**Technische Universität Hamburg**
**Institut für Mathematik**
**Lehrstuhl Numerische Mathematik**

**Jens-Peter M. Zemke**
**Jonas Grams**

## Mathematics of Neural Networks
## winter semester 2021/2022
## exercise sheet 2, solutions

**Exercise 1:** (2 points) Check $\mathsf{SoftAbs}_k(x)$ for monotonicity, symmetry, and asymptotic behavior and calculate the derivative.

**Solution:** Symmetry is shown by

$$\mathsf{SoftAbs}_k(x) = \frac{kx^2}{1+|kx|} = \frac{k(-x)^2}{1+|k(-x)|} = \mathsf{SoftAbs}_k(-x).$$

We further see that for positive $x$

$$\mathsf{SoftAbs}_k(x) = \frac{kx^2}{1+kx} = \frac{1}{k(kx+1)} - \frac{1}{k} + x$$

and for negative $x$

$$\mathsf{SoftAbs}_k(x) = \frac{kx^2}{1-kx} = -\frac{1}{k(kx-1)} - \frac{1}{k} - x.$$

Thus

$$a_+(x) = -\frac{1}{k} + x$$

is an oblique asymptote to $\mathsf{SoftAbs}_k(x)$ when $x$ tends to $+\infty$ and

$$a_-(x) = -\frac{1}{k} - x$$

is an oblique asymptote to $\mathsf{SoftAbs}_k(x)$ when $x$ tends to $-\infty$.
We calculate the derivative:

$$\frac{d}{dx}\mathsf{SoftAbs}_k(x) = \frac{d}{dx}\left(x \cdot \mathsf{SoftSign}_k(x)\right) = \frac{kx}{1+|kx|} + \frac{kx}{(1+|kx|)^2} = \frac{(1+|kx|)\cdot kx + kx}{(1+|kx|)^2}$$
$$= \frac{(2+|kx|)}{(1+|kx|)^2} \cdot kx$$

Hence $\mathsf{SoftAbs}_k(x)$ is strictly increasing for $x > 0$ and strictly decreasing for $x < 0$.

**Exercise 2:** (3 points)  Prove that any continous function

$$f : [a, b] \rightarrow [0, \infty)$$

can be approximated arbitrarily well by a neural network $N : \mathbb{R} \rightarrow \mathbb{R}$ with a single hidden layer comprising $n$ neurons for some $n \in \mathbb{N}$ based on ReLU activations.

Hint: Can you express the linear interpolant through the $n$ equidistant points

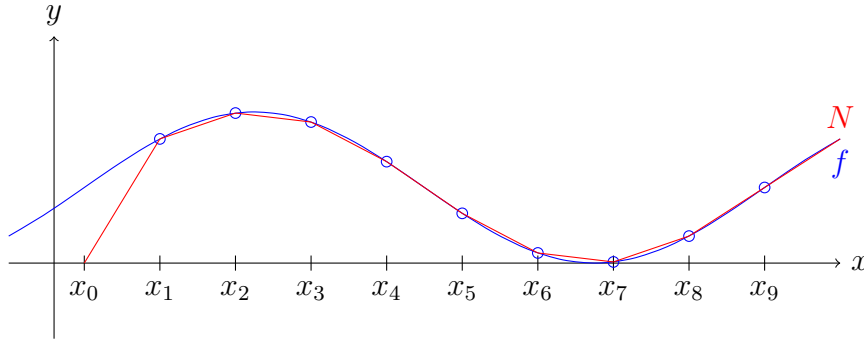$$a = x_1, x_2, \ldots, x_n = b \quad \text{and function values} \quad f_i = f(x_i), \quad i = 1, \ldots, n$$

using linear combinations of ReLU functions $\max(0, x + b_i)$ with biases $b_1 = -(a - h), b_2 = -a, b_3 = -(a + h), b_4 = -(a + 2h), \ldots, b_n = -(b - h)$, where $h = x_2 - x_1$?

**Solution:**   First let $n \in \mathbb{N}$, and $a = x_1, x_2, \ldots, x_n = b$ be equidistant points. These points devide the interval $[a, b]$ into $n - 1$ equal intervals $[x_i, x_{i+1}]$, $i = 1, \ldots n - 1$, each of length $h = \frac{1}{n-1}$. To construct a neural network with one hidden layer comprising $n$ neurons based on ReLU activations, we have to find the weights and biases for the network. Note that for $m, b \in \mathbb{R}$ we get

$$\mathsf{ReLU}(mx + b) = \max(0, mx + b) = m\max(0, x + b/m) = m\,\mathsf{ReLU}(x + b/m).$$

Thus we can set the weights of the hidden layer to 1, and only need to determine the weights for the output layer.

Now consider the following figure:



As described in the hint, we pick the biases of the hidden layer as

$$b_1 = -(a - h), b2 = -x_1 = -a, b_3 = -x_2 = -(a + h), \ldots, b_n = -x_{n-1} = -(a + (n - 2)h),$$

to possibly change the gradient after each intervall $[x_i, x_i + 1]$. Thus prior to $-b_1$ all ReLU functions $\max(0, x + b_i)$ are zero. Between $-b_1 = a - h$ and $-b_2 = a = x_1$ only one of these is non-zero, namely

$$\max(0, x + b_1) = x + b_1.$$

To achieve a possibly non-zero value $f(a) = f(x_1)$ at $x_1$ we choose a scalar multiple $m_1$ with

$$f(x_1) = m_1(x_1 + b_1) = hm_1 \quad \Rightarrow \quad m_1 = f(x_1)/h.$$

The function $m_1 \max(0, x + b_1)$ interpolates $f$ at $x_1$. Next we consider the point $x_2 = a + h$. To achieve $N(x_2) = f(x_2)$, we need to adjust the gradient on the interval $[x_1, x_2]$, and thus find a second scalar multiple $m_2$ with

$$f(x_2) = m_1(x_2 + b_1) + m_2(x_2 + b2) = 2hm_1 + hm_2.$$

Generally for $i \in \{1, \ldots, n\}$ we get

$$
\begin{aligned}
f(x_i) &= \sum_{j=1}^{i} m_j(x_i + b_j) \\
&= \sum_{j=1}^{i} m_j(a + (i-1)h - (a + (j-2)h)) \\
&= \sum_{j=1}^{i} m_j(i - j + 1)h.
\end{aligned}
$$

From these equations, we can derive the following uniquely solvable lower triangular system.

$$
h \begin{pmatrix} 1 & & & \\ 2 & 1 & & \\ \vdots & \ddots & \ddots & \\ n & \cdots & 2 & 1 \end{pmatrix} \begin{pmatrix} m_1 \\ m_2 \\ \vdots \\ m_n \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{pmatrix}.
$$

The solution of this system determines the weights for the output layer. The bias for the output layer can be set to 0. If we let $n \to \infty$ in the equidistant formulation, the ReLU-approximations will converge by continuity to the function $f$, since continuity at $x$ is given by

$$
\forall\, \epsilon > 0 \, \exists\, \delta > 0 \, : \, |f(x) - f(y)| < \epsilon \,\forall\, |x - y| < \delta.
$$

To show the convergence, let $\epsilon > 0$. We have to proof that we can find a neural network $N$ with

$$
|f(x) - N(x)| < \epsilon \quad \text{for all } x \in [a, b].
$$

Due to the continuity of $f$, we find a $\delta > 0$ with $|f(x) - f(y)| < \epsilon/2$ for all $|x - y| < \delta$.
Let $n \in \mathbb{N}$ with $h = \frac{1}{n-1} < \delta$. Let $N$ be a neural network as constructed above and let $i \in \{1, \ldots, n-1\}$.
Since $N$ is a linear polynomial with gradient $\frac{f(x_{i+1}) - f(x_i)}{h}$ on the interval $[x_i, x_{i+1}]$, we get

$$
N(x) = f(x_i) + \frac{f(x_{i+1}) - f(x_i)}{h}(x - x_i) \quad \text{for all } x \in [x_i, x_{i+1}].
$$

This indicates

$$
\begin{aligned}
|f(x) - N(x)| &= |f(x) - f(x_i) - \frac{f(x_{i+1}) - f(x_i)}{h}(x - x_i)| \\
&\leq |f(x) - f(x_i)| + |f(x_{i+1}) - f(x_i)| \frac{|x - x_i|}{h} \\
&\leq |f(x) - f(x_i)| + |f(x_{i+1}) - f(x_i)| < \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon
\end{aligned}
$$

for all $x \in [x_i, x_{i+1}]$. Since $i$ was chosen arbitrary, we get

$$
|f(x) - N(x)| < \epsilon \quad \text{for all } x \in [a, b].
$$

**Exercise 3:** (3+3 points)

   a) Add backpropagation to our feedforward neural network. You can use the given templates.

- Add an attribute `data` to our `ReLU` class that stores data for the backpropagation from the last evaluation.
- Add a method `backprop(self, delta)` to our `ReLU` class that performs the backpropagation $a'(\mathbf{z}) \circ \boldsymbol{\delta}$, where $\mathbf{z}$ can be recovered from `self.data`.
- Add attributes `__a` and `__z`, both initialised with `None`, to our `DenseLayer` class which store the input and the affine linear combination $\mathbf{Wa} + \mathbf{b}$ from the last evaluation.
- Add attributes `dW` and `db`, initialised with zeros, to our `DenseLayer` class that store the updates for the weights and biases.
- Add the method `backprop(self, delta)` to our `DenseLayer` class that computes the $\boldsymbol{\delta}$ for this layer, the derivatives of the cost function with respect to the weights and the biases of this layer, and returns $\mathbf{W}^{\mathsf{T}}\boldsymbol{\delta}$ for the previous layer.
- Implement the method `backprop(self, x, y)` in our `SequentialNet` class that computes and returns the derivatives of the cost function with respect to the weights and the biases for the training data `x` and `y` corresponding to one minibatch.

   b) Implement `train(self, x, y, batch_size=16, epochs=10)` that uses repeated calls to `backprop()` to carry out `epochs` epochs of SGD with batch size `batch_size` for the training data `x` and `y` corresponding to the whole training set. For this

- implement a class `SGD` which is initialised with a learning rate `eta` and has a method `update(self, data, ddata)` that performs a SGD update step on the data in the list `data` with update data in the list `ddata`,
- add an attribute `optim` to your `DenseLayer` class, which is is given on initialization and an instance of SGD by default,
- add a method `update(self)` to your `DenseLayer` class that updates the weights and biases for a layer using the `update` method from the optimizer given in `optim`.

**Solution:** See the following listings for a possible implementation.

```
 1  import numpy as np
 2
 3  class ReLU ():
 4
 5      def __init__ (self):
 6
 7          self.data = None
 8          self.name = 'ReLU'
 9
10      def wfun (self, m, n):
11          return np.random.randn(m,n)*np.sqrt(4/(m + n))
12
13      def evaluate (self, x):
14
15          self.data = x
16          return x.clip(min = 0)
17
18      def backprop (self, delta):
19
20          return (self.data >= 0) * delta
21
22  class Abs ():
```

```
23
24      def __init__(self):
25
26          self.data  = None
27          self.name  = 'Abs'
28
29      def wfun(self, m, n):
30          return np.random.randn(m,n)*np.sqrt(1/(m + n))
31
32      def evaluate(self, x):
33
34          self.data = x
35          return np.abs(x)
36
37      def backprop(self, delta):
38
39          return np.sign(self.data)*delta
```

```
1  import numpy as np
2
3  from activations import ReLU
4  from optimizers  import SGD
5
6  class DenseLayer:
7
8      def __init__(self,
9                     ni, # Number of inputs
10                    no, # Number of outputs
11                    afun = None, # Activationfunction for the layer
12                    optim = None
13      ):
14
15          self.ni   = ni
16          self.no   = no
17
18          if afun is None:
19              self.afun = ReLU()
20          else:
21              self.afun = afun
22
23          if optim is None:
24              self.optim = SGD()
25          else:
26              self.optim = optim
27
28          self.W    = self.afun.wfun(no, ni)
29          self.b    = np.zeros((no,1))
30          self.__z  = None
31          self.__a  = None
32          self.dW   = np.zeros_like(self.W)
33          self.db   = np.zeros_like(self.b)
34
35      def evaluate(self, a):
36
37          self.__a = a
38          self.__z = self.W @ self.__a + self.b
```

```
39            return self.afun.evaluate(self.__z)
40
41      def set_weights(self, W):
42
43            assert(W.shape == (self.no, self.ni))
44            self.W = W
45
46      def set_bias(self, b):
47
48            assert(b.size == self.no)
49            self.b = b
50
51      def backprop(self, delta):
52
53            delta   = self.afun.backprop(delta)
54            self.dW = delta @ self.__a.T
55            self.db = delta @ np.ones((delta.shape[1],1))
56
57            return self.W.T @ delta
58
59      def update(self):
60
61            self.optim.update([self.W, self.b],
62                              [self.dW, self.db])
```

```
1  import numpy as np
2
3  class SGD:
4
5      def __init__(self, eta=.1):
6
7            self.eta = eta
8
9      def update(self, data, ddata):
10
11           for p, dp in zip(data, ddata):
12
13               p -= self.eta * dp
```

```
1      def backprop(self, x, y):
2
3            delta = (self.evaluate(x) - y)/y.shape[1]
4
5            for layer in reversed(self.layers):
6
7                delta = layer.backprop(delta)
8
9      def train(self, x, y, batch_size=16, epochs=10):
10
11           n_data    = y.shape[1]
12           n_batches = int(np.ceil(n_data/batch_size))
13
14           for e in range(epochs):
15
16               p = np.random.permutation(n_data)
17               for j in range(n_batches):
```

```
18
19                      self.backprop(x[:,p[j*batch_size:(j+1)*batch_size]],
20                                     y[:,p[j*batch_size:(j+1)*batch_size]])
21
22              for layer in self.layers:
23                  layer.update()
```

**Exercise 4:** (2+2 points)  In multiclass classification the function $\mathsf{SoftMax} : \mathbb{R}^n \to (0,1)^n$,

$$\mathbf{p} := \mathsf{SoftMax}(\mathbf{x}) := \frac{e^{\mathbf{x}}}{\mathbf{e}^{\top} e^{\mathbf{x}}} = \begin{pmatrix} e^{x_1} \\ \vdots \\ e^{x_n} \end{pmatrix} / \left( \sum_{j=1}^{n} e^{x_j} \right),$$

which returns a probability distribution, is used as last layer together with the Cross Entropy Loss function

$$H(\mathbf{y}, \mathbf{p}) := -\sum_{j=1}^{n} y_j \ln(p_j)$$

as cost function. In training, the vector $\mathbf{y}$ will be a unit vector, $\mathbf{p}$ will be the outcome of the $\mathsf{SoftMax}$-layer. This is known as categorical cross entropy loss. Compute the derivative of the
   a) $\mathsf{SoftMax}$ function $\mathbf{p} = \mathsf{SoftMax}(\mathbf{x})$ with respect to $\mathbf{x}$,
   b) Cross Entropy Loss function combined with the $\mathsf{SoftMax}$ layer, $H(\mathbf{y}, \mathsf{SoftMax}(\mathbf{x}))$, with respect to $\mathbf{x}$.
How do you initialize backpropagation in this case?

**Solution:**
   a) To compute the Jacobi matrix

$$\frac{\partial \mathsf{SoftMax}(\mathbf{x})}{\partial \mathbf{x}}$$

of the $\mathsf{SoftMax}$ function we take a look at the elements of the Jacobian and use the quotient rule,

$$\frac{\partial \frac{e^{x_j}}{\sum_{k=1}^n e^{x_k}}}{\partial x_i} = \begin{cases} \frac{e^{x_i}\left(\sum_{k=1}^n e^{x_k}\right) - e^{x_i} e^{x_i}}{\left(\sum_{k=1}^n e^{x_k}\right)^2} = \frac{e^{x_i}}{\sum_{k=1}^n e^{x_k}} \cdot \frac{\sum_{k=1}^n e^{x_k} - e^{x_i}}{\sum_{k=1}^n e^{x_k}} = p_i(1 - p_i), & j = i, \\ \frac{-e^{x_j} e^{x_i}}{\left(\sum_{k=1}^n e^{x_k}\right)^2} = \frac{e^{x_j}}{\sum_{k=1}^n e^{x_k}} \cdot \frac{-e^{x_i}}{\sum_{k=1}^n e^{x_k}} = p_j(0 - p_i), & j \neq i. \end{cases}$$

Thus,

$$\frac{\partial \frac{e^{x_j}}{\sum_{k=1}^n e^{x_k}}}{\partial x_i} = p_j(\delta_{j,i} - p_i).$$

   b) The derivative of the Cross Entropy Loss function is given by

$$\frac{\partial H(\mathbf{y}, \mathbf{p})}{\partial x_i} = \frac{\partial - \sum_{j=1}^n y_j \ln(p_j)}{\partial x_i} = -\sum_{j=1}^{n} y_j \frac{\partial \ln(p_j)}{\partial x_i} = -\sum_{j=1}^{n} y_j \frac{\partial \ln(p_j)}{\partial p_j} \frac{\partial p_j}{\partial x_i}$$

$$= -\sum_{j=1}^{n} \frac{y_j}{p_j} p_j(\delta_{j,i} - p_i) = -y_i + p_i \sum_{j=1}^{n} y_j = p_i - y_i,$$

since $\mathbf{y}$ represents a probability distribution and thus the sum of all elements gives one. In vector form, the gradient is given by

$$\frac{\partial H(\mathbf{y}, \mathbf{p})}{\partial \mathbf{x}} = \mathbf{p} - \mathbf{y}.$$

To initialize the backpropagation we simply have to set $\boldsymbol{\delta}_{m-1} = \mathbf{p} - \mathbf{y}$. This is in contrast to the mean square error $\|\mathbf{a}_m - \mathbf{y}\|_2^2/2$, where we additionally had to multiply $\mathbf{a}_m - \mathbf{y}$ elementwise with the derivative of the activation function to obtain $\boldsymbol{\delta}_{m-1}$.

**Exercise 5:** (1+1+2 points)   Test your class `SequentialNet` or an implementation using TensorFlow/Keras with
   a) the data from exercise 1b) of the first exercise sheet with one hidden layer comprising 3, 4, 10, or 100 neurons and ReLU activation functions,
   b) the data given by

```
x = np.expand_dims(np.linspace(0, np.pi/2, 2000),0)
y = np.concatenate((np.sin(x), np.cos(x)))/2+.5
x_train, y_train = x[:,::2], y[:,::2]
x_test, y_test = x[:,1::2], y[:,1::2]
```

   with one hidden layer of 100 neurons and ReLU activation functions,
   c) MNIST with input layer of size 784, a hidden layer of size 100, and an output layer of size 10 with ReLU activation functions, either with the squared error cost function or a SoftMax layer followed by cross entropy loss. This data set is provided to you in the file `mnist.npz`. Further information can be found here:

   http://yann.lecun.com/exdb/mnist/

   Import the data set using `DATA = np.load('mnist.npz')` function followed by `x_train, y_train = DATA['x_train'], DATA['y_train']` and `x_test, y_test = DATA['x_test'], DATA['y_test']`.
   Note that you need to reshape the data before you can train your neural net.
If you want to use TF/Keras you can find a description on

   https://keras.io/getting-started/sequential-model-guide/

and a code skeleton in `TF_skeleton.py`.

**Solution:**   The hand-made solution is based on our Python class `SequentialNet`, the solution based on TF/Keras is based on the description mentioned in the exercise. The main difference is that TF/Keras is based on the *transposed* input and output. Additionally the mean squared error is used as a loss function, so compared to our cost function it is divided by the number of predictions. Because we didn't further specify it weights and biases in the TF/Keras version they are initialised using a uniform distribution. In the neural net based on our class a standard normal distribution is used.
   a) The implementation of the four neural nets of part a) is given in the following listing for our class `SequentialNet`, where we included the case of the absolute value as activation function,

```
1  import sys
2  import numpy                as np
3  import matplotlib.pyplot as plt
4
5  from networks     import SequentialNet
6  from activations import Abs
7  from optimizers   import SGD
8  from layers       import DenseLayer
```

```
 9
10  ###################################################
11  # Aufgabe 6a)                                     #
12  ###################################################
13
14  x = np.array([[0, 0, 1, 1],
15                [0, 1, 0, 1]])
16
17  y = np.array([[0, 0, 0, 1],
18                [0, 1, 1, 1],
19                [0, 1, 1, 0]])
20
21  bs, ep, eta = 1, 5000, .1
22
23  #######################################
24  # Tests with ReLU activation funtion  #
25  #######################################
26
27  layers = [DenseLayer(2, 3, optim=SGD(eta)),
28            DenseLayer(3, 3, optim=SGD(eta))]
29  netz = SequentialNet(layers)
30  netz.train(x, y, bs, ep)
31
32  ytilde = netz.evaluate(x)
33
34  print('y␣=\n', y)
35  print('ytilde␣=\n', ytilde)
36  print('error␣=', np.linalg.norm(y - ytilde))
37  print('n␣=␣3,␣ReLU')
38
39  layers = [DenseLayer(2, 4, optim=SGD(eta)),
40            DenseLayer(4, 3, optim=SGD(eta))]
41  netz = SequentialNet(2, layers)
42  netz.train(x, y, bs, ep)
43
44  ytilde = netz.evaluate(x)
45
46  print('y␣=\n', y)
47  print('ytilde␣=\n', ytilde)
48  print('error␣=', np.linalg.norm(y - ytilde))
49  print('n␣=␣4,␣ReLU')
50
51  layers = [DenseLayer(2, 10, optim=SGD(eta)),
52            DenseLayer(10, 3, optim=SGD(eta))]
53  netz = SequentialNet(2, layers)
54  netz.train(x, y, bs, ep)
55
56  ytilde = netz.evaluate(x)
57
58  print('y␣=\n', y)
59  print('ytilde␣=\n', ytilde)
60  print('error␣=', np.linalg.norm(y - ytilde))
61  print('n␣=␣10,␣ReLU')
62
63  layers = [DenseLayer(2, 100, optim=SGD(eta)),
64            DenseLayer(100, 3, optim=SGD(eta))]
```

```
65  netz = SequentialNet (2, layers)
66
67  print ('y␣=\n', y)
68  print ('ytilde␣=\n', ytilde)
69  print ('error␣=', np.linalg.norm (y - ytilde))
70  print ('n␣=␣100,␣ReLU')
71
72  #####################################
73  # Tests with Abs activation funtion #
74  #####################################
75
76  layers = [DenseLayer (2, 3, afun=Abs (), optim=SGD (eta)),
77            DenseLayer (3, 3, afun=Abs (), optim=SGD (eta))]
78  netz = SequentialNet (2, layers)
79  netz.train (x, y, bs, ep)
80
81
82  ytilde = netz.evaluate (x)
83
84  print ('y␣=\n', y)
85  print ('ytilde␣=\n', ytilde)
86  print ('error␣=', np.linalg.norm (y - ytilde))
87  print ('n␣=␣3,␣Abs')
88
89  layers = [DenseLayer (2, 4, afun=Abs (), optim=SGD (eta)),
90            DenseLayer (4, 3, afun=Abs (), optim=SGD (eta))]
91  netz = SequentialNet (2, layers)
92  netz.train (x, y, bs, ep)
93
94  ytilde = netz.evaluate (x)
95
96  print ('y␣=\n', y)
97  print ('ytilde␣=\n', ytilde)
98  print ('error␣=', np.linalg.norm (y - ytilde))
99  print ('n␣=␣4,␣Abs')
100
101 layers = [DenseLayer (2, 10, afun=Abs (), optim=SGD (eta)),
102            DenseLayer (10, 3, afun=Abs (), optim=SGD (eta))]
103 netz = SequentialNet (2, layers)
104 netz.train (x, y, bs, ep)
105
106 ytilde = netz.evaluate (x)
107
108 print ('y␣=\n', y)
109 print ('ytilde␣=\n', ytilde)
110 print ('error␣=', np.linalg.norm (y - ytilde))
111 print ('n␣=␣10,␣Abs')
112
113 layers = [DenseLayer (2, 100, afun=Abs (), optim=SGD (eta)),
114            DenseLayer (100, 3, afun=Abs (), optim=SGD (eta))]
115 netz = SequentialNet (2, layers)
116 netz.train (x, y, bs, ep)
117
118 ytilde = netz.evaluate (x)
119
120 print ('y␣=\n', y)
```

```
121  print('ytilde␣=\n', ytilde)
122  print('error␣=', np.linalg.norm(y - ytilde))
123  print('n␣=␣100,␣Abs')
```

The following listing gives the TF/Keras counterpart,

```
 1  import numpy as np
 2  import matplotlib.pyplot as plt
 3  import tensorflow as tf
 4  from time import sleep
 5
 6  # training data for exercise 1b), sheet 1
 7  x_train = np.array( [[0,0], [0,1], [1,0], [1,1]] )
 8  y_train = np.array( [[0,0,0],   # xor, and, or
 9                       [1,0,1],
10                       [1,0,1],
11                       [0,1,1]])
12
13  # overall parameters
14  bs, ep, eta = 1, 1000, .1
15  sleep_time = 2
16
17  ###########################################################
18  # use minimal net for exercise 1b), sheet 1 as example
19  ###########################################################
20  # TF/Keras model
21  model = tf.keras.models.Sequential([
22      tf.keras.layers.Dense(3, input_shape=(2,),
23                            activation='relu',
24                            dtype='float64'),
25      tf.keras.layers.Dense(3, activation='relu',
26                            dtype='float64')
27  ])
28
29  # TF/Keras training
30  sgd = tf.keras.optimizers.SGD(learning_rate=eta,
31                                momentum=0.0, nesterov=False)
32  model.compile(optimizer=sgd,
33                loss='mean_squared_error',
34                metrics=['accuracy'])
35  model.fit(x_train, y_train, epochs=ep, batch_size=bs)
36
37  y_tilde = model.predict(x_train)
38  print("y_train␣=", y_train)
39  print("y_tilde␣=", y_tilde)
40  print("error␣=", np.linalg.norm(y_train-y_tilde))
41  print("n␣=␣3")
42  sleep(sleep_time)
43
44  ###########################################################
45  # use net of exercise 1b), sheet 1 as example
46  ###########################################################
47  # TF/Keras model
48  model = tf.keras.models.Sequential([
49      tf.keras.layers.Dense(4, input_shape=(2,),
50                            activation='relu',
51                            dtype='float64'),
```

```
52       tf.keras.layers.Dense(3, activation='relu',
53                                 dtype='float64')
54 ])
55
56 # TF/Keras training
57 sgd = tf.keras.optimizers.SGD(learning_rate=eta,
58                                 momentum=0.0, nesterov=False)
59 model.compile(optimizer=sgd,
60                loss='mean_squared_error',
61                metrics=['accuracy'])
62 model.fit(x_train, y_train, epochs=ep, batch_size=bs)
63
64 y_tilde = model.predict(x_train)
65 print("y_train␣=", y_train)
66 print("y_tilde␣=", y_tilde)
67 print("error␣=", np.linalg.norm(y_train-y_tilde))
68 print("n␣=␣4")
69 sleep(sleep_time)
70
71 ##############################################################
72 # use 10 neurons in hidden layer for exercise 1b), sheet 1
73 ##############################################################
74 # TF/Keras model
75 model = tf.keras.models.Sequential([
76     tf.keras.layers.Dense(10, input_shape=(2,),
77                             activation='relu',
78                             dtype='float64'),
79     tf.keras.layers.Dense(3, activation='relu',
80                             dtype='float64')
81 ])
82
83 # TF/Keras training
84 sgd = tf.keras.optimizers.SGD(learning_rate=eta,
85                                 momentum=0.0, nesterov=False)
86 model.compile(optimizer=sgd,
87                loss='mean_squared_error',
88                metrics=['accuracy'])
89 model.fit(x_train, y_train, epochs=ep, batch_size=bs)
90
91 y_tilde = model.predict(x_train)
92 print("y_train␣=", y_train)
93 print("y_tilde␣=", y_tilde)
94 print("error␣=", np.linalg.norm(y_train-y_tilde))
95 print("n␣=␣10")
96 sleep(sleep_time)
97
98 ##############################################################
99 # use 100 neurons in hidden layer for exercise 1b), sheet 1
100 ##############################################################
101 # TF/Keras model
102 model = tf.keras.models.Sequential([
103     tf.keras.layers.Dense(100, input_shape=(2,),
104                             activation='relu',
105                             dtype='float64'),
106     tf.keras.layers.Dense(3, activation='relu',
107                             dtype='float64')
```

```
108  ])
109
110  # TF/Keras training
111  sgd = tf.keras.optimizers.SGD(learning_rate=eta,
112                                 momentum=0.0, nesterov=False)
113  model.compile(optimizer=sgd,
114                loss='mean_squared_error',
115                metrics=['accuracy'])
116  model.fit(x_train, y_train, epochs=ep, batch_size=bs)
117
118  y_tilde = model.predict(x_train)
119  print("y_train␣=", y_train)
120  print("y_tilde␣=", y_tilde)
121  print("error␣=", np.linalg.norm(y_train-y_tilde))
122  print("n␣=␣100")
```

We observe that we need multiple starts to find a solution. The implementation with the absolute value tends to find a solution more easily. With the chosen parameters for the batch size, the number of epochs, and the learning rate, the net either converges to almost machine precision or stagnates, the reason is the dying of ReLU neurons. Here, regularization would help.

b) The implementation of the neural net of part b) is given in the following listing for our class SequentialNet, where again we included the case of the absolute value as activation function,

```
 1  import numpy              as np
 2  import matplotlib.pyplot as plt
 3
 4  from networks     import SequentialNet
 5  from activations  import Abs
 6  from layers       import DenseLayer
 7  from optimizers   import SGD
 8
 9  ##############################################
10  # Aufgabe 6b)                               #
11  ##############################################
12
13  ####################################
14  # ReLU as activations function     #
15  ####################################
16
17  x = np.expand_dims(np.linspace(0, np.pi/2, 2000), 0)
18  y = np.concatenate((np.sin(x), np.cos(x)))/2 + .5
19
20  x_train, y_train = x[:,::2], y[:,::2]
21  x_test, y_test = x[:,1::2], y[:,1::2]
22
23  bs, ep, eta = 1, 100, .1
24
25  layers=[DenseLayer(1, 100, optim=SGD(eta)),
26          DenseLayer(100, 2, optim=SGD(eta))]
27  netz = SequentialNet(1, layers)
28  netz.train(x_train, y_train, bs, ep)
29
30  ytilde = netz.evaluate(x_test)
31
```

```
32  print('error␣=', np.linalg.norm(y_test - ytilde))
33
34  plt.title('approximate␣sine␣and␣cosine␣using␣ReLU')
35  plt.plot(x.T, y.T, '-')
36  plt.plot(x_test.T, ytilde.T, '-.')
37  plt.plot(x_test.T, np.log10(np.abs(y_test.T - ytilde.T)), '--')
38  plt.legend(['sin', 'cos',
39               'appr.␣sin', 'appr.␣cos',
40               'err.␣sin', 'err␣cos'])
41  plt.show()
42  ##########################################
43  # Abs as activation function            #
44  ##########################################
45
46  x = np.expand_dims(np.linspace(0, np.pi/2, 2000), 0)
47  y = np.concatenate((np.sin(x), np.cos(x)))/2 + .5
48
49  x_train, y_train = x[:,::2], y[:,::2]
50  x_test, y_test = x[:,1::2], y[:,1::2]
51
52  bs, ep, eta = 1, 100, .1
53
54  layers=[DenseLayer(1, 100, afun=Abs(), optim=SGD(eta)),
55           DenseLayer(100, 2, afun=Abs(), optim=SGD(eta))]
56  netz = SequentialNet(1, layers)
57  netz.train(x_train, y_train, bs, ep)
58
59  ytilde = netz.evaluate(x_test)
60
61  print('error␣=', np.linalg.norm(y_test - ytilde))
62
63  plt.title('approximate␣sine␣and␣cosine␣using␣Abs')
64  plt.plot(x.T, y.T, '-')
65  plt.plot(x_test.T, ytilde.T, '-.')
66  plt.plot(x_test.T, np.log10(np.abs(y_test.T - ytilde.T)), '--')
67  plt.legend(['sin', 'cos',
68               'appr.␣sin', 'appr.␣cos',
69               'err.␣sin', 'err␣cos'])
70  plt.show()
```

The following listing gives the TF/Keras counterpart,

```
 1  import sys
 2  import numpy as np
 3  import matplotlib.pyplot as plt
 4  import tensorflow as tf
 5
 6  # training data
 7  x = np.expand_dims(np.linspace(0, np.pi/2, 1000),0)
 8  y = np.concatenate((np.sin(x), np.cos(x)))/2+.5
 9  x_train = x.T
10  y_train = y.T
11
12  bs, ep, eta = 32, 100, .1
13  h_neurons = 100
14
15  ##################################################
```

```
16 # use FNN with one hidden layer comprising 100 neurons
17 #############################################################
18 # TF/Keras model
19 model = tf.keras.models.Sequential([
20     tf.keras.layers.Dense(h_neurons, input_shape=(1,),
21                           activation='relu',
22                           dtype='float64'),
23     tf.keras.layers.Dense(2, activation='relu',
24                           dtype='float64')
25 ])
26
27 # TF/Keras training
28 sgd = tf.keras.optimizers.SGD(learning_rate=eta,
29                               momentum=0.0, nesterov=False)
30 model.compile(optimizer=sgd,
31               loss='mean_squared_error',
32               metrics=['accuracy'])
33 model.fit(x_train, y_train, epochs=ep, batch_size=bs)
34
35 y_tilde = model.predict(x_train)
36 print("error␣=", np.linalg.norm(y_train-y_tilde))
37 plt.plot(x_train, y_train, '-')
38 plt.plot(x_train, y_tilde, '-.')
39 plt.plot(x_train, np.log10(np.abs(y_train-y_tilde)), '--')
40 plt.legend(['sin', 'cos',
41             'appr.␣sin', 'appr.␣cos',
42             'err.␣sin', 'err.␣cos'])
43 plt.show()
```

We observe that sometimes the approximation fails, which occurs more frequently for ReLU than for the abolute value. Typical resulting pictures for the chosen batch size, the number of epochs, and the learning rate in case of a working run are given in Figure 1.
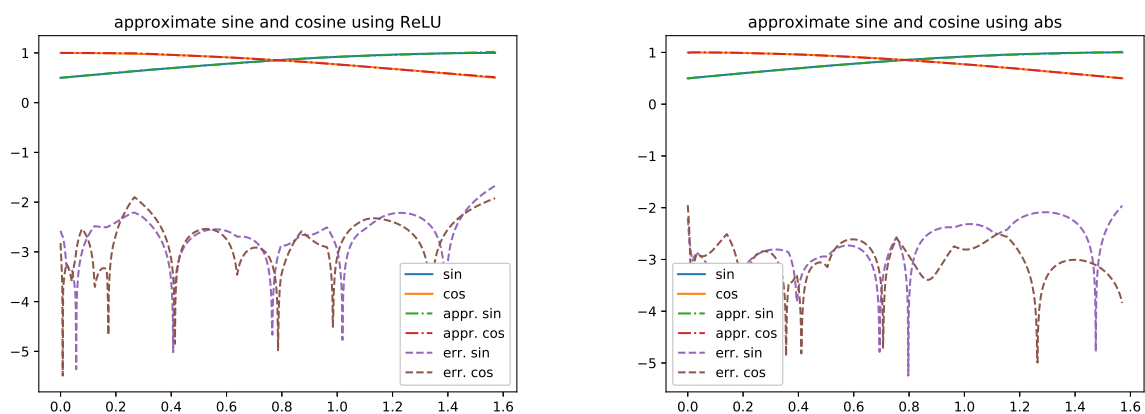


Figure 1: approximation of sine and cosine

c) The code for our class `SequentialNet` is given in the following listing,

```
1 import numpy          as np
2 import tensorflow     as tf
3 import matplotlib.pyplot as plt
4 from   random         import randrange
```

```
 5
 6  from networks    import SequentialNet
 7  from layers      import DenseLayer
 8  from optimizers import SGD
 9
10  DATA = np.load('mnist.npz')
11  x_train, y_train = DATA['x_train'], DATA['y_train'],
12  x_test, y_test   = DATA['x_test'], DATA['y_test']
13  x_train, x_test = x_train / 255.0, x_test / 255.0
14
15  x = x_train.reshape(60000, 784).T
16  I = np.eye(10)
17  y = I[:,y_train]
18
19  bs, ep, eta = 10, 10, .01
20
21  layers = [DenseLayer(784, 100, optim=SGD(eta)),
22            DenseLayer(100, 10,  optim=SGD(eta))]
23  netz = SequentialNet(784, layers)
24  netz.train(x, y, bs, ep)
25
26  y_tilde = netz.evaluate(x_test.reshape(10000, 784, 1))
27  guess   = np.argmax(y_tilde, 1).T
28  print('accuracy␣=', np.sum(guess == y_test)/100)
29
30  for i in range(4):
31
32      k = randrange(y_test.size)
33      plt.title('Label␣is␣{lb},␣guess␣is␣{gs}'.format(lb=y_test[k], gs=
            guess[0,k]))
34      plt.imshow(x_test[k], cmap='gray')
35      plt.show()
```

The TF/Keras counterpart is

```
 1  import numpy              as np
 2  import tensorflow         as tf
 3  import matplotlib.pyplot as plt
 4  from    random            import randrange
 5
 6  mnist = tf.keras.datasets.mnist
 7  (x_train, y_train), (x_test, y_test) = mnist.load_data()
 8  x_train, x_test = x_train / 255.0, x_test / 255.0
 9
10  x = x_train.reshape(60000, 784)
11  I = np.eye(10)
12  y = I[y_train,:]
13
14  bs, ep, eta = 10, 10, .01
15
16  model = tf.keras.models.Sequential([
17      tf.keras.layers.Dense(100, input_shape=(784,),
18                            activation='relu',
19                            dtype='float64'),
20      tf.keras.layers.Dense(10, activation='relu',dtype='float64')
21  ])
22
```

```
23  sgd = tf.keras.optimizers.SGD(learning_rate=eta, momentum=0.0, nesterov=
        False)
24
25  model.compile(optimizer=sgd, loss='mean_squared_error',metrics=['
        accuracy'])
26  model.fit(x, y, epochs=ep, batch_size=bs)
27
28  y_tilde = model.predict(x_test.reshape(10000, 784))
29  guess   = np.argmax(y_tilde, 1)
30  print('accuracy␣=', np.sum(guess == y_test)/100)
31
32  for i in range(4):
33
34      k = randrange(y_test.size)
35      plt.title('Label␣is␣{lb},␣guess␣is␣{gs}'.format(lb=y_test[k], gs=
            guess[k]))
36      plt.imshow(x_test[k], cmap='gray')
37      plt.show()
```

The training of the net results typically in a value around 90% accuracy on the test set, best values are around 95–97%. This best value is more easily obtained for the absolute value activation function. The snippet presented in the lecture used additionally dropout after the hidden layer as regularization.

**Exercise 6:** (2+2 points) In this exercise we consider part of the proof of the backpropagation. Let

$$\mathbf{z} = \mathbf{W}\mathbf{a} + \mathbf{b}, \quad \mathbf{z}, \mathbf{b} \in \mathbb{R}^m, \quad \mathbf{a} \in \mathbb{R}^n, \quad \mathbf{W} \in \mathbb{R}^{m \times n}$$

denote the affine mapping in one of the layers. Let $C$ be a cost function and let

$$\boldsymbol{\delta} = \frac{\partial C}{\partial \mathbf{z}} = \nabla_{\mathbf{z}} C \in \mathbb{R}^m$$

be given. Let the 3-tensor[1] $\mathbf{T} \in \mathbb{R}^{m \times m \times n}$ be defined by elements

$$t_{i,j}^{(\ell)} := \frac{\partial z_\ell}{\partial w_{i,j}}$$

for $\ell = 1, \ldots, m$, $i = 1, \ldots, m$, $j = 1, \ldots, n$.

a) Compute for $\ell = 1, \ldots, m$ the 2-tensor slices

$$\mathbf{T}^{(\ell)} := \begin{pmatrix} t_{1,1}^{(\ell)} & \cdots & t_{1,n}^{(\ell)} \\ \vdots & \ddots & \vdots \\ t_{m,1}^{(\ell)} & \cdots & t_{m,n}^{(\ell)} \end{pmatrix} \in \mathbb{R}^{m \times n}.$$

b) Show that the application of the 3-tensor $\mathbf{T}$ to the 1-tensor $\boldsymbol{\delta}$ gives the same 2-tensor as obtained using the elementwise chain rule,

$$\mathbf{T}\boldsymbol{\delta} := \sum_{\ell=1}^{m} \mathbf{T}^{(\ell)} \mathbf{e}_\ell^\mathsf{T} \boldsymbol{\delta} = \frac{\partial C}{\partial \mathbf{W}} \in \mathbb{R}^{m \times n}$$

---

[1]A 3-tensor is a natural generalization of a scalar (0-tensor), column vector (1-tensor), and matrix (2-tensor) to an additional dimension. Think of it as a cube of numbers, or as matrices stacked on top of each other towards the reader. Here, it makes sense to think of it as a row vector of matrices.

**Solution:**

a) We have

$$t_{i,j}^{(\ell)} = \frac{\partial z_\ell}{\partial w_{i,j}} = \frac{\partial \left( \sum_{k=1}^n w_{\ell,k} a_k + b_\ell \right)}{\partial w_{i,j}} = \begin{cases} a_j & \text{for } \ell = i, \\ 0 & \text{else.} \end{cases}$$

and thus

$$\mathbf{T}^{(\ell)} = \begin{pmatrix} 0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 0 \\ a_1 & \cdots & a_n \\ 0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 0 \end{pmatrix} \quad \leftarrow \quad \ell\text{-th row}$$

$$= \mathbf{e}_\ell \mathbf{a}^\mathsf{T}.$$

b) We need to calculate a tensor product. The matrix multiplication of two matrices $\mathbf{A}$ and $\mathbf{B}$

$$(\mathbf{A} \cdot \mathbf{B})_{i,j} = \sum_k a_{i,k} \cdot b_{k,j}$$

can be seen as being split in two parts. First we pair those entries of $\mathbf{A}$ and $\mathbf{B}$ where the second index of $\mathbf{A}$ and first index of $\mathbf{B}$ are identical. This is called *pairing*. Then we aggregate over the common index which is called *reducing*. One can also apply this logic to the matrix vector multiplication. Here we match the second index of $\mathbf{A}$ with the only index of $\mathbf{x}$:

$$(\mathbf{A} \cdot \mathbf{x})_i = \sum_k a_{i,k} \cdot x_k$$

Applying this in our case by pairing the index denoted by $\ell$ with the only index of $\delta$ and then reducing we get

$$(\mathbf{T}\boldsymbol{\delta})_{i,j} = \sum_{\ell=1}^m t_{i,j}^{(\ell)} \cdot \delta_\ell = a_j \cdot \delta_i \Rightarrow \mathbf{T}\boldsymbol{\delta} = \boldsymbol{\delta}\mathbf{a}^\mathsf{T}.$$

Using the hint in the exercise

$$\mathbf{T}\boldsymbol{\delta} = \sum_{\ell=1}^m \mathbf{T}^{(\ell)} \mathbf{e}_\ell^\mathsf{T} \boldsymbol{\delta} = \sum_{\ell=1}^m \mathbf{e}_\ell \mathbf{a}^\mathsf{T} \delta_\ell = \boldsymbol{\delta}\mathbf{a}^\mathsf{T}$$

we get the same result.