# Lab 7

## Hüseyin Anil Gündüz

## 2022-07-21

Welcome to the seventh lab, which is focused on convolutions and convolutional neural networks. The first exercise shows how to train CNNs with Keras, while the second exercise is about implementing convolutions for black-and-white images. The third exercise is about computing the gradients of the convolution operator.

## Exercise 1

In this exercise, we will learn how to build CNNs in Keras to classify images.

CNNs are a special type of neural network inspired by the structure of the visual cortex in animals. They can be applied to a wide range of tasks such as image recognition, time-series analysis, sentence classification, etc. Two key features that differentiate CNNs from fully connected nets are:

1. Local connections: Each neuron in a convolutional layer is only connected to a subset of the neurons in the previous layer.
2. Shared weights: Each convolutional layer consists of multiple filters and each filter consists of multiple neurons. All the neurons in a given filter share the same weights but each of these neurons is connected to a different subset of the neurons in the previous layer.

CNNs consistently outperform all other models in machine vision tasks such as image recognition, object detection, etc.

### Classifying hand-written digits

We will be working with is the MNIST dataset (included in Keras). It consists of 28x28 pixels, grayscale images of hand-written digits and their associated labels (0 to 9). The training set contains 60,000 images and the test set contains 10,000. Let's load the data:

```r
library(keras)

mnist = dataset_mnist()

#train_images and train_labels form the training set
train_images = mnist$train$x
train_labels = mnist$train$y

#test_images and test_labels from the test set
test_images = mnist$test$x
test_labels = mnist$test$y
```

The images are encoded as 3D arrays of integers from 0 to 255, and the labels are 1D arrays of integers from 0 to 9.

### Build the network

A CNN typically consists of a series of convolutional and pooling layers followed by a few fully-connected layers. The convolutional layers detect important visual patterns in the input, and the fully-connected layers then classify the input based on the activations in the final convolutional/pooling layer. Each convolutional layer consists of multiple filters. When the CNN is trained, each filter in a layer specializes in identifying patterns in the image that downstream layers can use.

To create a convolutional layer in Keras, call the `layer_conv_2d` function, and specify the the number of filters in the layer (`filters` parameter), the size of the filters (`kernel_size` parameter), and the activation function to use (`activation` parameter).

Pooling layers are used to downsample intermediate feature maps in the CNN. Keras has multiple options for pooling layers, but today we will only use `layer_max_pooling_2d` which takes a `pool_size` argument for the size of the pooling window.

```r
model = keras_model_sequential() %>%
  layer_conv_2d(
    # TODO use 32 filters of size 3x3 and relu activation.
    # Don't forget the `input_shape` parameter
  ) %>%
  layer_max_pooling_2d(
    # TODO use a window of size 2
  ) %>%
  layer_conv_2d(
    # TODO use 64 filters of size 3x3 and relu activation
  ) %>%
  layer_max_pooling_2d(
    # TODO use a window of size 2
  ) %>%
  layer_conv_2d(
    # TODO use 64 filters of size 3x3 and relu activation
  ) %>%
  layer_flatten() %>%
  layer_dense(
    # TODO use 64 units with relu activation
  ) %>%
  layer_dense(
    # TODO use 10 units and softmax activation
  )
```

Let's take a look at what we've built so far:

```r
summary(model)
```

You can see that the output of every `layer_conv_2d` and `layer_max_pooling_2d` is a 3D tensor of shape (`height, width, channels`). For example, the output of the first layer is a tensor of shape (26, 26, 32). Note that the width and height dimensions shrink as you go deeper in the network. The number of channels is controlled by the `filters` parameter of the convolutional layers. Also, as you can see, the (3, 3, 64) outputs are flattened into vectors of shape $576 = 3 \cdot 3 \cdot 64$ before going through two dense layers.

### Data preprocessing

Before training this model, we have to preprocess the data by scaling the inputs. Recall that the inputs are integer arrays in which the elements take values between 0 in 255. It is standard practice to scale the inputs so that the elements take values between 0 and 1. This typically helps the network train better.

```r
# Reshape and rescale train_images and test_images.
train_images = array_reshape(train_images, c(60000, 28, 28, 1))
train_images = train_images / 255

test_images = array_reshape(test_images, c(10000, 28, 28, 1))
test_images = test_images / 255
```

### Compile,train and evaluate the model

```r
# TODO compile the model using the rmsprop optimizer,
# the sparse categorical cross-entropy loss, and the accuracy metric.
```

```
# TODO Train the model for 5 epochs with a batch size of 64,
# and use 20% of the images as validation

evaluate(
  # TODO Evaluate the model on the test data
)
```

The accuracy of our model on MNIST is quite good!

## Exercise 2

In this exercise we are going to implement convolution on images, without worrying about stride and padding, and test it with the Sobel filter. There are two Sobel filters: $G_x$ detects horizontal edges and $G_y$ detects vertical edges.

$$G_x = \begin{vmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{vmatrix} \qquad G_y = \begin{vmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{vmatrix} = G_x{}^T \tag{1}$$

Can you explain why and how these filters work?

In order to get the image $E$ with the edges, we convolve $G_x$ and $G_y$ with the input image $I$, to obtain the degree of horizontal and vertical "borderness" of each pixel. We then combine these values (separately for each pixel) with an L2 norm:

$$E = \sqrt{(G_x * I)^2 + (G_y * I)^2} \tag{2}$$

As a reference, this is the result we want to obtain:

```
library(OpenImageR)

# we only keep the first channel, as the image is already black and white
img = readImage("einstein.jpg")[,,1]

# NB: you can resize the image to speed up subsequent operations
# img = resizeImage(img, floor(ncol(img) / 4), floor(nrow(img) / 4))

# define the two filters
sobel_x = matrix(c(-1, 0, 1, -2, 0, 2, -1, 0, 1), nrow = 3)
sobel_y = matrix(c(-1, -2, -1, 0, 0, 0, 1, 2, 1), nrow = 3)

apply_sobel <- function(conv_fn) {
  # convolve both filters on the image
  conv_x = conv_fn(img, sobel_x)
  conv_y = conv_fn(img, sobel_y)

  # combine the two convolutions
  conv = sqrt(conv_x^2 + conv_y^2)

  # normalize maximum value to 1
  conv / max(conv)
}


result = apply_sobel(convolution)
grid::grid.raster(result)
```

We now implement our version of convolutions. For an input matrix $\mathbf{X}$ of size $r(\mathbf{X}) \times c(\mathbf{X})$ and a kernel $\mathbf{K}$ of size $r(\mathbf{K}) \times c(\mathbf{K})$, the result of the convolution is $\mathbf{Y} = \mathbf{K} * \mathbf{X}$ with $r(\mathbf{Y}) = r(\mathbf{X}) - r(\mathbf{K}) + 1$, $c(\mathbf{Y}) = c(\mathbf{X}) - c(\mathbf{K}) + 1$, and elements:

$$y_{ij} = \sum_{k=1}^{r(\mathbf{K})} \sum_{l=1}^{c(\mathbf{K})} x_{i+k-1,j+l-1} \cdot k_{kl} \tag{3}$$

for $1 \leq i \leq r(\mathbf{Y})$ and $1 \leq j \leq c(\mathbf{Y})$.

You now have to implement a function that computes $y_{ij}$ given the image, the kernel, $i$ and $j$.

```
compute_convolution_at_position <- function(i, j, image, kernel) {
  # TODO compute and return the convolution at row i and column j with the formula above
}


our_conv <- function(image, kernel) {
  result_rows = (
    # TODO compute the number of rows of the result
  )

  result_cols = (
    # TODO compute the number of columns of the result
  )

  # perform the convolution
  vec = apply(expand.grid(1:result_rows, 1:result_cols), 1, function(pos) {
    compute_convolution_at_position(pos[1], pos[2], image, kernel)
  })

  # reshape to a matrix
  matrix(vec, nrow = result_rows, ncol = result_cols)
}

our_result = apply_sobel(our_conv)
grid::grid.raster(our_result)
```

If you did everything correctly, this image should match the image above.

## Exercise 3

Recall that the convolution $\mathbf{Y} = \mathbf{K} * \mathbf{X}$ has elements

$$y_{ij} = \sum_{k=1}^{r(\mathbf{K})} \sum_{l=1}^{c(\mathbf{K})} x_{i+k-1,j+l-1} \cdot k_{kl} \tag{4}$$

Now consider $\mathbf{X}$ and $\mathbf{Y}$ to be the input and output of a convolutional layer with filter $\mathbf{K}$. For simplicity, we focus on a single channel; actual convolution layers in CNN perform this operation several times with different learnable filters.

Imagine this convolution is a hidden layer of the neural network, with $\mathbf{X}$ being the input from the previous layer, and $\mathbf{Y}$ the pre-activation output to the next layer. Then, we can define the loss function in terms of $\mathbf{Y}$, i.e. $\mathcal{L} = f(\mathbf{Y})$, where $f$ includes the activation, all the following layers, and the classification/regression loss.

Show that:

$$\frac{\partial \mathcal{L}}{\partial k_{kl}} = \sum_{i=1}^{r(\mathbf{Y})} \sum_{j=1}^{c(\mathbf{Y})} \frac{\partial \mathcal{L}}{\partial y_{ij}} \cdot x_{i+k-1,j+l-1} \tag{5}$$

Then show that

$$\frac{\partial \mathcal{L}}{\partial x_{ij}} = \sum_{k=L_k}^{U_k} \sum_{l=L_l}^{U_l} \frac{\partial \mathcal{L}}{\partial y_{ab}} k_{kl} \tag{6}$$

With

$$a = i - k + 1 \tag{7}$$
$$b = j - l + 1 \tag{8}$$
$$L_k = \max(1, i - r(\mathbf{X}) + r(\mathbf{K})) \tag{9}$$
$$L_l = \max(1, j - c(\mathbf{X}) + c(\mathbf{K})) \tag{10}$$
$$U_k = \min(r(\mathbf{K}), i) \tag{11}$$
$$U_l = \min(c(\mathbf{K}), j) \tag{12}$$

As you can see, the gradient of the input is obtained by convolving the same filter with the gradient of the output, with some care at the borders.

Hint: it is easier to analyze convolutions in one dimension with a small example, then generalize the result to two dimensions and arbitrary filter/image size.

Now, write a function that computes $\partial \mathcal{L}/\partial x_{ij}$, with $\mathcal{L} = \sum_{i,j} y_{ij}^2$ and $\mathbf{K} = G_x$.

```r
conv_gradient_wrt_input <- function(dloss_dy, kernel) {
  image_rows = (
    # TODO compute the number of rows of the original image
  )

  image_cols = (
    # TODO compute the number of columns of the original image
  )

  vec = apply(expand.grid(1:image_rows, 1:image_cols), 1, function(pos) {
    i = pos[1]
    j = pos[2]

    # TODO compute and return the gradient at row i and column j
  })

  matrix(vec, nrow = image_rows, ncol = image_cols)
}


result = our_conv(img, sobel_x)
input_gradient = conv_gradient_wrt_input(2 * result, sobel_x)
grid::grid.raster((
  input_gradient - min(input_gradient)
) / (
  max(input_gradient) - min(input_gradient)
))
```

We can verify this gradient is correct for a single pixel with finite differences:

```
eps = 1e-6

i = floor(runif(1, 1, nrow(img) + 1))
j = floor(runif(1, 1, ncol(img) + 1))

# add epsilon to position i,j and convolve
img[i, j] = img[i, j] + eps
conv_pos = our_conv(img, sobel_x)

# remove epsilon to position i,j and convolve
img[i, j] = img[i, j] - 2 * eps
conv_neg = our_conv(img, sobel_x)

# undo modification to the image
img[i, j] = img[i, j] + eps

# compute the difference of the losses
# NB: we sum the differences to get a more accurate result
empirical_gradient = sum(conv_pos^2 - conv_neg^2) / (2 * eps)

# compare empirical and analytical gradients
c(empirical_gradient, input_gradient[i, j])
```

If you did everything correctly, these two numbers should be the same.

Now, can you guess what image *maximizes* the loss we just defined? We can find this through gradient *ascent*:

```
maxim = matrix(runif(81), nrow = 9)

losses = lapply(1:250, function(i) {
  # TODO convolve `sobel_x` with `maxim` and compute the loss

  # TODO compute the gradient of the loss, and modify `maxim` accordingly

  loss
})

plot(1:length(losses), losses)

grid::grid.raster(
  (maxim - min(maxim)) / (max(maxim) - min(maxim)),
  interpolate = FALSE)
```