Prof. Dr. Andreas Hotho, Albin Zehe, Jonas Kaiser
Data Science Chair

08.05.2024

# 2. Assignment in "Machine Learning for Natural Language Processing"

Summer Term 2024

# 1 General Questions

1. Name three common activation functions and their derivatives!

2. Given two word embeddings $E$ and $E'$ for a corpus of words, how could you rate the quality of the embeddings relative to one another? Is there a possibility to give an absolute, global rating for word embeddings?

   ? **Something to think about**
3. Suggest a combination of GloVe and FastText!

# 2 Neural Networks

## Bias Trick

Formally, a feed-forward layer in a neural network is defined by a number of parameters consisting of a weight matrix $W$ and a bias vector $b$. The output $y$ that is passed to the activation function is then calculated as

$$y = Wx + b.$$

In practice, we want to keep the number of vectorised operations as small as possible, to enable maximum training/prediction speed (remember that GPUs can parallelise the matrix operations). Therefore, these two parameter sets are often combined into one

matrix

$$W' = W\|b = \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1m} & b_1 \\ w_{21} & w_{22} & \cdots & w_{2m} & b_2 \\ \vdots & & \ddots & & \vdots \\ w_{n1} & w_{n2} & \cdots & w_{nm} & b_n \end{pmatrix}$$

1. How does the input $x$ need to be modified to enable computing the output $y$ in a single operation?

2. How is the output $y$ calculated using the modified $W'$ and input $x'$?

3. Proof that this always yields the same result as $Wx + b$!

# 3 Python

In the lecture, some commonly used optimisers for neural networks were introduced. Implement functions for

1. `sgd_update`

2. `nesterov_momentum_update`

3. `adam_update`

in Python. Each function should

- take as input

  - the current value of *one* input parameter,

  - a function returning the gradient regarding this parameter, and

  - the necessary hyper-parameters such as the learning rate, and

- return the new value for this parameter after the update.

You do *not* have to implement backpropagation in this assignment!