiadające im oczekiwane dane wyjściowe.

Nieznane początkowe wartości warstwy ukrytej \mathbf{h}_0 przyjmuje się za zerowe. Dla uczenia sieci neuronowych istnieją różne formy funkcji straty (optymalizowanej funkcji celu) ale dla RNN przyjmuje się zwykle entropię krzyżową.

Równania sieci o indeksie czasu t :

$$y_k(t) = \text{softmax} \left(\sum_{r=1}^{n_H} w_{kr} h_r(t) + b_k \right)$$
$$h_j(t) = \text{sigmoid} \left(\sum_{s=0}^{n_X} v_{js} x_s(t) + \sum_{r=0}^{n_H} u_{jr} h_r(t-1) + b_j \right)$$

BPTT - uczenie sekwencji sieci RNN

Back_Propagation_Through_Time($x[]$, $s[]$, K , N)

// $x[t]$ to wektor wejściowy w chwili t ; $y[t]$ to oczekiwany wektor wyjściowy
// w chwili t ; K to długość sekwencji sieci RNN; N to długość sekwencji
// danych uczących

POWTARZAJ do spełnienia warunku końca (np. liczba epok)

$h[0:K-1] := 0$;

FOR $t = 1, \dots, N - K - 1$ DO //

wejście_sieci $\leftarrow (h[0:K-1], x[t], x[t+1], \dots, x[t+K])$

$y[t+K] := \text{propagacja_wprzód}(x[], h[])$

$e := L_{CE}(s[t+K], y[t+K])$; // funkcja straty

$(\Delta W, \Delta V[K], \Delta V[K-1], \dots, \Delta V[1], \Delta U[K], \Delta V[K-1], \dots, \Delta U[1]) =$
propagacja_wstecz_błędu(e);

$\Delta V = \Delta V[K] + \Delta V[K-1] + \dots + \Delta V[1]$; // suma zmian wag

$\Delta U = \Delta U[K] + \Delta U[K-1] + \dots + \Delta u[1]$; // suma zmian wag

$(W, V, U) = \text{modyfikuj_wagi}(\Delta W, \Delta V, \Delta U)$

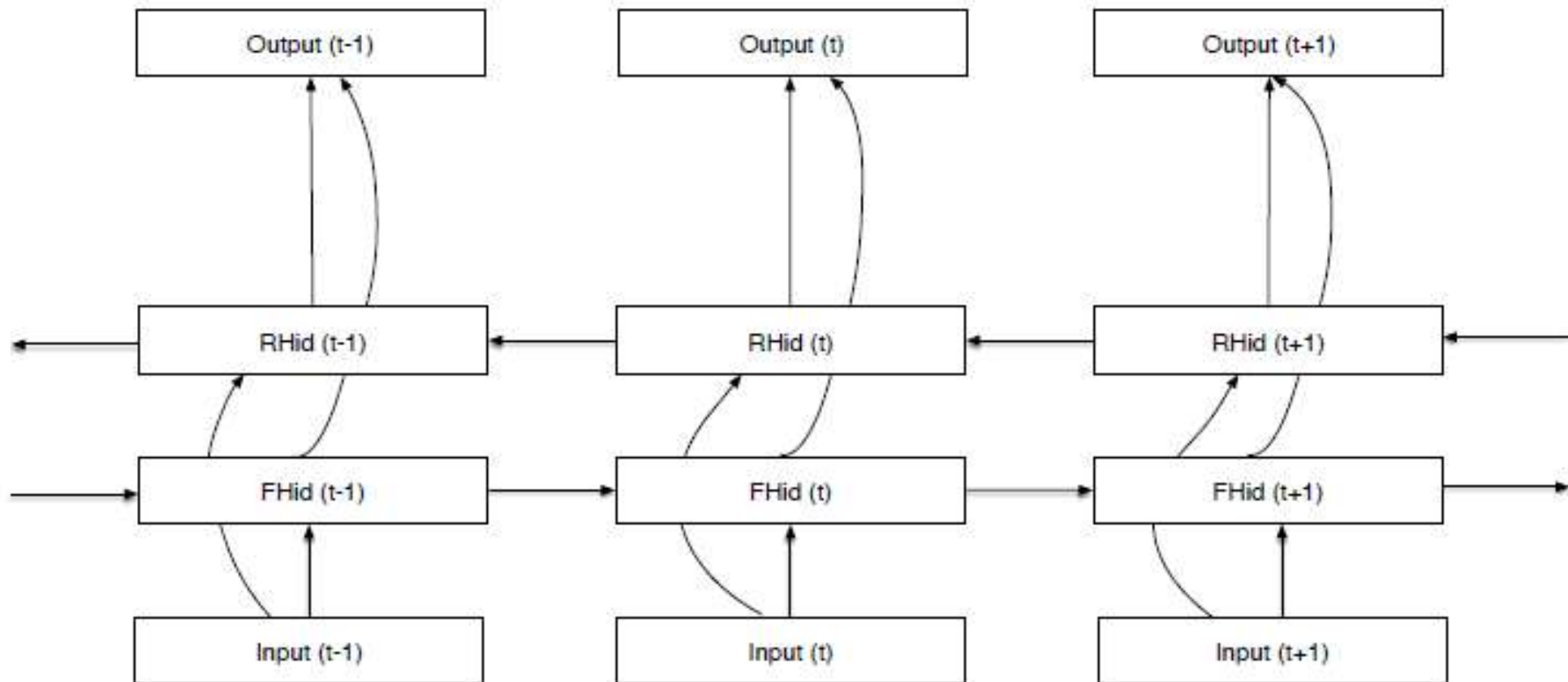
$h[0:K-1] := h[1:K]$; // pamięć dla następnej iteracji

END

END
MSI

Dwukierunkowa sekwencja sieci RNN

Wzajemne połączenie dwóch kolejnych warstw ukrytych sieci RNN – daje możliwość uwzględnienia zależności dwukierunkowych w czasie.



4.2 GRU RNN

Gated recurrent unit (GRU)

Dodatkowo występują komórki pamięci przeznaczone do sterowania zapamiętywaniem informacji przez elementy rekurencyjne. Są to:

- bramka aktualizacji (*update gate*) z_t ,
- bramka „zerowania” (*reset gate*) r_t

Aktualizacja neuronu w warstwie ukrytej (rekurencyjnej) ma postać:

$$\mathbf{h}_t = (1 - z_t) \odot \mathbf{h}_{t-1} + z_t \odot \widetilde{\mathbf{h}}_t$$

gdzie $\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1} + \mathbf{b}_z)$,

$$\widetilde{\mathbf{h}}_t = \tanh(\mathbf{W} \mathbf{x}_t + \mathbf{U}(r_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_h)$$

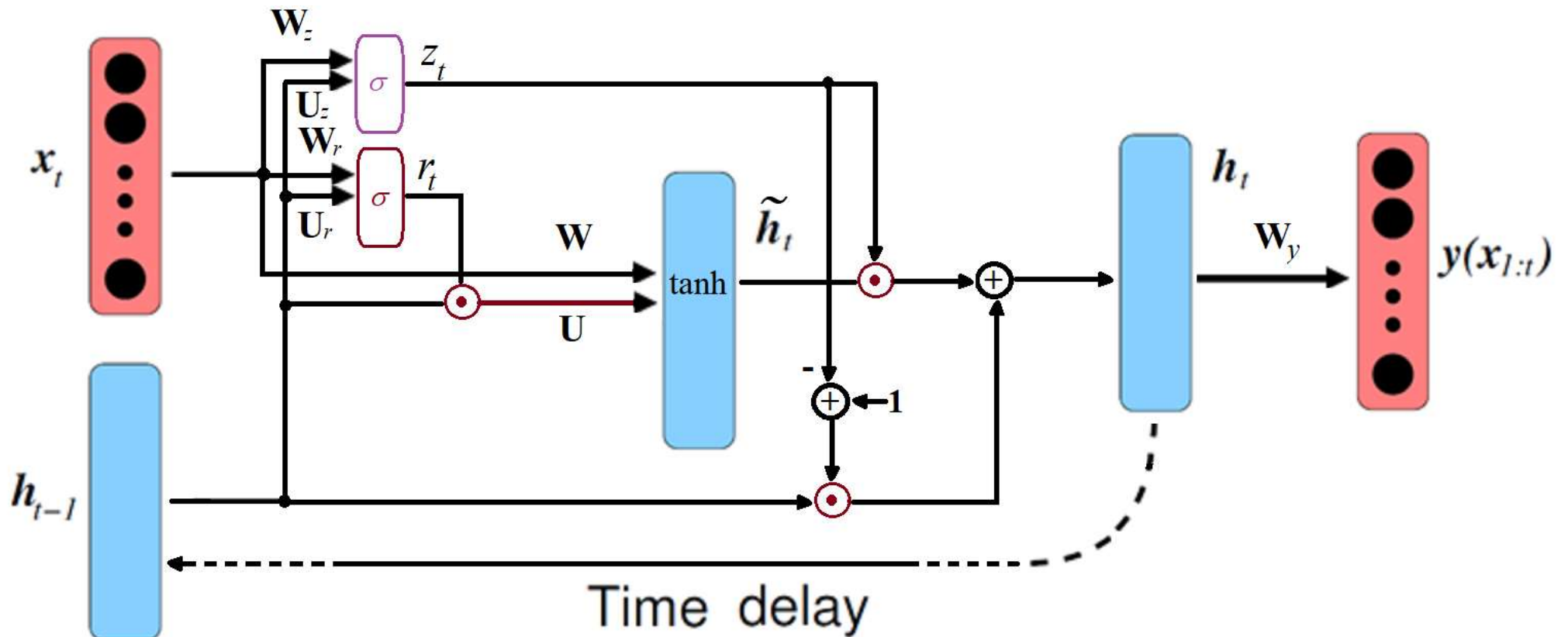
$$\mathbf{r}_t = \sigma(\mathbf{W}_r \mathbf{x}_t + \mathbf{U}_r \mathbf{h}_{t-1} + \mathbf{b}_r)$$

\odot oznacza mnożenie *Hadamarda* (mnożenie punkt-po-punkcie),
 $\sigma()$ jest zwykłą funkcją sigmoidalną.

Warstwa wyjściowa nie zmienia się: $\mathbf{y}(\mathbf{x}_{1:t}) = \mathbf{f}^f(\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y)$

GRU RNN

Gated recurrent unit (GRU)



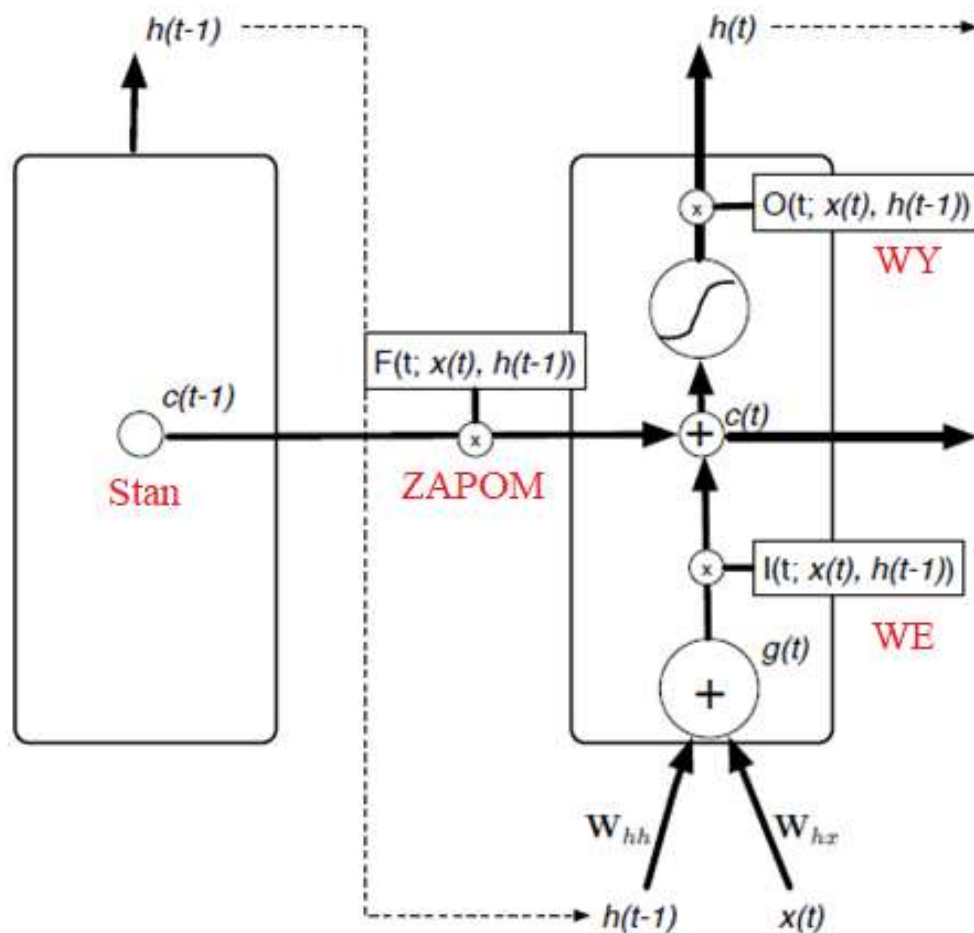
4.3 Sieć LSTM

Specyficznym rozwiązaniem sekwencji sieci rekurencyjnych stosowanym w modelowaniu sekwencji czasowych jest **sieć LSTM** (ang. *Long Short-Term Memory*):

- Każdy element warstwy ukrytej posiada wewnętrzną pamięć (**stan C**) - steruje ona usuwa problem zanikania gradientu przy długich opóźnieniach, umożliwia podtrzymywanie odległych informacji.
- Wprowadza **bramki** dla selekcji informacji. Wyróżnia się w niej bramkę wejściową, bramkę wyjściową i bramkę zapominania:
 - **Bramka wejściowa I** – określa to, która *nowa* wartość wpływa na wynik,
 - **Bramka wyjściowa O** – określa to, które wartości są używane do obliczeń wartości funkcji aktywacji,
 - **Bramka zapominania F** – określa to, które wartości są zapominane.

Funkcje sieci LSTM w warstwie ukrytej

Stan, wagi połączeń, bramki i ich funkcje:



$$\begin{aligned}
 I(t) &= \sigma(W_{ix}x(t) + W_{ih}h(t-1) + b_i) \\
 F(t) &= \sigma(W_{fx}x(t) + W_{fh}h(t-1) + b_f) \\
 O(t) &= \sigma(W_{ox}x(t) + W_{oh}h(t-1) + b_o) \\
 g(t) &= W_{hx}x(t) + W_{hh}h(t-1) + b_h \\
 c(t) &= F(t) \circ c(t-1) + I(t) \circ g(t) \\
 h(t) &= O(t) \circ \tanh(c(t))
 \end{aligned}$$

σ - funkcja sigmoidalna
 \circ - mnożenie punktowe

Pytania

1. Przedstawić funkcje aktywacji stosowane w głębokich sieciach neuronowych.
2. Omówić budowę sieci CNN.
3. Przedstawić warianty sieci rekurencyjnej RNN