

1 Convolutional Networks

By: Jussi Martin

1.1 Preface

This blogpost is written for the *Image and Video Statistics* study group held by the AIHelsinki Academia. It is based on the Chapter 9 of the book *Deep Learning* by Ian Goodfellow, Yoshua Bengio and Aaron Courville, which is available online at <http://www.deeplearningbook.org>.

We will assume the reader to be familiar with basic components of artificial neural networks, such as weights and biases. Some mathematical preliminaries are introduced in the text. Our goal is to introduce basics of convolutional neural networks and explain some of their properties.

Convolutional networks are particularly useful in image classification, since they can learn to be insensitive to local translations in the input. For example, convolutional network can learn to classify image of a face correctly despite individual faces having slight variation in the positions of facial features.

Besides the classificational aspects, convolution is also effective in the sense that it uses fewer connections and shared parameters, which reduces the computational cost and memory requirements and improves the statistical efficiency. In deep networks this difference is significant.

1.2 Definition

We begin by introducing the mathematical notion of convolution. Let $x(t)$ and $w(t)$ two functions then we define their convolution $x * w$ as

$$(x * w)(t) = \int x(a)w(t - a)da.$$

In practical applications x could be some signal depending on the time t and w a filter applied to it.

Obviously the data for which the convolution is applied usually is obtained only as a discrete input, in these cases we define the convolution by using summation:

$$(x * w)(t) = \sum_{a=-\infty}^{\infty} x(a)w(t - a).$$

Moreover, the data might be two dimensional, which is the case with images, then we use the following definition:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n)$$

where I and K represent two images and i and j are the pixel coordinates.

In practice the input grid is finite, the filter is defined on a smaller grid and usually only applied when its coordinates are fully contained in the input grid. One special property of convolution is that it is commutative, that is

$$(K * I) = (K * I).$$

Which is relatively easy to derive. Many machine learning libraries actually defined convolution as

$$S(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n),$$

which is also known as cross-correlation or convolution without kernel-flipping. It is not commutative, but from practical point of view this does not make any difference since commutativity is not needed and the network can learn its parameters with respect to either of these convolutions.

We say that a neural network is convolutional if it uses convolution in at least one of its layers.

1.3 Convolutional Networks Structure and Output

Here we give an example layout for convolutional networks structure. We will also consider what type of outputs convolutional network can have. Individual layers are presented later.

Typically the first layer in a convolutional network is convolutional. It is followed by nonlinearity. Often after this there is some type of pooling layer, usually max pooling. This type of structure can then be repeated multiple times. The final layers are fully connected and responsible for the one hot encoding of the output.

In the previous scheme the output is an array with one entry equal to one and others zero, representing a membership of a given class. Other optional output types could be for example, an array with estimated probabilities that the image has an object in a given class, a tensor (tensors will be explained in the next section) which tells estimated probabilities that a given pixels is part of an object in a given class.

1.4 Tensors and Convolutional Layer

We adopt the Deep Learning book's convention of defining tensors as multidimensional arrays of real numbers.

Namely, 0-D tensors are just real numbers, 1-D tensors are arrays

$$T = (T_1, \dots, T_n)$$

of real numbers, 2-D tensors are arrays

$$T = ((T_{1,1}, \dots, T_{1,n}), \dots (T_{m,1}, \dots, T_{m,n}))$$

of arrays of real numbers (with mutually equal length), and so on.

Basically our tensors can be viewed as D-dimensional grids of real numbers.
Example in 2-D:

$$((1, 2, 3), (4, 5, 6), (7, 8, 9)) \simeq \begin{array}{c|c|c} 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \end{array}$$

In general case the sizes of the axes do not need to match.

1.5 (submerge to previous, when ready)

Typically convolution network has several channels which may correspond to different colors in the image or they may originate from same input image which is processed in different ways in parallel.

In the following we assume that that input grid has spatial shape of $a \times a$ and the convolution kernel K has the spatial shape of $b \times b$. Both of these are positive integers and

$$a = bn \quad \text{for some } n > 1.$$

Furthermore, the shapes are independent of the channel been used.

Output of the convolutional layer at pixel (j, k) of the channel i is

$$Z_{i,j,k} = \sum_{l,m,n} V_{l,j+m-1,k+n-1} K_{i,l,m,n}$$

where $K_{i,l,m,n}$ is the weight of the connection from input in channel j to output in channel i , both having row index k and column index l . The sum is taken over the index set

$$[1, \dots, ?] \times [1, \dots, ?] \times [1, \dots, ?]$$

output $Z_{i,j,k}$ is only defined when the convolution is valid. This means that
With stride:

$$Z_{i,j,k} = c(K, V, s) = \sum_{l,m,n} V_{l,(j-1) \times s + m, (k-1) \times s + n} K_{i,l,m,n}$$

*Tiled (optional?)

*Local connections (optional?)

*Formulas, etc.

... traditional fully connected layer can be seen as other extreme case, with zero stride and kernel grid size been equal to the input grid size.

1.6 Nonlinearities

Some of the typical nonlinearities used in convolutional networks are:

Rectified linear unit also known as *ReLU*, defined by

$$f(x) = \max(x, 0).$$

Logistic function also known as *sigmoid*, defined by

$$f(x) = \frac{1}{1 + \exp(-x)}.$$

Hyperbolic tangent function also known as *tanh*, defined by

$$f(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}.$$

A nonlinearity is applied after convolutional layer. Its output is defined as

$$T_{i,j,k} = f(Z_{i,j,k} + b_{i,j,k})$$

where $Z_{i,j,k}$ is the output of the convolution and $b_{i,j,k}$ is the bias unit, which is often shared between units in the same channel, that is

$$b_{i,j,k} = b_i \quad \text{for all } j \text{ and } k$$

where j and k are the pixel indices and i is the channel index.

In some case however, for example if objects in a set of images are known to be in the center of the images, it may be more convenient to allow the biases to vary with respect to location.

1.7 (Max) Pooling

Typically pooling is used for down sampling, which means that the grid size of the data is reduced.

$$\begin{array}{|c|c|c|c|} \hline 0 & 1 & 3 & 9 \\ \hline 4 & 8 & 6 & 2 \\ \hline 7 & 4 & 2 & 0 \\ \hline 4 & 3 & 5 & 1 \\ \hline \end{array} \mapsto \begin{array}{|c|c|} \hline 8 & 9 \\ \hline 7 & 5 \\ \hline \end{array}$$

Example in 2-D: max pooling with 2×2 kernel and stride = 2.

Another benefit of max pooling with down sampling is that it is not so sensitive to small translations. If we take to grids, which are small translations of each other, and compare the pooled grid, we see that many of the values are the same.

Pic: Without Shift, Shifted

1.8 Zero Padding

One drawback of using convolutional layers is that they shrink size of the data grid, which opposes limit to the depth of the network. This can be compensated by adding artificial zeros to the boundary of the input layer. This way we can prevent the sizes of the layers shrinking and obtain deeper networks.

The method described above is known as zero padding. One way to do it is to add just enough zeros that the layer sizes stays the same. Other more

extreme option is to add so many zeros that the convolution can visit every neuron k times, where k is the size of the convolution kernel.

Pic: 2-D, original layer and layer with minimal padding

Pic: 2-D, original layer and layer with maximal padding

One may notice that zero padding makes the outputs near the boundary to have less information from the inputs, especially if there is maximal amount of padding used. This may be compensated to some extent if the biases are not shared.

1.9 Receptive Field

For a given node the nodes in the previous layers which are directly connected to it are known as receptive field.

Although the receptive field for single convolutional layer is small due to sparse connections, it grows when there are more layers added to the network and it can even cover the whole input grid if the network has enough layers.

1.10 Learned Invariances and Feature Detectors

Examples, Explanation (no pics)

1.11 Further Topics (Optional ?/ rename as Training?)

In this last section we will mention briefly some of the topics presented in the Chapter 9 of the Deep Learning book that we have not yet covered.

One thing we did not yet cover is the training of convolutional network. It turns out that ... (in order to keep this text concise ... OR: since we assume the reader to be familiar with basics of regular networks ...) we skip this.

An other topic of possible further interest is the use of recurrent structures in a convolutional network, for this we also refer to consult the Deep Learning book.