**Ludwig-Maximilians-Universität München**
<span></span>Munich, 18.11.2024
**Institut für Informatik**
Prof. Dr. Matthias Schubert
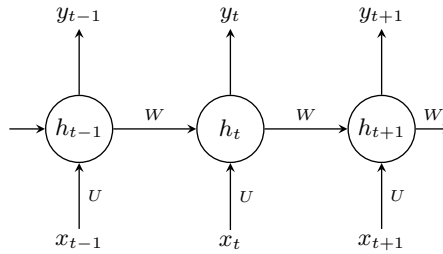Maximilian Bernhard, Niklas Strauß

## Deep Learning and Artificial Intelligence
WS 2024/25

## Exercise 5: Recurrent Neural Networks

### Exercise 5-1 Backpropagation through Time

Consider the following RNN:



Each state $h_t$ is given by:

$$h_t = \sigma(W h_{t-1} + U x_t), \qquad \sigma(z) = \frac{1}{1 + e^{-z}}$$

Let $L$ be a loss function defined as the sum over the losses $L_t$ at every time step until time $T$: $L = \sum_{t=0}^{T} L_t$, where $L_t$ is a scalar loss depending on $h_t$.

In the following, we want to derive the gradient of this loss function with respect to the parameter $W$.

(a) Given $\mathbf{y} = \sigma(W\mathbf{x})$ where $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{x} \in \mathbb{R}^d$ and $W \in R^{n \times d}$. Derive the Jacobian $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$!

$$\left[\text{Solution: } \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \text{diag}(\sigma'(\cdot))\, W\right]$$

(b) Derive the quantity $\frac{\partial L}{\partial W} = \sum_{t=0}^{T} \sum_{k=0}^{t} \frac{\partial L_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W}$!

### Exercise 5-2 Vanishing/Exploding Gradients in RNNs

In this exercise, we want to understand why RNNs are especially prone to the Vannishing/Exploding Gradients problem and what role the eigenvalues of the weight matrix play. Consider part b) of exercise 5-1 again.

(a) Write down $\frac{\partial L}{\partial W}$ as expanded sum for $T = 3$. You should see that if we want to backpropagate through $n$ timesteps, we have to multiply the matrix $diag(\sigma')W$ n times with itself.

(b) Remember that any diagonalizable (square) matrix $M$ can be represented by its eigendecomposition $M = Q\Lambda Q^{-1}$ where $Q$ is a matrix whose $i$-th column corresponds to the $i$-th eigenvector of $M$ and $\Lambda$ is a diagonal matrix with the corresponding eigenvalues placed on the diagonals.[1].

Proof by induction that for such a matrix the product $\prod_{i=1}^{n} M$ can be written as: $M^n = Q\Lambda^n Q^{-1}$!

---

[1]Every eigenvector $v_i$ satisfies the linear equation $M v_i = \lambda_i v_i$ where $\lambda_i = \Lambda_{ii}$
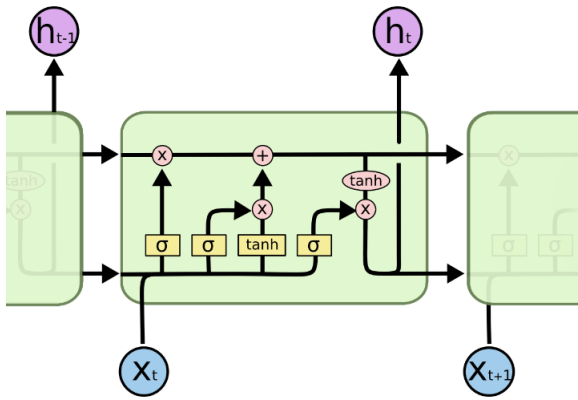
(c) Consider the weight matrix $W = \begin{pmatrix} 0.58 & 0.24 \\ 0.24 & 0.72 \end{pmatrix}$. Its eigendecomposition is:

$$W = Q\Lambda\,Q^{-1} = \begin{pmatrix} 0.8 & 0.6 \\ -0.6 & 0.8 \end{pmatrix} \begin{pmatrix} 0.4 & 0 \\ 0 & 0.9 \end{pmatrix} \begin{pmatrix} 0.8 & -0.6 \\ 0.6 & 0.8 \end{pmatrix}.$$

Calculate $W^{30}$! What do you observe? What happens in general if the absolute value of all eigenvalues of $W$ is smaller than 1? What happens if the absolute value of any eigenvalue of $W$ is larger than 1? What if all eigenvalues are 1?

## Exercise 5-3    LSTMs

Recall the elements of a module in an LSTM and the corresponding computations, where $\odot$ stands for pointwise multiplication. [2]



$$f_t = \sigma(W_f h_{t-1} + U_f x_t)$$
$$i_t = \sigma(W_i h_{t-1} + U_i x_t)$$
$$o_t = \sigma(W_o h_{t-1} + U_o x_t)$$
$$\tilde{C}_t = \tanh(W_c h_{t-1} + U_c x_t)$$
$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$$
$$h_t = o_t \odot \tanh(C_t)$$

(a) What do the gates $f_t$, $i_t$ and $o_t$ do?

(b) Which of the quantities next to the figure are always positive?

Let's now try to understand how this architecture approaches the vanishing gradients problem. To calculate the gradient $\frac{\partial L}{\partial \theta}$, where $\theta$ stands for the parameters $(W_f, W_o, W_i, W_c)$, we now have to consider the cell state $C_t$ instead of $h_t$. Like $h_t$ in normal RNNs, $C_t$ will also depend on the previous cell states $C_{t-1}, ...C_0$, so we get a formula of the form:

$$\frac{\partial L}{\partial W} = \sum_{t=0}^{T} \sum_{k=1}^{t} \frac{\partial L}{\partial C_t} \frac{\partial C_t}{\partial C_k} \frac{\partial C_k}{\partial W}. \quad [3]$$

(c) We know that $\frac{\partial C_t}{\partial C_k} = \prod_{i=k+1}^{t} \frac{\partial C_t}{\partial C_{t-1}}$. Let $f_t = 1$ and $i_t = 0$ such that $C_t = C_{t-1}$ for all $t$.

What is the gradient $\frac{\partial C_t}{\partial C_k}$ in this case?

---

[2]For a good explanation on LSTMs you can refer to http://colah.github.io/posts/2015-08-Understanding-LSTMs/

[3]The real formula is a bit more complicated since $C_t$ also depends on $f_t$, $i_t$ and $\tilde{C}_t$, which in turn all depend on $W$, but this can be neglected.