

9 DeepLearning

Our treatment of neural networks in the last chapter was basically state of the art in the 1980 and 1990s. However, neural networks have recently receiving considerable attention under the name of deep learning. Deep learning basically refer to neural networks with many layers, which will of course lead to more complex models and hence it is likely that they can outperform simpler models for complex tasks. While this was already clear in the 1990s, we have only been able to train such networks in more complex tasks in the last few years. This advancement is due to several factors, but specifically through the availability of large supervised datasets, the enormous increase of computational power in particular with the help of graphical processor units (GPUs) that are applicable to matrix operations, and also thanks to the development and of some more sophisticated applications of regularization techniques. Part of making such deep networks successful for image classification is the resurgence of convolutional networks which are now an integral part of deep networks and with which we start our discussion.

9.1 Filters and 1d convolution

Let us start by running a simple linear perceptron on a pattern identification problem. This consists of some training vectors with ten features. These training vectors can also be seen as a time series of some measurements where each feature value represents a time point t . An example training set is shown in Fig. 9.1. These training pattern are generated in the following way. All the features values are uniformly drawn between 0 and 0.3 except the last one which is always 1 to represent the bias. However, 10 of the training patterns have specific values for three time points, namely $x(2) = 1$, $x(2) = 0.5$, and $x(2) = 1$.

The program which generates the training examples and trains the perceptron is shown below.

```
from pylab import *

np=100
X=0.3*rand(11,np); X[10,:]=1
a=array([1.0,0.5,1.0])
X[2:5,0:10]=array([a,a,a,a,a,a,a,a,a,a]).T

Y=zeros(np); Y[0:10]=1

# model specifications
Ni=11; No=1;
```

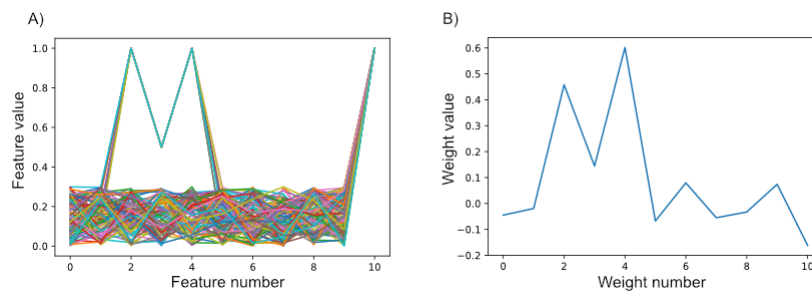


Fig. 9.1 Example of a pattern identification problem. A) the training patterns are mainly noise except for some pattern that have consistently larger values for features 2, 3 and 4. B) The weights of a linear perceptron after training on the patterns on the left. The weights represent a filter for the M-shaped pattern at position 3.

```
#parameter and array initialization
Ntrials=1000
wo=rand(No, Ni); dwo=zeros(wo.shape)
do=zeros(No)

for trial in range(Ntrials):
    y=wo@X # prediction for all samples
    do=(Y-y) # delta term
    dwo=0.9*dwo+do@X.T # momentum
    wo=wo+0.01*dwo # weight update
    error[trial]=sum(abs(Y-y)>0.2)/np

plot(wo[0, :])
figure(); plot(X)
```

We can learn two important point from this example. The first is that weights represent (learned) feature detectors. This is an important way of thinking about these weights, and plotting the weights will be a good way of visualizing some result of training. Of course, these weights look for patterns at specific locations. To detect the same pattern at a different location we would have to train a new node with examples of the pattern at this new location. This seems to be a lot of training and a lot of weights. In many cases we want to look for patterns that can occur at any places. For example, if we are looking for a coffee mug than this can be placed at many locations in our visual field. Having feature detectors for specific locations does there seem to be a inefficient way of doing pattern search.

A better way to do location-invariant pattern recognition is to use convolutions. Let us illustrate the idea before formally defining this operation. For this we will use again the specific M-shaped pattern in a series of data points. Let say that we know that we look for this M-shaped pattern, and to detect this we use a **filter** f which is a short vector describing this pattern, namely $f(0) = 1$, $f(1) = 0.5$, and $f(2) = 1$. To apply this filter we simply multiply this filter with the corresponding signal. Let's say

the signal is the vector $\mathbf{x} = (0010.5100000)$. Let's place the filter at the beginning of this signal and let's multiply the components and then add them together. This results in the value of 1. Placing the filter starting at the second position and calculating this sum of the products gives a value of 1.5. After placing the filter at each possible location we get the sequence (112.511000) . The operation of multiplying the filter with the signal and adding the terms gives us some form of overlap measure between the pattern and the filter. A large value indicates that the pattern is present and where it is present. The operation of multiplying and adding while shifting the filter is called a **convolution**. The resulting filtered signal represents how much and where a pattern is present in the original signal. For example, if we have a signal with a shifted location, say $\mathbf{x} = (0000010.5100)$, we get a filtered signal (000112.511) .

You might have noticed that the filtered signal is smaller than the original signal. Indeed, we are losing some components in the filtered signal proportional to the size of the filter as the filter has a finite extent and can not be placed with the beginning of the filter at the end of the original signal. There are different ways of handling this. For example, we could just add some zeros to the original signal. This is called zero **padding**. Or we could repeat some of the previous entries such as continuing with the start of the signal when we reach the end. Or we can just accept that the resulting filter signal is slightly smaller. Sometimes we are even interested in getting smaller versions of the signal which somewhat represent a compressed representation. One way of doing this is to not just shift the filter by one step in each iteration of the convolution, but use larger steps such as 2 or more. This is called a **stride**. A stride value of two would half the size of the filtered signal.

Before we move to higher dimensions, let us formalize somewhat our discussion and define the above described convolution for a 1-dimensional discrete signal mathematically as

$$(f * x)(t) = \sum_{t'=0}^T f(t')x(t+t'). \quad (9.1)$$

We can also formulate this for a continuous signal. Since we can then not start the signal at zero, as the signal runs formally from $-\infty$ to $+\infty$, it is more common to place the filtered value at the center of the filter

$$(f * x)(t) = \int_{-\infty}^{\infty} f(t')x(t-t')dt'. \quad (9.2)$$

Note that we consider here a filter that is reversed compared to the discrete definition. This is historically the more common way to define the convolution operation. A function like the filters appearing in an integral are mathematically called a kernel.

9.2 2-d convolution and image processing

In the last section we looked at convolutions for 1-dimensional data. It is straight forward to generalize this to 2-dimensional data. Mathematically, a complete n-dimensional convolution could be written as

$$(f * g)(x_1, x_2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x'_1, x'_2)g(x_1 - x'_1, x_2 - x'_2)dx'_1dx'_2, \quad (9.3)$$

where we used the letter g for the signal, which is now an n dimensional array. Mathematically, this corresponds to a tensor. Note that the difference between an tensor and a n -dimensional array is usually that the mathematical construct of a tensor (or a matrix in 2 dimensions or a vector in 1 dimensions) represent the structure of an array together with some operations such as adding and multiplying such object.

In the following we will illustrate higher dimensional convolutions in a discrete 2-dimensional case, that of processing greyscale images. Let us first make a simple example of a greyscale image, say a matrix of the 4x4 matrix where each entry corresponds to a pixel with a grey scale value between 0 and 10. An example image is shown on the left in Fig. 9.2.

$$\begin{pmatrix} 7 & 8 & 3 & 2 \\ 9 & 6 & 2 & 3 \\ 8 & 9 & 3 & 1 \\ 5 & 7 & 2 & 2 \end{pmatrix} \odot \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix} = \begin{pmatrix} 2 & 9 & 0 \\ 2 & 10 & 1 \\ -3 & 11 & 2 \end{pmatrix}$$

Fig. 9.2 Example of edge detector.

We apply to it a filter with 1s on the left column and -1s on the right column. Applying this filter means first to overlay the filter with the image in the upper left corner, multiplying the overlapping components and adding them all up. This results in the value two which we noted in the upper left corner of the resulting filtered image. Next we move the filter one position to the right and perform the same operation of component-wise multiplication and addition, which gives the value of 9. We repeat this procedure until we have placed the filter over all possible places in the image.

It is interesting to observe the resulting filtered image. What stands out in this image is a horizontal column in the middle. This would correspond to a dark line which corresponds to an edge in the original image. The filter is therefore a form of an edge detector, specifically for horizontal edges. An example on a real image is shown in Fig. 9.3. The first figure shows a color image in JPG format that we can read and display with the following Python code

```
from pylab import *
img=imread('photo.jpg'); imshow(img)
print('image_shape', img.shape)
```

This image has 720*1280 pixels, where the color is represented in three channels. We pick the first channel, which is displayed as the second image, and then apply the simple edge filter which results in the third image. This filtered image shows mainly horizontal lines. In a similar way we can also design edge detectors for vertical edges.

Of course, our filter is very small and only shows edges with significant changes between two consecutive edges. There are therefore better designs of edge detectors, such as the Canny edge detector. These techniques combine such gradient filters with Gaussian smoothing and removal of some spurious cases. Also, a continuous version

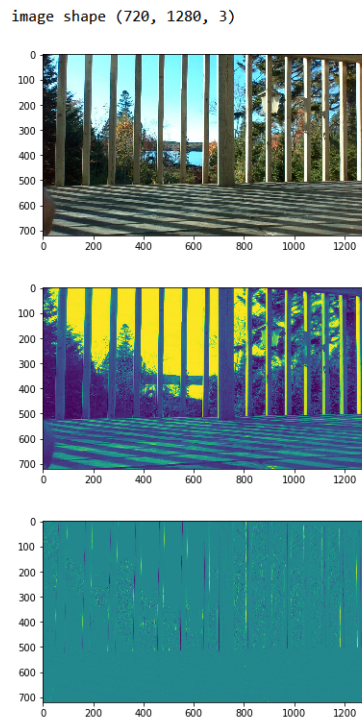


Fig. 9.3 Example of simple edge detector on the first component of a JPG image.

of edge filters is for example described by Gabor functions such as the ones shown in Fig. 9.4a and b. A Gabor function is described by a sinusoidally-modulated Gaussian,

$$f(u, v) = e^{-\frac{u^2 + \gamma v^2}{2\sigma^2}} \cos\left(\frac{2\pi}{\lambda}u + \varphi\right). \quad (9.4)$$

The example of a 64^2 pixel filter with parameters $\gamma = 0.5$, $\sigma = 10$, $\lambda = 32$, and $\varphi = \pi/2$ is shown in Fig. 9.4a. This filter can also be rotated with a rotation matrix

$$\begin{pmatrix} x \\ y \end{pmatrix} \leftarrow \begin{pmatrix} \cos(\varphi) & \sin(\varphi) \\ -\sin(\varphi) & \cos(\varphi) \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \quad (9.5)$$

as shown in Fig. 9.4b for $\varphi = \pi$. Interestingly, such functions describe some of the neurons in the primary visual cortex of primates. Detecting edges seems therefore a good first step to process images, a fact that we will encounter again in later discussion.

Attachments area

9.3 Convolutions with tensors

We already discussed the case of a 1-dimensional convolution and a 2-dimensional convolution. It is easy to generalize this to n -dimensional data. Mathematically, an n -dimensional convolution could be written as

A. Gabor function with $\phi = \pi/2$ B. Rotated version of A

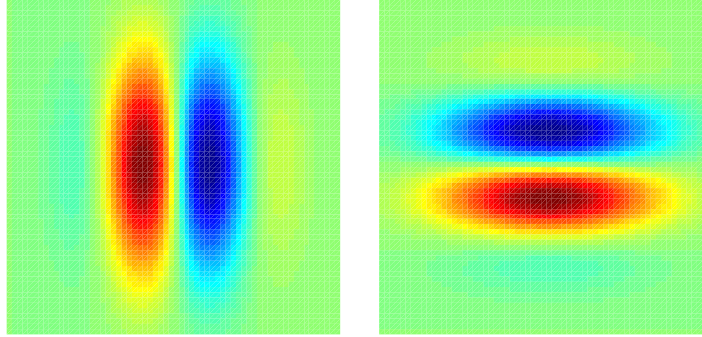


Fig. 9.4 Example of Gabor functions for (a) vertical and (b) horizontal edge detection.

$$(f * g)(x_1, \dots, x_n) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x'_1, \dots, x'_n)g(x_1 - x'_1, \dots, x_n - x'_n)dx'_1 \dots dx'_n, \tag{9.6}$$

where we used the letter g for the signal, which is now an n -dimensional array, which is mathematically a tensor. For example, a 3-dimensional tensor would be a cube, and convolving them with another 3-dimensional tensor (cube) would result in another 3-dimensional tensor (cube).

However, such higher dimensional operations also do not always make sense even though our input data might be a tensor. A prominent example is a color image. A color image can be represented with three 2-dimensional intensity maps, one representing the red components, one the green components, and one the blue components (see Fig. 9.6). This is what is behind the RGB representation of a color image. While these data can be represented by a 3-dimensional tensor, convolving over the color channels does not really make sense. Indeed, convolutions are only useful when we are looking for patterns in the relation between neighboring entries in a tensor. While some objects could be color sensitive, it is more common to just add or average the information from the different color channels. We can do this in two ways. Before we used already one in which we just averaged over the color channels that resulted in a grey-scale image. Another way is to allow first the filtering of each channel and then average. These two possible ways lead to different results if we allow different filters for each channel. Mathematically, the second way of averaging over $(n-m)$ channels after applying m -dimensional convolutions of the $(n-m)$ channels is given by

$$(f * g)(x_1, \dots, x_m) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x'_1, \dots, x'_n) \dots \tag{9.7}$$

$$\dots g(x_1 - x'_1, \dots, x_m - x'_m, x'_{m+1}, \dots, x'_n) dx'_1 \dots dx'_n, \tag{9.8}$$

Thus, if we have one edge filter, say a horizontal edge detector, and apply this to a color image, we get one grey image that averages the lines of the r, g, and b channels of the color image. If we are using another filter, say vertical edge detector, then we would get a second filtered image. The different filtered images can be viewed as **different**

channels generated by a **filter bank**. Of course, we can use filter banks that have many more filters in them. This is illustrated in the following figure.

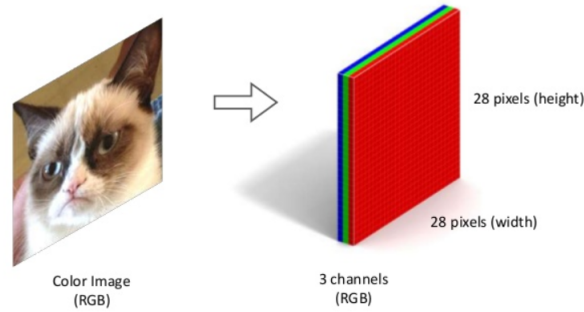


Fig. 9.5 Representation of color images with color channels.

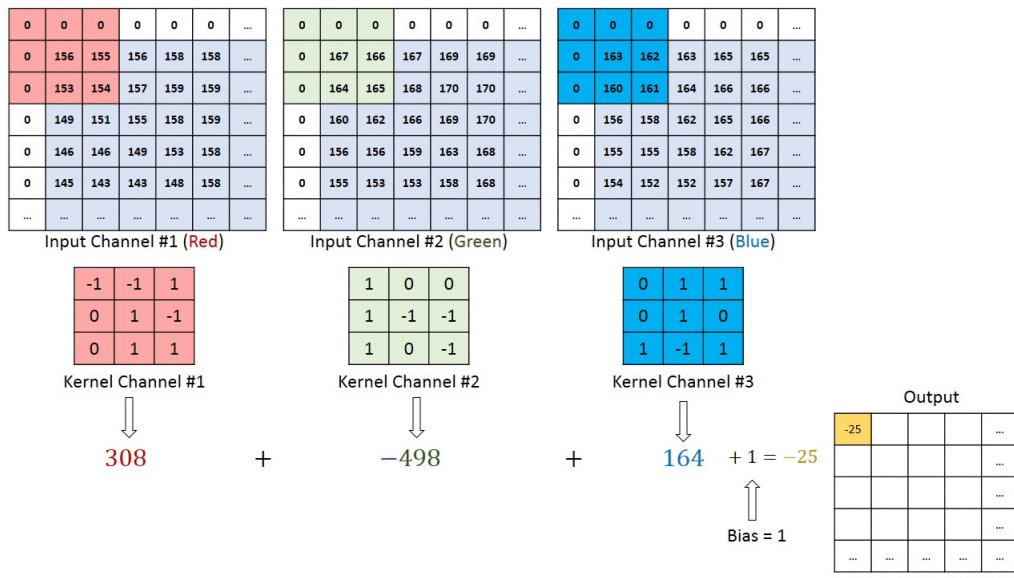


Fig. 9.6 Convolutions with tensors and channels.

9.4 CNN

The idea of designing position invariant feature detectors in neural networks was first described by Fukushima in 1980. Fukushima worked at this time for NHK (the Japanese public broadcaster) together with physiologists as NHK was interested to understand the mechanisms of human vision. It was well known since the early 1960s

form the experiments by Hubel and Wiesel that some neurons in the primary visual cortex (the first stage of visual processing in the cortex) are edge detectors as already mention above.

Edge detectors are also the workhorse of computer vision, and we discussed in the first section how such filters are implemented with convolutions. The neural networks that we discussed before had to learn individual weights to each pixel location. Even if this network would learn to represent an edge detector, such detectors have to be learned for each location in an image since edges could usually appear in all locations. So another way of thinking about convolution is that a neuron (specific filter) is applied to every possible location in the image. This leads us to a **convolutional neural network (CNN)**.

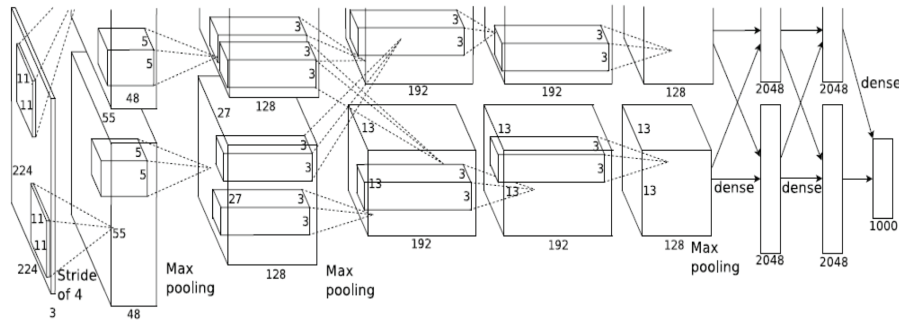


Fig. 9.7 A famous implementation (so called AlexNet) of a convolutional neural network for image classification Representation (Krizhevsky, Sutskever, Hinton, 2012).

A famous example called AlexNet is shown that was used to classify a big database of images from many different classes is shown in Fig. 9.7. This network takes three dimensional images (e.g. RGB values of pixels) and applies a layer of filters to it. This specific network divided the pathways though the network into to stream in order to facilitate the computation with two GPUs, but at this point it is only necessary to look at one stream. In this example, the inputs are of size 224x224 with three channels (RGP). We then perform the following layered computation on this input:

convolution The first level filters are of size 11x11, so that the resulting image is of size 55x55. This layer also consists of 48 different filters so that we end up with 48 filters.

gain function At this stage we are using a gain (or activation) function that is not explicitly shown in the figure but is sometimes indicated as separate computing step. It is common to use the RELU activation function for reasons discussed later.

pooling If we would only apply convolutions on convolutions, then we would end with filtered images of filtered images. However, what we really hope to achieve are high level representation such as nodes that represent class labels. Such labels are likely not to depend on individual pixels and rather represent a highly compressed summary of an image. To help with this we add **pooling layers** after each convolution layer. A pooling operation is usually just taking the average or the maximum of the

responses in a certain area of the filtered image, thus compressing the images down by only considering the average or maximal feature represented by the filter in this area.

These three steps are the typical components of a convolutional layer of a CNN, and these steps are repeated to build a deep convolutional network. At the end we use a **fully connected layer (MLP)** to gather all the information and make the final classification based on the features extracted by the network.

While we have just outlined the operation of this network, an important part of using this network is of course the training. Indeed, for our further discussion it is good to realize that the filters are not chosen by hand but are learned from examples. This training was trained with the back propagation algorithm. This is fairly straight forward except when back propagating through the pooling layers. One approach is thereby to give all credit for the error (average pooling), or just change the winning unit (max pooling).

9.5 Tensorflow

To implement deep networks we use the Google's Tensorflow package. At this time you should follow the tensorflow tutorial

https://www.tensorflow.org/get_started/get_started

where some of the structure of Tensorflow is described. The please follow the MNIST tutorial

https://www.tensorflow.org/get_started/mnist/beginners