

Problem description

The following problem is proposed:

6x6 images are generated as one of two distinct classes.

- Class 0 represents noise as i.i.d. - Uniform(0w, high), where low and high are parameters that determine the noise level.
 - Class 1 takes class 0 and inserts a 3x3 square by setting each element in a 3x3 grid to a predefined value, starting from a random spot in the image.
- The problem then becomes to create a binary classification network that can separate these two.

Given

- The network architecture

Layer	Input dim	Output dim	Other
Input	6x6	6x6	Class 0 or 1
2D convolve	6x6	4x4	Kernel: 3x3 with stride 1 and no padding. Activation: ReLU
Maxpool	4x4	3x3	Kernel: 2x2 with stride 2 and no padding
Flatten	2x2	4x1	
FCNN	4x1	1	Activation: Sigmoid

which has a total of 15 parameters, decomposed into 2 biases from the CNN and FCNN, 9 kernel weights in the CNN and 4 weights in the FCNN.

- An input matrix X and a label y , where $y = 0$ when X is all uniform and $y = 1$ when X contains a square.

Find

- The intermediate values in both the CNN and FCNN part
- Post-activation for a forward pass
- Loss for a forward pass
- Gradient of the loss function w.r.t. the network parameters

Imports and helper functions

```
In [1]: import numpy as np

import matplotlib.pyplot as plt
from dataclasses import dataclass
from ...future... import annotations
from itertools import product

class Matrix2D():

    def __init__(self, array: np.ndarray):
        self._init(array.shape) = 2, "Input array must be 2D!"

    def self_data: np.ndarray = array
    self._row: int = array.shape[0]
    self._col: int = array.shape[1]

    def __getitem__(self, key):
        return self.data[key]

    def __setitem__(self, key, val):
        self.data[key] = val

    @staticmethod
    def gaussian(center, ndim, sigma = 1):

        xx, yy = np.meshgrid(
            np.linspace(0, ndim-1, ndim) - center[0],
            np.linspace(0, ndim-1, ndim) - center[1]
        )

        return Matrix2D(np.exp(-0.5 * (np.square(xx) + np.square(yy)) / np.square(sigma)))

    @staticmethod
    def uniform_noise(ndim: int, intensity: float = 1.0) -> Matrix2D:
        np.random.seed(0)
        return Matrix2D(np.random.uniform(intensity, intensity, (ndim, ndim)))

    @staticmethod
    def square(ndim: int, ndim_square: int) -> Matrix2D:
        matrix = Matrix2D(uniform_noise(ndim, ndim))

        # Select top-left corner as draw square from
        top_left_row = np.random.randint(low=0, high = ndim, matrix - ndim_square)
        top_left_col = np.random.randint(low=0, high = ndim, matrix - ndim_square)

        matrix[top_left_row:top_left_row + ndim_square, top_left_col:top_left_col + ndim_square] = 0.0

    return matrix

def draw(self, title: str = "") -> None:
    """Display a visualization of the matrix values"""
    See: https://stackoverflow.com/questions/4887753/display-matrix-values-and-colormap
    fig, ax = plt.subplots()
    ax.imshow(self.data, cmap=plt.cm.Blues)

    if self._row == 1:
        for col in range(self._col):
            ax.text(col, 0, f'{self.data[0, col]:.2f}', va='center', ha='center')
        self._col = 1
    elif for row in range(self._row):
        ax.text(0, row, f'{self.data[row, 0]:.2f}', va='center', ha='center')
        for col, row in np.ndindex(self.data.shape):
            ax.text(col, row, f'{self.data[row, col]:.2f}', va='center', ha='center')

    ax.set_title(title)
    ax.set_xlabel(title)

def convolve2d(self, kernel: Matrix2D, bias: float = 0, stride: int = 1) -> Matrix2D:
    """Convolve the 2D matrix with a 2D kernel (blue bias)"""
    See: http://www.songho.ca/dsp/convolution/convolution2d_example.html for the math
    """
    assert stride != 0, "Stride cannot be zero!"

    output_width: int = (self._row - ndim) // stride + 1
    output_height: int = (self._col - ndim) // stride + 1

    output_image = np.zeros((output_height, output_width))
    indices = []

    for row, col in product(range(0, output_width, stride), range(0, output_height, stride)):
        output_image[row, col] = self_data[row:row+kernel.nrow, col:col+kernel.ncol].flatten() \
            * kernel.data.flatten() \
            + bias

    return Matrix2D(output_image)

def maxpool(self, ndim: int = 1, stride: int = 1):
    """Dimension reduction using the max of a neighborhood defined by ndim"""
    See: https://computer-graphics.stackexchange.com/questions/13189/backlog-through-max-pooling-layers
    for the intuition on why the indices of the maximum values for each pool is needed
    """
    assert stride != 0, "Stride cannot be zero!"

    output_width: int = (self._row - ndim) // stride + 1
    output_height: int = (self._col - ndim) // stride + 1

    output_image = np.zeros((output_height, output_width))
    indices = []

    for row, col in product(range(output_height), range(output_width)):
        rows = slice(row * ndim, (row+1) * ndim)
        cols = slice(col * ndim, (col+1) * ndim)
        pool = self.data[rows, cols]

        output_image[row, col] = np.max(pool)
        indices.append(index)

    return Matrix2D(output_image, indices)

def flatten(self):
    return Matrix2D(self.data.flatten().reshape(self.nrow*self._col, 1))

@dataclass
class Network():
    learning_rate: float = 0.1

    kernel: Matrix2D = None
    conv_bias: float = 0.0

    weights: np.ndarray = np.empty(0)
    from_bias: float = 0.0

    backward: bool = False
    maxpool_indices: np.ndarray = np.empty(0)

    # parameters - hardcoded for this solution and will be None until the first forward pass
    z1: Matrix2D = None
    a1: Matrix2D = None
    a2: Matrix2D = None
    a3: Matrix2D = None
    z4: float = None
    a4: float = None

    def fcm_forward(self, x):
        return self.weights * x + self.fcm_bias

    def relu(self, value):
        """ReLU activation"""
        return value if value > 0 else 0

    if self.backward:
        return act * 1

    return act * abs(value)

def leaky_relu(self, value):
    """Leaky ReLU as a solution to dying ReLU"""
    a = 0.01

    if self.backward:
        return np.where(value > 0, 1, a)

    return np.where(value > 0, value, value * a)

def sigmoid(self, value):
    """Sigmoid activation"""
    act = np.exp(value) / (1 + np.exp(value))

    if self.backward:
        return act * (1 - act)

    return act

def binary_cross_entropy(self, y: float, p: float) -> float:
    return -(y * np.log(p) + (1 - y) * np.log(1 - p))

def forward_pass(self, X, y):
    """Full single forward pass, given an input image X and the class y it belongs to"""
    self.backward = False

    self.z1 = X.convolve2d(self.kernel, bias=self.conv_bias, stride=1) # TODO: Move to parameter list
    self.a1 = Matrix2D(self.relu(self.z1.data))

    self.a2, self.maxpool_indices = self.a1.maxpool(ndim=2, stride=2) # TODO: Move to parameter list
    self.a3 = self.a2.flatten()

    self.z4 = self.fcm_forward(self.a3.data) [0, 0]
    self.a4 = self.sigmoid(self.z4)

    return None if y is None else self.binary_cross_entropy(y, self.a4)

def backward_pass(self, x, y):
    """Full single backward pass, given an input image X and the class y it belongs to"""
    self.backward = True

    delta = self.binary_cross_entropy(y, self.a4) * self.sigmoid(self.z4)
    self.fcm_bias += self.learning_rate * delta # update bias for the FCNN

    delta = delta * self.a4.data
    self.weights = self.learning_rate * delta.T # update weights for the FCNN

    delta_temp = np.zeros_like(self.a1.data)
    for max_indices, gradient in zip(self.maxpool_indices, delta):
        delta_temp[max_indices] = gradient

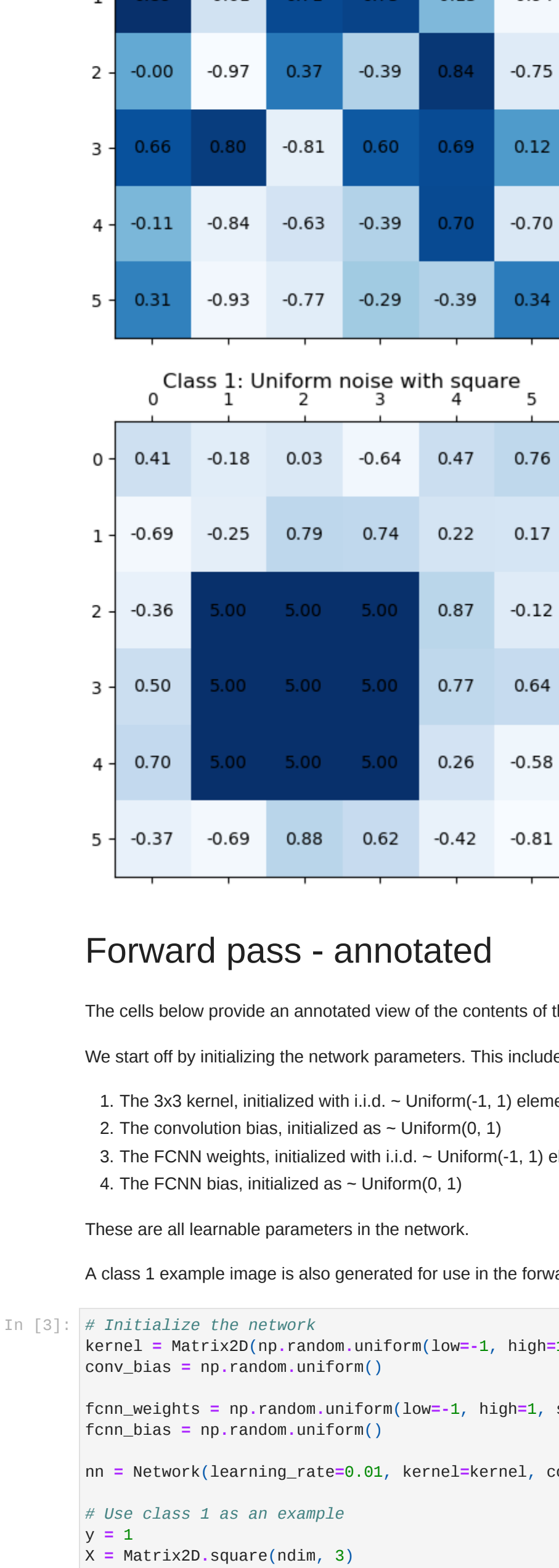
    delta = self.relu(self.z1.data) * delta_temp

    grad_conv_bias = X.data.sum()
    grad_kernel = X.convolve2d(Matrix2D(delta))

    # update weights and bias for convolution
    self.kernel.data += self.learning_rate * grad_kernel.data
    self.conv_bias += self.learning_rate * grad_conv_bias

def predict(self, X):
    """Full single forward pass, without needing the class y"""
    self.forward_pass(X, None)
    print(f"prob. for class 0 (no square): {self.a4}")
    print(f"prob. for class 1 (square): {self.a4}")
```

Example of data and usage of the Network class



Forward pass - annotated

The cells below provide an annotated view of the contents of the forward() function of the Network class

We start off by initializing the network parameters. This includes:

- The 3x3 kernel, initialized with i.i.d. - Uniform(0,1) elements
- The convolution bias, initialized as - Uniform(0,1)
- The FCNN weights, initialized with i.i.d. - Uniform(0,1) elements
- The FCNN bias, initialized as - Uniform(0,1)

These are all learnable parameters in the network.

A class 1 example image is also generated for use in the forward and backward pass example.

```
In [3]: # Initialize the network
kernel = Matrix2D(np.random.uniform(low=-1, high=1, size=(3, 3)))
conv_bias = np.random.uniform()

from_weights = np.random.uniform(low=-1, high=1, size=(1, 4))
from_bias = np.random.uniform()

nn = Network(learning_rate=0.01, kernel=kernel, conv_bias=conv_bias, weights=from_weights, from_bias=from_bias)

# Use class 1 as an example
y = 1
X = Matrix2D(square(ndim, 3))
```

Convolution

Since the kernel is initialized to random values, trying to assign any meaning to the output of the convolution layer at this stage does not make sense. It is instead better to understand how the values in the convolved output image are generated. See [this](#) [gif](#) to get a visualization of the convolution operation.

```
In [4]: z1 = X.convolve2d(kernel, bias=conv_bias, stride=1)
z1.draw("class 1: image convolved with a 3x3 kernel, stride=1")

Class 1: image convolved with a 3x3 kernel, stride=1
0 1 2 3
0 -0.29 -0.80 7.97 -0.29
1 0.86 3.73 4.83 -2.06
2 6.32 1.36 1.78 1.16
3 0.37 0.60 1.51 -0.14
```

ReLU

The ReLU function applies $\max(0, \text{value})$ element-wise which adds a non-linearity to the network. [It](#) lists some of the benefits of having a ReLU activation, mainly sparsity and fast training times. While neither of these are problematic for such a small-scale problem as here, it yields an extra step to discuss and perform calculations on.

```
In [5]: a1 = Matrix2D(nn.relu(z1.data))
a1.draw("a1 = ReLU(z1)")

a1 = ReLU(z1)
0 1 2 3
0 0.00 0.00 7.97 0.00
1 0.86 3.73 4.83 0.00
2 6.32 1.36 1.78 1.16
3 0.37 0.60 1.51 0.00
```

Maxpool

Our binary classifier needs a fully-connected neural network between the augmented input matrix and the output. However, due to the sparsity introduced by the ReLU operation, and for scalability concerns, it is reasonable to decimate the augmented image to a lower dimension, and this is what maxpool does. Mathematically, it is similar to convolution (a sliding window (kernel) is applied to the image. However, instead of weighing the image elements with the kernel data, the maximum image value inside a given window is extracted to the output of the maxpool operation. [This](#) [video](#) is a quick and dirty introduction to maxpooling.

Instead of requiring 16 weights in our FCNN, we now only need 4 - which is good - from both the perspective of parsimony and computational demand (again, the latter is not really a concern for us.)

```
In [6]: a2, maxpool_indices = a1.maxpool(ndim=2, stride=2)
a2.draw("z2 = a1.maxpool(ndim=2, stride=2)")

z2 = a1.maxpool(ndim=2, stride=2)
0 1
0 3.73 7.97
1 6.32 1.78
```

FCNN

The FCNN consists of four weights, one for each cell of the maxpool output. In addition, a bias term is introduced to avoid being restricted to only being able to model functions where $f(0) = 0$

The pre-activation value is calculated by the dot-product between the weights and the input vector, plus bias.

Since the purpose of the network is to discriminate between two classes, a sigmoid activation is applied to the input. This ensures that the output of the network is in the range (0, 1), where values closer to 0 corresponds to the network modeling a given input as class 0, and vice versa.

To assess the performance of the network for a given (known) input, the binary cross-entropy / log-loss function is used, since the output of the sigmoid function may be interpreted as a probability. Log-loss in the binary case is a sum of two parts, each corresponding to the probability of misclassifying that class. Thus, a perfect predictor would have a log-loss of zero.

```
In [7]: matrix2d = a2.flatten()
a2.draw("a2: a2.flatten()")

# FCNN output, post-activation: (a4)
print(f"FCNN output, post-activation: {a4}")

a4 = nn.sigmoid(z4)
print(f"FCNN output, post-activation: {a4}")

print(f"FCNN output, post-activation: {a4}")
print(f"FCNN output, post-activation: {a4}")

FCNN output, post-activation: -4.352467795139205
cross-entropy / log-loss: 1.89802814902024
```

Forward pass done - now what?

Through the forward pass, we have accumulated intermediate values, both in the form of scalars (FCNN pre- and post-activation, log-loss) and matrices (convolved input, maxpool). The point of doing this is to set up for the backward pass, which ultimately leads to a way for us to update the network parameters in such a way as to increase the immediate performance of the network.

Backward pass

The point of the backward pass is to compute the derivative (gradient) of the loss function with respect to the trainable parameters in the network. The intuition for this comes from the fact that, if we are able to tune the network parameters in such a way as to decrease numerical value of the loss function, then we would obtain networks that yield "better performance" (as modelled by the loss function).

NOTE: The notation in the math equations match the naming of the intermediate values of the forward pass. Scroll back and forth if you lose track!

```
In [8]: # Implementation detail, to ensure that the numerical derivatives of the
# functions in the network are used
nn.backward = True

FCNN

For the FCNN layer, we need to calculate the numerical derivative of the loss function w.r.t the weights and the bias of the FCNN, which boils down to smart (or straight-forward) applications of the chain rule.
```

For the weights:

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z_4} \frac{\partial z_4}{\partial w}$$

and the bias

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial z_4} \frac{\partial z_4}{\partial b}$$

Noting that

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial z_1} \frac{\partial L}{\partial z_1} \sigma'(z_1) = \delta$$

where $\frac{\partial L}{\partial z_1}$ is the derivative of the loss function w.r.t the input, and that $\frac{\partial z_1}{\partial z_1} = 1$

We get

$$\frac{\partial L}{\partial w} = \delta \cdot a_1, \quad \frac{\partial L}{\partial b} = \delta$$

NOTE: The "delta" defined in the code from this point on is the delta as defined here, but propagated through the layers of the network, which boils down to smart (or straight-forward) applications of the chain rule.

```
In [9]: delta_from_bias = nn.binary_cross_entropy(y, a4) * nn.sigmoid(z4)
delta_from_weights = delta_from_bias * a3.data
print(f"delta_from_weights: {delta_from_weights.T[0]}")

FCNN bias: -0.9872885318022069
FCNN weights: [-3.626822 -7.86394154 -6.2382538 -1.76069339]
```

Flattening and maxpool

The flattening operation can be reversed without any calculations, as it is a simple reshape of the input vector to the FCNN. The output below shows the reshaped 2x2 derivative.

```
In [10]: delta = Matrix2D(delta_from_weights.reshape(a2.data.shape))
delta.draw("Pre-flattened/post-maxpooled gradient")

Pre-flattened/post-maxpooled gradient
0 1
0 -3.63 -7.86
1 -6.24 -1.76
```

Maxpool

Calculating the gradients through the maxpool operation is a bit more involved. It hinges on the fact that, since all non-max values for each window are discarded entirely, the max values are the only ones which we need to calculate the gradient for. That is, why the maxpool function borrowed from the published example stores the indices of the max values, so we can insert the gradients from the decimated matrix (above) into a matrix with the same size as the original convolved input. The non-max indices will remain zero, but this is again not a problem since these gradients do not affect the loss.

```
In [11]: delta_temp = np.zeros_like(a1.data)
for max_indices, gradient in zip(maxpool_indices, delta_from_weights):
    delta_temp[max_indices] = gradient

delta = delta_temp
matrix2d(delta).draw("Post-convolve/pre-maxpool gradient")

Post-convolve/pre-maxpool gradient
0 1 2 3
0 -6.24 -6.24 -6.24 -6.24
1 -1.76 -1.76 -1.76 -1.76
2 0.00 0.00 0.00 0.00
3 0.00 0.00 0.00 0.00
```

ReLU

The gradients propagate through the ReLU activation, which luckily is much simpler than the maxpool, since ReLU is either constant or linear. Mathematically, the gradient at this layer is simply (again using the chain rule)

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial z_1} \frac{\partial L}{\partial z_1} \sigma'(z_1)$$

```
In [12]: # Derivative through the ReLU activation, which is simple (since relu is either const 0 or linear)
nn.relu(z1).draw("ReLU derivative")
delta = nn.relu(z1.data) * delta
matrix2d(delta).draw("Pre-ReLU gradient")
```

```
Pre-ReLU gradient
0 1 2 3
0 -6.00 -6.00 -6.24 -6.00
1 -1.76 -1.76 -1.76 -6.00
2 0.00 0.00 0.00 0.00
3 0.00 0.00 0.00 0.00
```

Convolution

Calculating the derivative of the loss function w.r.t the convolution kernel does not seem very easy at first, however, by utilizing the fact that convolution can be used to express derivatives, we may use the result

$$\frac{\partial L}{\partial K} = X * \frac{\partial L}{\partial z_1}$$

```
In [13]: # Derivative through the convolution operation
grad_b1 = delta.sum()
grad_kernel = X.convolve2d(Matrix2D(delta))
grad_kernel.draw("Gradient w.r.t. the convolution kernel")

Gradient w.r.t. the convolution kernel
0 1 2 3
0 -57.61 -23.82 -0.89
1 -57.61 -23.82 -0.70
2 -31.87 -9.91 -6.10
```

Updating network parameters

As mentioned previously, the point of going through all these steps is to update the parameters of the network in such a manner as to decrease the numerical value of the loss function. We have now calculated gradients of the loss function w.r.t. the parameters in the network, and since the gradient is a vector pointing in the steepest (positive) direction at a given point on the manifold spanned by the network parameters, we may negate the gradient to find the direction of steepest descent, i.e. do a step of gradient descent. The learning rate, i.e. how far in the direction of steepest descent we perturb the network parameters, is a tunable hyperparameter. Too small, and the learning process takes an infeasibly long time, but too high and the process may overshoot minima and become unstable.

Here η is set arbitrary since "it just works"

```
In [14]: # Update weights and biases for convolution
learning_rate = 0.01

nn.weights += learning_rate * delta_from_weights.T # update weights for the FCNN
nn.from_bias += learning_rate * delta_from_bias # update bias for the FCNN

kernel.data += learning_rate * grad_kernel.data
conv_bias += learning_rate * grad_b1

In [15]: nn.backward = False
print(f"Network parameters: X = {kernel.data.flatten()}, b1 = {conv_bias}, w = {nn.weights[0]}, b2 = {nn.from_bias}")

0.6480985 -1.0259834 0.26448544, b1 = 0.55870414 0.90432152 1.3937715 -0.72684268 -0.76925995
0.65469765 -0.65156923 -0.37699351, b2 = 0.139231121285798644
```

Experimentation and discussion

I really thought that it could be interesting to see how softmax compares to sigmoid as the FCNN output activation function in this case. However, since softmax is a generalization of the sigmoid function for N classes, I realized that it would yield identical performance. Nonetheless, I went ahead and wrote a quick and dirty implementation, and indeed it using the code found in the Automatic training section. As expected, this did not alter performance at all, and so I removed it during cleanup of the code.

A problem that I ran into before doing such a symmetric distribution for class 0, was a problem with the convolved input matrix having all negative values, which then yields an all-zero matrix after the ReLU activation, which in turn degrades performance substantially. Additionally, training is almost certain to halt, since the gradient will also be zero. After reading up on the problem of "dead ReLU's" found two alternatives:

- "Leaky ReLU" with a slight positive gradient for negative values
- "Exponential ReLU" where instead of a hard cutoff at zero, the gradient decays exponentially after zero.

Leaky ReLU is implemented in the Network class, and proved to be effective in reducing ReLU death empirically when I was playing around with training. However, it comes at the cost of another hyperparameter, and as such I decided to simply lower the learning rate as a way to avoid ReLU deaths. Due to this, I did not implement the exponential ReLU.

Automatic training

I am entirely neglecting the holy spirit of Validation here, but this is very much on purpose since this example is so deprived to begin with :)

```
In [17]: kernel = Matrix2D(np.random.uniform(low=-1, high=1, size=(3, 3)))
conv_bias = np.random.uniform()

nn.weights = np.random.uniform(low=-1, high=1, size=(1, 4))
nn.from_bias = np.random.uniform()

nn = Network(learning_rate=0.001, kernel=kernel, conv_bias=conv_bias, weights=from_weights, from_bias=from_bias)

# Create Matrix2D for class 0 and 1
ndim = 6

# Train
# Learning samples to use
classes = np.concatenate((np.ones(n_training_samples // 2), np.zeros(n_training_samples // 2)))
np.random.shuffle(classes)

for y in classes:
    if y == 0:
        X = Matrix2D(uniform_noise(ndim))
    else:
        X = Matrix2D(square(ndim, 3))

    nn.forward_pass(X, y)
    nn.backward_pass(X, y)

# Test
X = Matrix2D(uniform_noise(ndim))
print("Given a matrix without a square:")
nn.predict(X)

# Project a new sample and see how the network performs
X = Matrix2D(square(ndim, 3))
print("Given a matrix with a square:")
nn.predict(X)

Given a matrix without a square:
prob. for class 0 (no square): 0.924885373640924
prob. for class 1 (with square): 0.07511462635907576
Given a matrix with a square:
prob. for class 0 (no square): 0.9254289447516034
prob. for class 1 (square): 0.07457105524839656
```