
Mathematics of Neural Networks
 winter semester 2021/2022
 exercise sheet 4, solutions

Exercise 1: (2 points) Let $\mathbf{i} = (i_1, \dots, i_d)$ be a multi-index. Prove that multi-dimensional convolution

$$(\mathbf{X} * \mathbf{F})(n_1, \dots, n_d) = \sum_{\mathbf{i}} \mathbf{X}(n_1 - i_1, \dots, n_d - i_d) \cdot \mathbf{F}(i_1, \dots, i_d)$$

is commutative, i.e., $\mathbf{X} * \mathbf{F} = \mathbf{F} * \mathbf{X}$.

Solution: We define the multi-index $\mathbf{j} = (j_1, \dots, j_d)$ and set $\mathbf{j} = \mathbf{n} - \mathbf{i}$. Thus, $\mathbf{i} = \mathbf{n} - \mathbf{j}$ and

$$\begin{aligned} (\mathbf{X} * \mathbf{F})(n_1, \dots, n_d) &= \sum_{\mathbf{i}} \mathbf{X}(n_1 - i_1, \dots, n_d - i_d) \cdot \mathbf{F}(i_1, \dots, i_d) \\ &= \sum_{\mathbf{j}} \mathbf{X}(j_1, \dots, j_d) \cdot \mathbf{F}(n_1 - j_1, \dots, n_d - j_d) = (\mathbf{F} * \mathbf{X})(n_1, \dots, n_d). \end{aligned}$$

Exercise 2: (4 points) In Lecture 4 we showed that

$$\frac{\partial C}{\partial \mathbf{f}} = \widetilde{(\mathbf{a} * \boldsymbol{\delta})} = \tilde{\mathbf{a}} * \boldsymbol{\delta},$$

for the case of a 1D kernel \mathbf{f} (see lecture notes, p. 93¹), where the tilde represents a flipping of the entries. We now consider a 2D kernel $\mathbf{F} \in \mathbb{R}^{2 \times 2}$ with input $\mathbf{A} \in \mathbb{R}^{3 \times 3}$. The convolution $\mathbf{A} * \mathbf{F}$ is the given by

$$\begin{pmatrix} y_{1,1} & y_{1,2} \\ y_{2,1} & y_{2,2} \end{pmatrix} = \begin{pmatrix} a_{1,1} & a_{1,2} & a_{1,3} \\ a_{2,1} & a_{2,2} & a_{2,3} \\ a_{3,1} & a_{3,2} & a_{3,3} \end{pmatrix} * \begin{pmatrix} f_{1,1} & f_{1,2} \\ f_{2,1} & f_{2,2} \end{pmatrix} = \begin{pmatrix} a_{1,1} & a_{1,2} & a_{1,3} \\ a_{2,1} & a_{2,2} & a_{2,3} \\ a_{3,1} & a_{3,2} & a_{3,3} \end{pmatrix} \star \begin{pmatrix} f_{2,2} & f_{2,1} \\ f_{1,2} & f_{1,1} \end{pmatrix},$$

or in vector notation by

$$\begin{pmatrix} y_{1,1} \\ y_{1,2} \\ y_{2,1} \\ y_{2,2} \end{pmatrix} = \left(\begin{array}{ccc|ccc|ccc} f_{2,2} & f_{2,1} & 0 & f_{1,2} & f_{1,1} & 0 & 0 & 0 & 0 \\ 0 & f_{2,2} & f_{2,1} & 0 & f_{1,2} & f_{1,1} & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & f_{2,2} & f_{2,1} & 0 & f_{1,2} & f_{1,1} & 0 \\ 0 & 0 & 0 & 0 & f_{2,2} & f_{2,1} & 0 & f_{1,2} & f_{1,1} \end{array} \right) \cdot \begin{pmatrix} a_{1,1} \\ a_{1,2} \\ a_{1,3} \\ a_{2,1} \\ a_{2,2} \\ a_{2,3} \\ a_{3,1} \\ a_{3,2} \\ a_{3,3} \end{pmatrix}.$$

¹Version from 26.10.2021

- a) How does the Δ from the backpropagation look like in this case?
 b) Prove the equality

$$\widetilde{(\mathbf{A} * \tilde{\Delta})} = \tilde{\mathbf{A}} * \Delta.$$

Solution:

- a) The input is a matrix $\mathbf{A} \in \mathbb{R}^{3 \times 3}$, the kernel is $\mathbf{F} \in \mathbb{R}^{2 \times 2}$, thus the output for a complete overlap of kernel and input is a matrix $\mathbf{Y} \in \mathbb{R}^{2 \times 2}$. This implies that the convolution maps a layer from $\mathbb{R}^{3 \times 3}$ to the next layer in $\mathbb{R}^{2 \times 2}$, $\mathbf{Z} = \mathbf{A} * \mathbf{F} + \mathbf{B} \in \mathbb{R}^{2 \times 2}$, so the backpropagation will start with a $\Delta \in \mathbb{R}^{2 \times 2}$ in the new layer. The update of the kernel is based on $\mathbf{A} \in \mathbb{R}^{3 \times 3}$ and $\Delta \in \mathbb{R}^{2 \times 2}$ and gives a matrix in $\mathbb{R}^{2 \times 2}$. The matrix Δ and the flipped matrix $\tilde{\Delta}$ have the entries

$$\Delta = \begin{pmatrix} \delta_{1,1} & \delta_{1,2} \\ \delta_{2,1} & \delta_{2,2} \end{pmatrix}, \quad \tilde{\Delta} = \begin{pmatrix} \delta_{2,2} & \delta_{2,1} \\ \delta_{1,2} & \delta_{1,1} \end{pmatrix}.$$

- b) We switch from matrices to vectors and use the commutativity of the convolution,

$$\begin{pmatrix} z_{1,1} \\ z_{1,2} \\ z_{2,1} \\ z_{2,2} \end{pmatrix} = \left(\begin{array}{ccc|ccc} f_{2,2} & f_{2,1} & 0 & f_{1,2} & f_{1,1} & 0 \\ 0 & f_{2,2} & f_{2,1} & 0 & f_{1,2} & f_{1,1} \\ 0 & 0 & 0 & f_{2,2} & f_{2,1} & 0 \\ 0 & 0 & 0 & 0 & f_{2,2} & f_{2,1} \end{array} \begin{array}{c} 0 \\ 0 \\ f_{1,2} \\ f_{1,1} \\ 0 \\ 0 \end{array} \right) \begin{pmatrix} a_{1,1} \\ a_{1,2} \\ a_{1,3} \\ a_{2,1} \\ a_{2,2} \\ a_{2,3} \\ a_{3,1} \\ a_{3,2} \\ a_{3,3} \end{pmatrix} + \begin{pmatrix} b_{1,1} \\ b_{1,2} \\ b_{2,1} \\ b_{2,2} \end{pmatrix}$$

$$= \begin{pmatrix} a_{2,2} & a_{2,1} & a_{1,2} & a_{1,1} \\ a_{2,3} & a_{2,2} & a_{1,3} & a_{1,2} \\ a_{3,2} & a_{3,1} & a_{2,2} & a_{2,1} \\ a_{3,3} & a_{3,2} & a_{2,3} & a_{2,2} \end{pmatrix} \begin{pmatrix} f_{1,1} \\ f_{1,2} \\ f_{2,1} \\ f_{2,2} \end{pmatrix} + \begin{pmatrix} b_{1,1} \\ b_{1,2} \\ b_{2,1} \\ b_{2,2} \end{pmatrix} = \mathbf{T}\mathbf{f} + \mathbf{b}.$$

The derivative of \mathbf{z} with respect to \mathbf{f} is thus given by \mathbf{T}^\top and

$$\begin{pmatrix} \partial C / \partial f_{1,1} \\ \partial C / \partial f_{1,2} \\ \partial C / \partial f_{2,1} \\ \partial C / \partial f_{2,2} \end{pmatrix} = \frac{\partial C}{\partial \mathbf{f}} = \frac{\partial \mathbf{z}}{\partial \mathbf{f}} \frac{\partial C}{\partial \mathbf{z}} = \mathbf{T}^\top \boldsymbol{\delta} = \begin{pmatrix} a_{2,2} & a_{2,3} & a_{3,2} & a_{3,3} \\ a_{2,1} & a_{2,2} & a_{3,1} & a_{3,2} \\ a_{1,2} & a_{1,3} & a_{2,2} & a_{2,3} \\ a_{1,1} & a_{1,2} & a_{2,1} & a_{2,2} \end{pmatrix} \begin{pmatrix} \delta_{1,1} \\ \delta_{1,2} \\ \delta_{2,1} \\ \delta_{2,2} \end{pmatrix}.$$

We rewrite the latter convolution as matrix-vector product with the vector of the elements of the matrix \mathbf{A} ,

$$\begin{pmatrix} a_{2,2} & a_{2,3} & a_{3,2} & a_{3,3} \\ a_{2,1} & a_{2,2} & a_{3,1} & a_{3,2} \\ a_{1,2} & a_{1,3} & a_{2,2} & a_{2,3} \\ a_{1,1} & a_{1,2} & a_{2,1} & a_{2,2} \end{pmatrix} \begin{pmatrix} \delta_{1,1} \\ \delta_{1,2} \\ \delta_{2,1} \\ \delta_{2,2} \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 & \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} \\ 0 & 0 & 0 & \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} & 0 \\ 0 & \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} & 0 & 0 & 0 \\ \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} a_{1,1} \\ a_{1,2} \\ a_{1,3} \\ a_{2,1} \\ a_{2,2} \\ a_{2,3} \\ a_{3,1} \\ a_{3,2} \\ a_{3,3} \end{pmatrix}$$

We note that flipping the vector of the elements of \mathbf{A} results in the vector of the doubly flipped activated element $\tilde{\mathbf{A}}$,

$$\mathbf{A} = \begin{pmatrix} a_{1,1} & a_{1,2} & a_{1,3} \\ a_{2,1} & a_{2,2} & a_{2,3} \\ a_{3,1} & a_{3,2} & a_{3,3} \end{pmatrix} \rightsquigarrow \begin{pmatrix} a_{1,1} \\ a_{1,2} \\ a_{1,3} \\ a_{2,1} \\ a_{2,2} \\ a_{2,3} \\ a_{3,1} \\ a_{3,2} \\ a_{3,3} \end{pmatrix} \Leftrightarrow \begin{pmatrix} a_{3,3} \\ a_{3,2} \\ a_{3,1} \\ a_{2,3} \\ a_{2,2} \\ a_{2,1} \\ a_{1,3} \\ a_{1,2} \\ a_{1,1} \end{pmatrix} \rightsquigarrow \begin{pmatrix} a_{3,3} & a_{3,2} & a_{3,1} \\ a_{2,3} & a_{2,2} & a_{2,1} \\ a_{1,3} & a_{1,2} & a_{1,1} \end{pmatrix} = \tilde{\mathbf{A}}.$$

The two ways of writing the derivative of the cost function C with respect to the elements of \mathbf{f} are given by flipping the rows of both sides or to flipping the rows of the vector of \mathbf{A} 's elements and the columns of the block Toeplitz matrix of the δ elements. The flipping of the rows on both sides results in

$$\widetilde{\frac{\partial C}{\partial \mathbf{f}}} = \begin{pmatrix} \partial C / \partial f_{2,2} \\ \partial C / \partial f_{2,1} \\ \partial C / \partial f_{1,2} \\ \partial C / \partial f_{1,1} \end{pmatrix} = \begin{pmatrix} \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} & 0 & 0 & 0 & 0 \\ 0 & \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} & 0 & 0 & 0 \\ 0 & 0 & 0 & \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} & 0 \\ 0 & 0 & 0 & 0 & \delta_{1,1} & \delta_{1,2} & 0 & \delta_{2,1} & \delta_{2,2} \end{pmatrix} \begin{pmatrix} a_{1,1} \\ a_{1,2} \\ a_{1,3} \\ a_{2,1} \\ a_{2,2} \\ a_{2,3} \\ a_{3,1} \\ a_{3,2} \\ a_{3,3} \end{pmatrix},$$

which corresponds to

$$\widetilde{\frac{\partial C}{\partial \mathbf{F}}} = \mathbf{A} * \tilde{\mathbf{\Delta}} = \mathbf{A} \star \mathbf{\Delta}.$$

The flipping of the elements of \mathbf{A} and the columns of the block Toeplitz matrix results in

$$\frac{\partial C}{\partial \mathbf{f}} = \begin{pmatrix} \partial C / \partial f_{1,1} \\ \partial C / \partial f_{1,2} \\ \partial C / \partial f_{2,1} \\ \partial C / \partial f_{2,2} \end{pmatrix} = \begin{pmatrix} \delta_{2,2} & \delta_{2,1} & 0 & \delta_{1,2} & \delta_{1,1} & 0 & 0 & 0 & 0 \\ 0 & \delta_{2,2} & \delta_{2,1} & 0 & \delta_{1,2} & \delta_{1,1} & 0 & 0 & 0 \\ 0 & 0 & 0 & \delta_{2,2} & \delta_{2,1} & 0 & \delta_{1,2} & \delta_{1,1} & 0 \\ 0 & 0 & 0 & 0 & \delta_{2,2} & \delta_{2,1} & 0 & \delta_{1,2} & \delta_{1,1} \end{pmatrix} \begin{pmatrix} a_{3,3} \\ a_{3,2} \\ a_{3,1} \\ a_{2,3} \\ a_{2,2} \\ a_{2,1} \\ a_{1,3} \\ a_{1,2} \\ a_{1,1} \end{pmatrix},$$

which corresponds to

$$\frac{\partial C}{\partial \mathbf{F}} = \tilde{\mathbf{A}} * \mathbf{\Delta} = \tilde{\mathbf{A}} \star \tilde{\mathbf{\Delta}}.$$

Thus we get the desired equality

$$(\widetilde{\mathbf{A} * \tilde{\mathbf{\Delta}}}) = \frac{\partial C}{\partial \mathbf{F}} = \tilde{\mathbf{A}} * \mathbf{\Delta}.$$

Exercise 3: (4 points) In preparation for convolutional networks change the current implementation of `DenseLayer` to treat *derivatives* in place of *gradients*:

- Initialize with transposed weights and biases given by `np.zeros(no)`.
- Evaluate using $\mathbf{z} = \mathbf{a}\mathbf{W} + \mathbf{b}$, $a(\mathbf{z})$.
- Backpropagate using

$$\frac{\partial C}{\partial \mathbf{W}} = \mathbf{a}^\top \boldsymbol{\delta}, \quad \frac{\partial C}{\partial \mathbf{b}} = \boldsymbol{\delta},$$

now by abuse of notation denoting the *derivatives*, and output $\boldsymbol{\delta}\mathbf{W}^\top$.

- Which other parts do you have to change accordingly?

Solution: A possible solution can be found in the following listing:

```

1 class DenseLayer:
2
3     def __init__(self,
4                   ni, # Number of inputs
5                   no, # Number of outputs
6                   afun = None, # Activationfunction for the layer
7                   optim = None,
8                   initializer = None,
9                   kernel_regularizer = None
10                ):
11
12         self.ni    = ni
13         self.no    = no
14
15         if afun is None:
16             self.afun = ReLU()
17         else:
18             self.afun = afun
19
20         if optim is None:
21             self.optim = SGD()
22         else:
23             self.optim = optim
24
25         self.kernel_regularizer = kernel_regularizer
26
27         if initializer is None:
28             self.initializer = RandnAvarage()
29         else:
30             self.initializer = initializer
31
32         self.W      = self.afun.factor * self.initializer.wfun(ni, no)
33         self.b      = np.zeros(no)
34         self.__z    = None
35         self.__a    = None
36         self.dW     = np.zeros_like(self.W)
37         self.db     = np.zeros_like(self.b)
38
39     def evaluate(self, a):
40
41         self.__a = a
42         self.__z = self.__a @ self.W + self.b
43         return self.afun.evaluate(self.__z)
44

```

```

45     def set_weights(self, W):
46
47         assert(W.shape == (self.ni, self.no))
48         self.W = W
49
50     def set_bias(self, b):
51
52         assert(b.size == self.no)
53         self.b = b
54
55     def backprop(self, delta):
56
57
58         delta = self.afun.backprop(delta)
59         self.dW = self.__a.T @ delta
60         self.db = np.sum(delta, 0)
61
62         if not self.kernel_regularizer is None:
63             self.kernel_regularizer.update(self.W, self.dW)
64
65         return delta @ self.W.T
66
67     def update(self):
68
69         self.optim.update([self.W, self.b],
70                           [self.dW, self.db])

```

Note that the calls of `backprop` in the `train` method have to be changed, too.

```

1     def train(self, x, y, batch_size=16, epochs=10):
2
3         n_data = y.shape[0]
4         n_batches = int(np.ceil(n_data/batch_size))
5
6         for e in range(epochs):
7
8             p = np.random.permutation(n_data)
9             for j in range(n_batches):
10
11                 self.backprop(x[p[j*batch_size:(j+1)*batch_size],:],
12                               y[p[j*batch_size:(j+1)*batch_size],:])
13
14             for layer in self.layers:
15                 layer.update()

```

Exercise 4: (4 points)

a) Add a class `DropoutLayer` to `layers.py`:

(i) Initialize with the following attributes:

- `p`: the probability $p \in [0.5, 0.8]$ with which neurons are kept, default: `p = .5`.
- `mask`: a scaled version of a vector with 0's and 1's, specifying the dropped neurons, set to `None` on initialization.
- `trainable`: a Boolean value given on initialization (**True**: the network is trained, neurons are dropped, and the output is scaled; **False**: the network is not trained and the weights remain unchanged).

(ii) Add the following methods:

- `evaluate(self, a)`: If the network is trained, i.e., `trainable` is `True`, neurons are dropped and the input gets scaled. If the network is not trained, i.e., `trainable` is `False`, the input is returned.

Hint: Implement the dropping of neurons by setting `self.mask` using `numpy.random.binomial(1, self.p, size=shape)`.

The parameter `shape` is a tuple $(1, m)$ with m being the number of given input vectors. This ensures the same mask for the whole minibatch. The scaling factor is chosen as $1/(1 - p)$.

- `backprop(self, delta)`: backpropagation by dropping neurons as given in `self.mask`.
- `update`: perform the update after backpropagation. It shouldn't do anything since there are no parameters to update.

b) Test your code by expanding exercise 5c) from sheet 2. Use the logistic function as activation function. Compare the result for different probabilities $p \in [0.5, 0.8]$ to a network without dropout.

Note that a higher learning rate is suggested for networks with dropout. Try different sizes for minibatches and different learning rates as well. Note that a higher learning rate is recommended for layers with dropout, use, e.g., `10*eta` as learning rate for layers w/ dropout, if `eta` is the learning rate for a layer w/o dropout.

Solution: A possible implementation can be found in the following listing:

```

1 class DropoutLayer:
2
3     def __init__(self, p=.5, trainable=True):
4
5         self.p = p
6         self.trainable = trainable
7         self.mask = None
8
9     def evaluate(self, a):
10
11         if self.trainable:
12
13             shape = (1,) + a.shape[1:] # same mask for minibatch
14             scale = 1/(1-self.p)
15             self.mask = scale * \
16                 np.random.binomial(1, self.p, size=shape)
17             drop_a = a * self.mask
18             return drop_a
19         else:
20             return a
21
22     def backprop(self, delta):
23
24         delta *= self.mask
25         return delta
26
27     def update(self):
28         pass

```

If a layer with dropout is evaluated for testing, the learned weights lead to a higher output since the weights were learned for thinned layers. Thus the output is typically scaled by p . Another approach, which we used for our implementation, is to scale the input while the network is trained.

A possible solution for part b) can be found in the following listing.

```

1 import numpy as np
2 from random import randrange
3 import matplotlib.pyplot as plt
4
5 from networks import SequentialNet
6 from layers import *
7 from optimizers import *
8 from activations import *
9
10 DATA = np.load('mnist.npz')
11 #DATA = np.load('fashion_mnist.npz')
12 x_train, y_train = DATA['x_train'], DATA['y_train']
13 x_test, y_test = DATA['x_test'], DATA['y_test']
14 x_train, x_test = x_train / 255.0, x_test / 255.0
15
16 categories = ['T-shirt/top', 'Trouser', 'Pullover', 'Dress', 'Coat', 'Sandal',
17              ', 'Shirt', 'Sneaker', 'Bag', 'Ankleboot']
18
19 x = x_train.reshape(60000, 784)
20 I = np.eye(10)
21 y = I[y_train,:]
22
23 bs, ep, eta = 1000, 10, .001
24
25 print('Ohne Dropout')
26 layers = [DenseLayer(784, 100, afun=Logistic(), optim=SGD(eta)),
27           DenseLayer(100, 10, afun=Logistic(), optim=SGD(eta))]
28 netz = SequentialNet(784, layers)
29 netz.train(x, y, bs, ep)
30
31 y_tilde = netz.evaluate(x_test.reshape(10000, 784))
32 guess = np.argmax(y_tilde, 1).T
33 print('accuracy=', np.sum(guess == y_test)/100)
34
35 for i in range(4):
36     k = randrange(y_test.size)
37     plt.title('Label is {lb}, guess is {gs}'.format(lb=y_test[k], gs=guess[k]))
38     plt.imshow(x_test[k], cmap='gray')
39     plt.show()
40
41 prop = np.linspace(.5, .8, 7, True)
42 for p in prop:
43     print('Mit Dropout, p=', p)
44     layers = [DenseLayer(784, 100, afun=Logistic(), optim=SGD(eta*10)),
45               DropoutLayer(p),
46               DenseLayer(100, 10, afun=Logistic(), optim=SGD(eta*10))]
47     netz = SequentialNet(784, layers)
48     netz.train(x, y, bs, ep)

```

```
49 y_tilde = netz.evaluate(x_test.reshape(10000,784))
50 guess   = np.argmax(y_tilde, 1)
51 print('accuracy=', np.sum(guess == y_test)/100)
52
53
54 for i in range(4):
55
56     k = randrange(y_test.size)
57     plt.title('Label is {lb}, guess is {gs}'.format(lb=y_test[k], gs=
58             guess[k]))
59     plt.imshow(x_test[k], cmap='gray')
    plt.show()
```