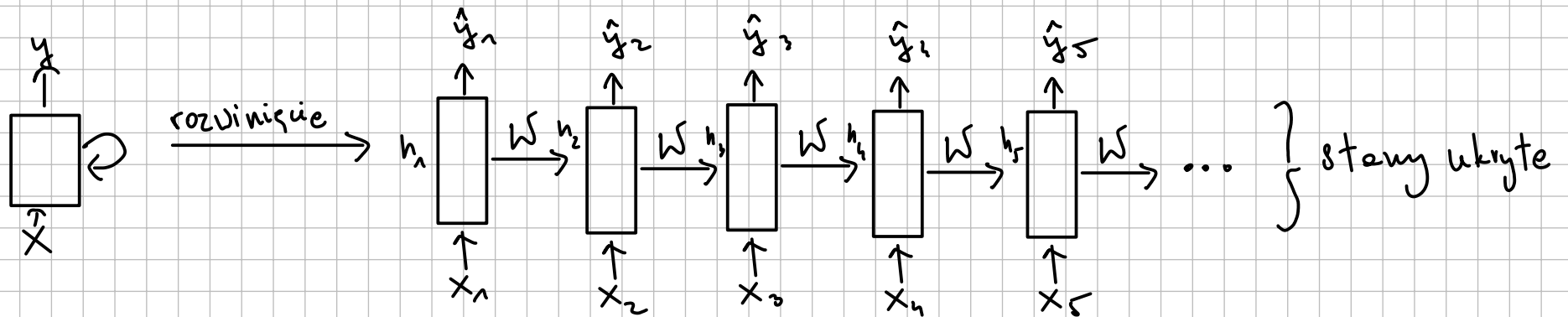


RNN (Recurrent neural network)

- Rozwiązanie problemu wejścia o zmiennej długości
- Główna idea - użycie tej samej macierzy wag W w wielu krokach



Bardziej dokładny schemat modelu (językowego)

x_i - słowa, zakodowane jako

one-hot-encoding

"word embedding"

[przedstawienie wektorowe]

$$e_i = E x_i$$

stany ukryte

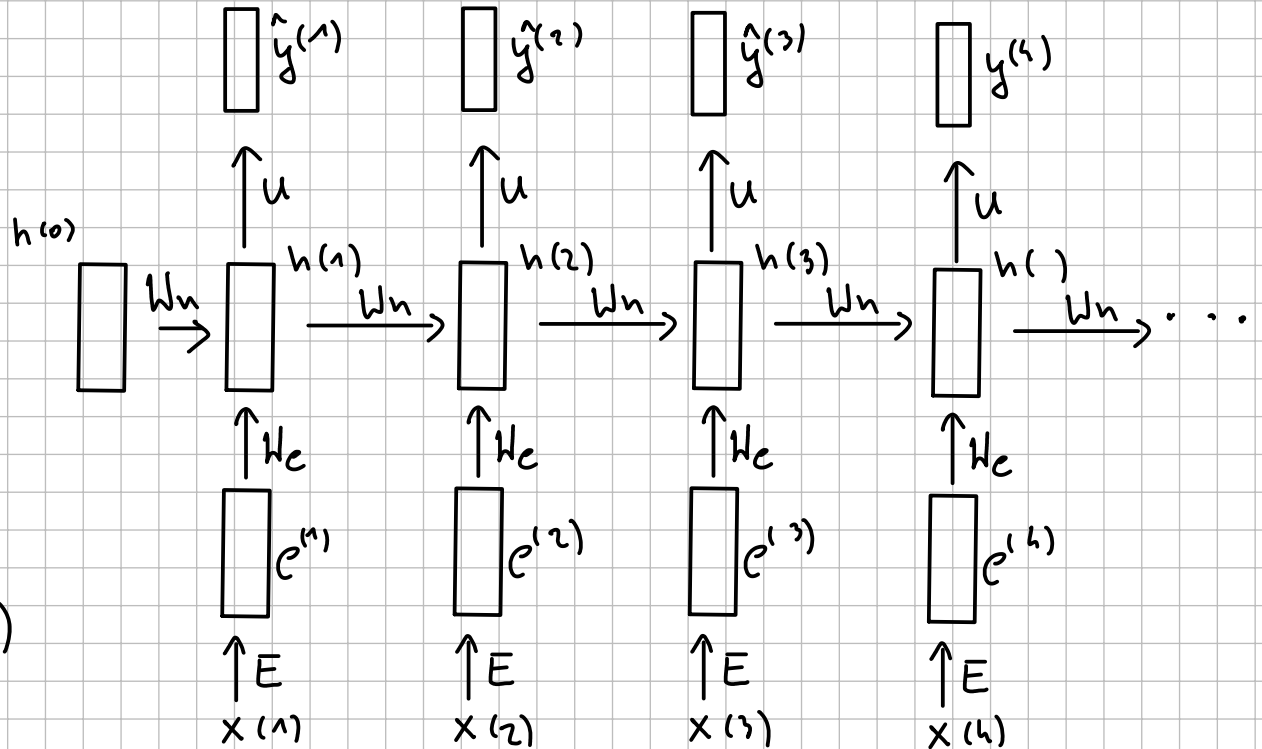
$h(0)$ - stan początkowy,

najczęściej wektor zer

$$h(i) = \sigma(W_h h(i-1) + W_e e(i) + b_1)$$

stany wyjściowe

$$y(i) = \text{softmax}(U h(i) + b_2)$$



Zalety i wady RNN

Zalety:

- Wejście może być dowolnie długie
- W teorii w czasie t może mieć informacje z wielu kroków wcześniej
- Ta sama macierz W używana w każdym kroku
 - ilość parametrów do trenowania nie zależy od długości wejścia
 - stałość tego jak traktujemy wejście

Wady:

- Część rekurencyjna sprawia że obliczenia są długie
- W praktyce w czasie t model "nie pamięta" zbyt dobrze dalekiej przeszłości

Trenowanie RNN

- Duże wejście \rightarrow np. książka, tekst z wikipedii
- Przepuszczamy przez RNN i obliczamy $\hat{y}(i)$ dla każdego kroku
- Funkcja kosztu - Entropia krzyżowa

Dla jednego kroku:

$$J(t)(\theta) = CE(y(i) - \hat{y}(i)) = - \sum_{w \in V} y_w(t) \log \hat{y}_w(t) = - \log \hat{y}_{x(t+1)}(t)$$

\uparrow one-hot

Całokłowy koszt:

$$J(\theta) = \frac{1}{T} \sum_{t=0}^T J(t)(\theta) = -\frac{1}{T} \sum_{t=0}^T \log \hat{y}_{x(t+1)}(t)$$

W praktyce nie przechodzimy przez cały tekst na raz, tylko obliczamy $J(\theta)$

dla mniejszego podzbioru (np 16 zdań), obliczamy gradienty

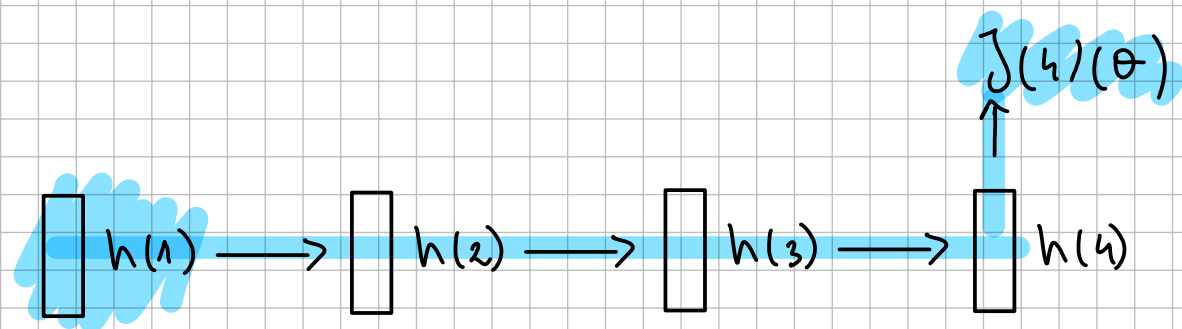
aktualizujemy wagi i kontynuujemy

Backpropagacja

$$\frac{\partial J(t)}{\partial W_n} = \sum_{i=1}^t \frac{\partial J(t)}{\partial W_n} \Big|_i \quad \leftarrow \text{"Backpropagation through time"}$$

- Problem \rightarrow jak zdanie / tekst jest zbyt długi to trwałoby to długo!
- rozwiązanie \rightarrow "ucięta" backpropagacja w czasie, przechodzimy w tył nie przez całe zdanie, ale przez z góry ustaloną liczbę kroków, np 20

Problem znikającego/wybuchającego gradientu



Obliczamy gradient $J(4)(\theta)$ względem $h(1)$, przy pomocy metody Łańcucha:

$$\frac{\partial J(4)(\theta)}{\partial h(1)} = \frac{\partial h(2)}{\partial h(1)} \cdot \frac{\partial h(3)}{\partial h(2)} \cdot \frac{\partial h(4)}{\partial h(3)} \cdot \frac{\partial J(4)(\theta)}{\partial h(4)}$$

→ wielokrotne mnożenie → małe liczby znikają do 0

→ duże liczby eksplodują

Przykład dla $\sigma(x) = x$ (funkcja identycznościowa)

$$h(t) = \sigma(W_n h(t-1) + W_x \dot{x}(t) + b_1)$$

$$\frac{\partial h(t)}{\partial h(t-1)} = \text{diag}(\sigma'(W_n h(t-1) + W_x \dot{x}(t) + b_1)) W_n = I W_n = W_n$$

⇓

$$\frac{\partial J(i)}{\partial h(j)} = \frac{\partial J(i)}{\partial h(i)} \cdot \prod_{j < t \leq i} \frac{\partial h(t)}{\partial h(t-1)} = \frac{\partial J(i)}{\partial h(j)} \cdot (W_n)^l \quad [i-j=l]$$

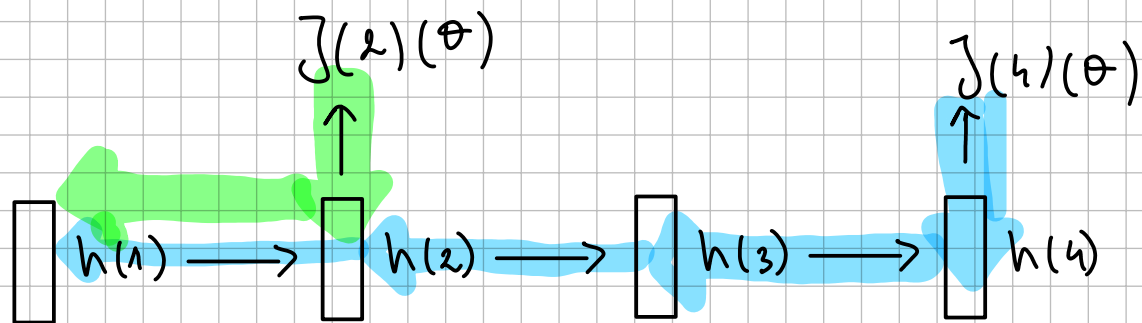
↑ może mieć wykładniczo

dla W_n te wartości własne $\lambda_1, \dots, \lambda_n < 1$: [warunek dostateczny ale nie konieczny]

$$\frac{\partial J(i)}{\partial h(j)} (W_n)^l = \sum_{i=1}^n c_i (\lambda_i)^l q_i \approx 0 \text{ dla dużych } l$$

[Dla bardziej skomplikowanych σ dzieje się to samo tylko dla $\lambda_i < 1$]
które zależą od tego jakie wybieramy σ

Skutki znikającego gradientu



- Wpływu odległego sygnału znika bo rzędy wielkości gradientów mogą być zupełnie różne
- Model nie potrafi "pamiętać" informacji z dalekiej przeszłości

Przykład:

"Kiedy Asia chciała wydrukować bilet, okazało się że nie ma tuszu w drukarce. Asia poszła do sklepu kupić toner, był akurat na promocje.

Wzięła nowy toner do drukarki i wydrukowała _____"

25 kroków różnicy! → znikający gradient sprawia, że wptyw

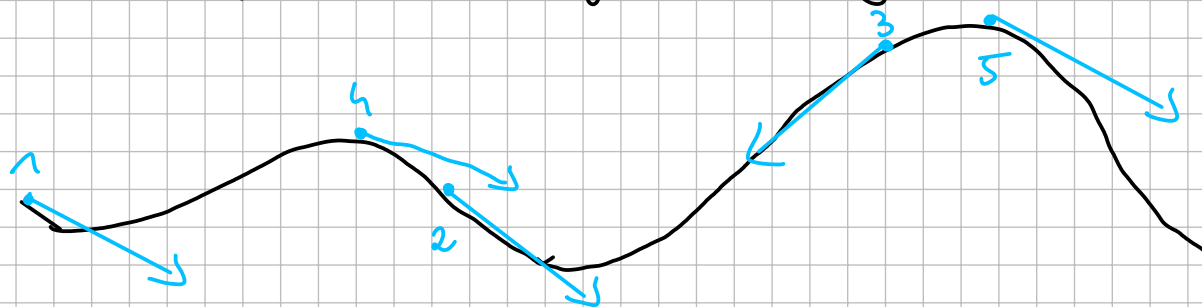
słowa "bilet" z tekstu dawno jest znikomy i nie potrafimy

"zapamiętać" informacji z dalekiej przeszłości

Skutki eksplodującego gradientu

$$\theta^{\text{new}} = \theta^{\text{old}} + \alpha \underbrace{\nabla_{\theta} J(\theta)}_{\text{gradient}}$$

- jak gradient jest zbyt duży, to robimy zbyt duże kroki i błędzimy



Rozwiązanie problemu eksplodującego gradientu

- "Zmniejszenie kroku przy zachowaniu kierunku chodzenia"

Pseudokod:

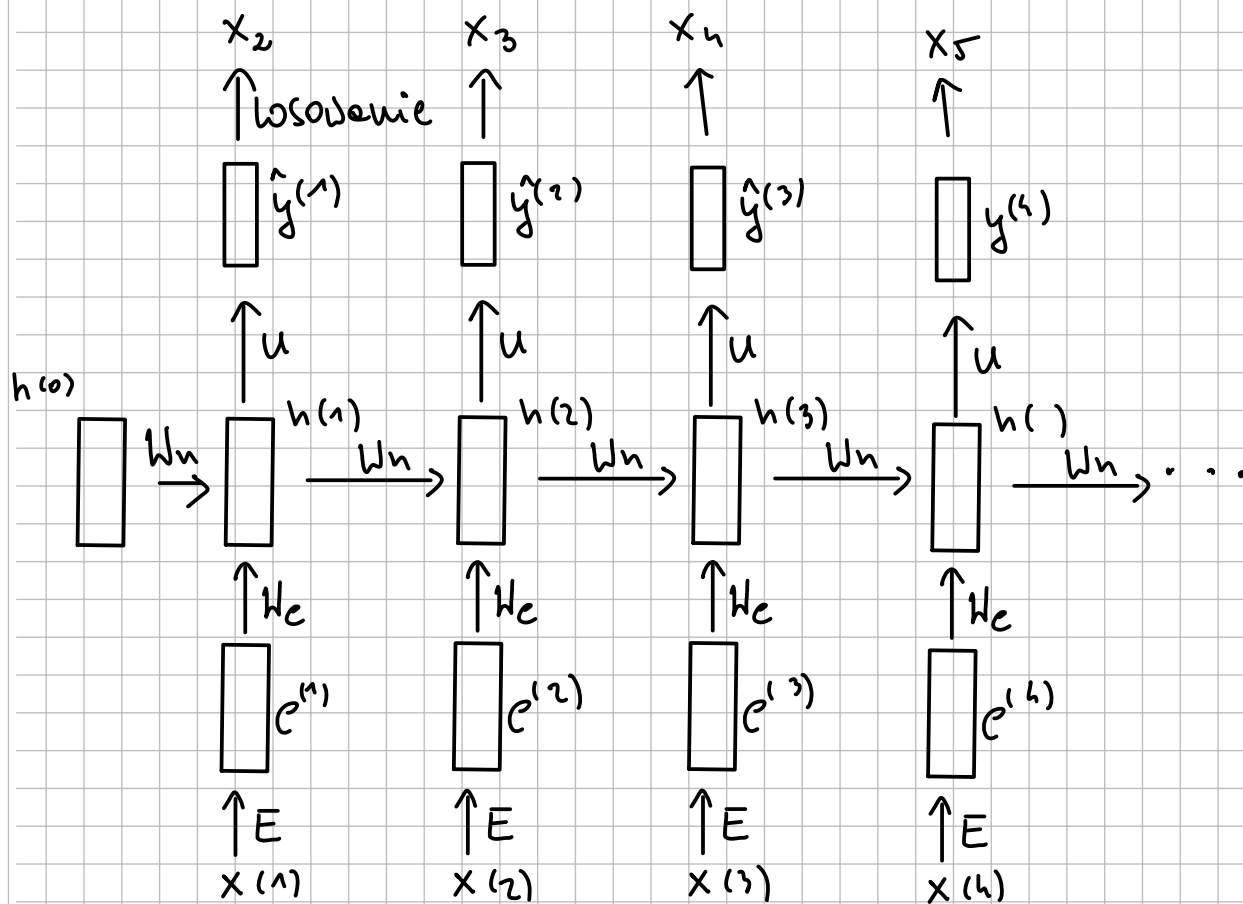
$$\cdot \hat{g} \leftarrow \frac{\partial \epsilon}{\partial \theta}$$

· if $\|\hat{g}\| > \tau$ then

$$\hat{g} \leftarrow \frac{\tau}{\|\hat{g}\|} \hat{g}$$

end if

Generowanie tekstu przez RNN



Wartości wylosowane
w 1. kroku traktujemy
jako 2. element wejścia
itd.