

Introduction

NN (perceptron) depuis (60-70?) aujourd'hui, bc d'applications, CNN Notre but est alors de comprendre et savoir comment mettre en oeuvre, NN qui s'exécute dans un temps raisonnable

Table of contents

1 Classical Neural Networks (Multilayer Perceptron)	2
1.1 Training phase	2
1.2 Optimization using gradient descent	3
1.3 Activation function	3
1.4 Softmax	3
1.5 Forward propagation	3
1.6 Backward propagation	3
2 Convolutional Neural Networks	4
2.1 The convolutional neuron	5
2.1.1 Where is the convolution?	6
2.1.2 Multiple channels	8
2.1.3 The stride	8
2.1.4 Zero padding	9
2.2 The pooling 'neuron'	9
2.3 Backward Propagation	9
3 Python Implementation	9
4 Conclusion	9

1 Classical Neural Networks (Multilayer Perceptron)

NN can be used for solving a wide array of problems in computer science, arguably the most common of which lies in computer vision. As such, for the purpose of demonstration, we will explain the structural motivation of a typical NN, the Multi-Layered Perceptron (MLP), in the context of computer vision.

Say for instance you have an image containing an animal and you want to identify what animal it is, or to classify the image using a NN. The network itself can then be considered a function with a number of inputs equal to the number of pixels in the image times three in the case of an RGB. The number of outputs is in this case the number of animal types. This is obviously an extremely complicated function and finding the exact function is near impossible. A more intelligent approach is to approximate f with composition of several easily computed functions, such as linear functions (affine) (Figure 1). The structure is loosely based on neuroscience, hence neurons, weights, biases. A simple network might consist of two such layers with Y number of neurons each. This means (Y) number of adjustable parameters! With the structure defined, we have yet to determine the parameters. This is where the training phase enters.

1.1 Training phase

The idea is simple. We start out with a large dataset with classified examples, called a training set and we adjust our network to fit these data. The hope is then that the network will be able to classify new images with these same parameters. In order to quantify the performance of the NN during the training phase, we introduce a cost function. Knowing the desired output with the corresponding input, we look to minimize this cost function by adjusting the parameters of the NN i.e. its weights and biases. This can be seen as an optimization problem with thousands of inputs (in our case, Y weights and biases). For demonstration purposes, we will explain the concept of gradient descent, an intuitive and common optimization algorithm.

1.2 Optimization using gradient descent

We have a function f that we want to minimize while avoiding calculating all its possible values. The gradient of f is an operation that tells us which inputs to adjust and by what weight to adjust them in order to increase f by as much as possible in a given point. Visually, it can be considered the direction in the input space that increases the output the most. In our case, we want to decrease f by as much as possible so, assuming f is continuous, we take the negative value of the gradient. As such, the gradient descent algorithm consists of, at each iteration, calculating f and $\text{grad}f$, adjusting the input by “stepping in the direction of $\text{grad}f$ ” (Figure 2) and repeating until we have found a local minimum or until the value of f is sufficiently low.

In the context of training a neural network, we have yet to introduce a way to calculate this crucial gradient of the cost function. To do so, we first need to determine a mathematical expression of the neural network itself. In fact, the network typically consists of more than just linear functions.

1.3 Activation function

There are a few different motivations behind an activation function. Firstly, it introduces non-linearity to a problem that is manifestly not linear. For optimization purposes, we also want this function to be continuous. Better still if we know the expression of its derivate. Both the sigmoid function defined as σ and the arctan function have the desired properties, though the simple ReLU (Figure 3) is shown to be more efficient for optimization (source) and thus more common in practice.

1.4 Softmax

Before evaluating a loss function, we transform the output layer in a way that their values are positive and their sum is equal to 1. In this way, each value in the output layer can be interpreted as the probability that a given input is classified as such. While in the training phase, we obviously know the desired output to each input, thus we want that specific “neuron” to have a value close to 1, and all the others to be close to 0. The loss function, which quantifies “how close” the network’s guess is to a given output, can have different definitions depending on the problem. A common loss function is the so called ilogit defined as ilogit . In order to quantify how well the network performs on the entire training set, we want to evaluate the loss function on each data-input. The cost function can then be defined as cost . Or the average of losses.

1.5 Forward propagation

We now have an expression for each component in the neural network. Let’s see how they are assembled in the forward propagation (Figure 4).

1.6 Backward propagation

Backward propagation, finding the gradient. As previously mentioned, the NN can be seen as a composition of functions of several functions of which we can easily find the derivate. Intuitively, in order to find the gradient of the cost function with respect to the weights and biases, we can just apply the chain rule (Figure 5). As we can see, the output layer is dependent on its preceding layer, weights and biases, which, in turn, depends on the previous parameters. This creates a cascade of gradients calculated from the output- to the input layer, hence the name backward propagation.

After iterating over the training set, we now have the gradients needed for the optimization algorithm, which in turn allows us to adjust the weights and biases in order to, iteration by iteration, reach a local minimum. In practice, this approach performs relatively well, though there are improvements to be made. In the following chapters, we will examine how convolutional neural networks differs from the classic MLP.

2 Convolutional Neural Networks

A convolutional neural network is an evolution of a classical multilayer perceptron network. Recall the basic principle underpinning how a normal neuron in a neural network is supposed to work. The neuron is supposed to ‘look’ for features in its input data. If the neuron ‘thinks’ that those features are present in the input data it ‘fires’. Otherwise the neuron does not fire.

In a classical multilayer perceptron network this is implemented in the following way. Each neuron contains a vector of weights, a bias and an activation function. The input to the neuron—which must be a vector of equal length to the neuron’s own weight vector—is combined with the neuron’s weight vector using the dot product. The neuron’s bias is added onto the result which in turn is passed to the activation function which finally determines if the neuron ‘fires’ or not.

Using several layers of neurons one can achieve quite remarkable results using this implementation of a neural network. However, a multilayer perceptron is inherently limited. The major problem is that neurons in these kinds of networks only accept input that is in the form of a vector. This means that for applications where it is not natural for the input to be in a vector format, say image recognition, the input first has to be translated to a vector format. Usually this results in a loss of information contained in the input. In the typical case of image recognition, the input is in the form of one or more arrays of two dimensions. For a multilayer perceptron to treat this input, the images has to be ‘flattened’ into a vector of one dimension before it can be passed on to the network. This procedure eliminates some of the pixel relations in the image. To deduce this, consider the process of reconstructing a flattened image. If the image’s dimensions prior to being flattened is not known, it is impossible, without the aid of pattern recognition, to reconstruct the image and be sure the reconstruction is equal to the original image.

To fix this problem, we can use a ‘simple’ solution. Instead of having the neuron contain a vector of weights, let it have an array of weights. If we change the neuron’s vector of weights into an array, we also need to change the operation that is used to combine the weights with the input (which in a neuron of a multilayer perceptron is the dot product). There are two things to consider here. The purpose of the weights is to look for features or *patterns* in the input—by emphasizing or deemphasizing certain aspects—and the operation must reflect this purpose of the weights. Furthermore, the result of the operation should be a single number which, in a sense, represents the neurons ‘initial’ confidence that the feature it is looking for is present in the input. The operation which does both of these things is the *Hadamard product*. The Hadamard product can be viewed as an extension of the dot product to two dimensional arrays. It combines two arrays—of the same dimensions—by multiplying corresponding entries together and summing the results. Which is precisely what the dot product does with two vectors.

The Hadamard product (Hp) of two matrices A, B with entries $a_{i,j}, b_{i,j}$ of equal dimensions is

$$\text{Hp}(A, B) = \sum_j \sum_i a_{i,j} \cdot b_{i,j} \quad (1)$$

Another thing we need to take into account is the dimensions of our new weight array. If we were to proceed analogically to how a multilayer perceptron works, the array should have the same size as the input to the neuron. ‘Connecting’ each value in the input to an individual weight in the neuron. But this approach means that each neuron, in principle, looks for a single feature in the

entire input at once. A more refined approach, is to let the neuron's weight array be smaller in size than its input. Instead of applying the weight array (using the Hadamard product) to the entire input at once, we apply it (still using the Hadamard product) to portions of the input separately. Intuitively, this means that the weight arrays are 'scanned' across the entire input image. This allows the neuron to be trained to look for a single feature, such as a 'sharp edge' or a 'round corner', in multiple areas of the input. The result of this method will no longer be a single number representing how 'initially' confident the neuron is that a particular feature is present in the input **as a whole**. Rather, the result becomes a *feature map*. Another two dimensional array which represents how 'initially' confident the neuron is that a particular feature is present in **specific locations** of the input.

Going one step further, we can allow the neuron to treat inputs of not only a single two dimensional array, but several two dimensional arrays. A typical example where the input would consist of several interlinked two dimensional arrays is RGB images. An RGB image consists of three arrays of pixel values (numbers) that describe how red, green and blue an image is in each pixel. The number of interlinked two dimensional arrays present in the input, is known as the number of *channels* that the input has. In order for our neuron to treat inputs with more than one channel we let the neuron have as many channels as the input. That is to say, we equip the neuron with a weight array for each channel in the input. For each channel the weight arrays are applied to the input (using Hadamard) and the result in each channel is combined to form a final single feature map.

A neural network which makes use of layers of neurons of this kind, is a convolutional neural network. The multi-channel two dimensional arrays of weights inside each such neuron is called the neuron's *kernel* or *filter*. Why are these neural networks called convolutional neural networks? That will be explained in the next section.

2.1 The convolutional neuron

Let us start this section with a simple convolutional neuron. The neuron's kernel consists of a single weight matrix (one channel), some bias and some activation function. Consequently, the input that this neuron accepts is any one channel two dimensional array of size greater than its kernel.

The process of calculating the neurons 'initial' confidence that the feature it is looking for is present in the input, is illustrated by Figure 1 on page 6. Figure 1 shows how the kernel 'scans' the entire input array to compute the feature map. Let us construct the general formula for this feature map.

We will denote the feature map, kernel and input as F , K and M respectively. They are two dimensional arrays and their individual entries will be denoted as $F(i, j)$ where i is the row index and j the column index. The indexes will all start at zero (e.g. $M(1, 2)$ is the entry in the second row, third column of M . $M(0, 0)$ is the *first* entry). Horizontal and vertical lengths of arrays will be denoted by h and v with the associated array as a subscript (e.g. h_K is the horizontal length of the kernel). A typical MNIST image of dimensions 28x28 as the input, will have $h_M = v_M = 28$, its indexes will range in $[0..v_M - 1] \times [0..h_M - 1]$.

Let us say that the kernel scans its input one column or one row at a time. In which case, the output feature map will have dimensions

$$\begin{aligned} h_F &= h_M - h_K + 1 \\ v_F &= v_M - v_K + 1 \end{aligned} \tag{2}$$

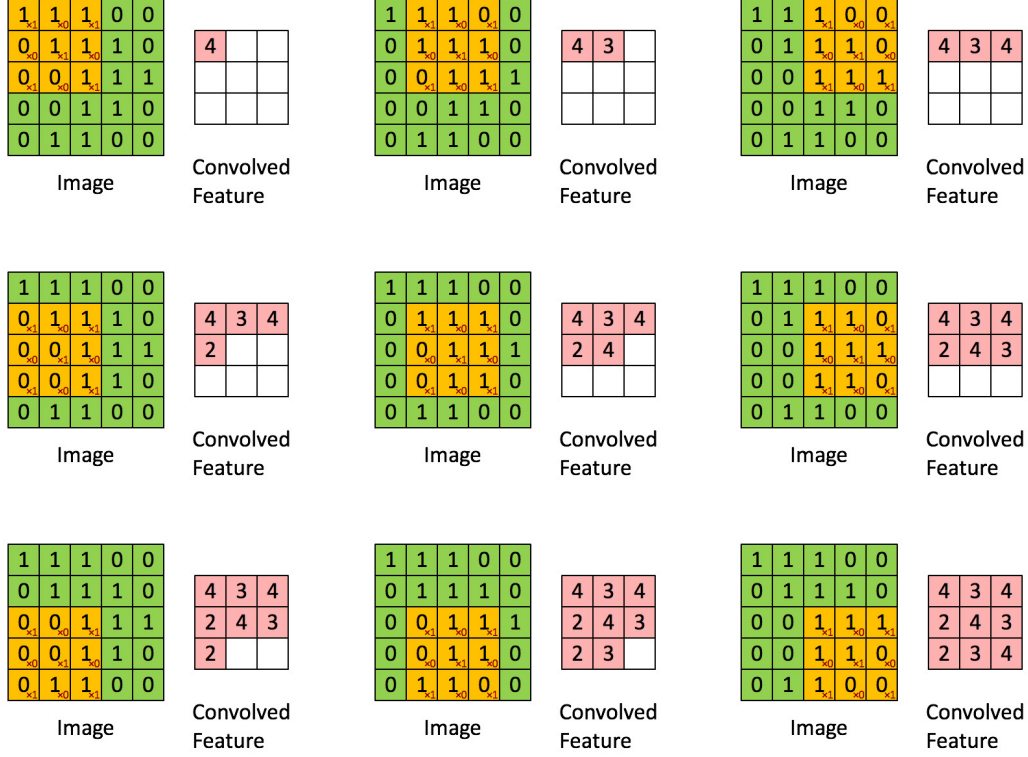


Figure 1 The basic forward operation of a convolutional layer

The entries in the initial F^* (before activation) are produced by applying the Hadamard product (Eq. (1)) to each area of dimensions $h_K \times v_K$ in the input. Thus the entry at row $i \in [0..v_F - 1]$ and column $j \in [0..h_F - 1]$ is computed as

$$F^*(i, j) = \sum_{i=0}^{v_K-1} \sum_{j=0}^{h_K-1} K(i, j) \cdot M(i + i, j + j) \quad (3)$$

If we add a bias b to this filter, the formula becomes

$$F^*(i, j) = \sum_{i=0}^{v_K-1} \sum_{j=0}^{h_M-1} \left(K(i, j) \cdot M(i + i, j + j) + b \right) \quad (4)$$

Lastly, the final output feature map (F) of the filter is produced by sending all of these entries through an activation function (ReLU or sigmoid in practice).

$$F(i, j) = \text{activation} \left(\sum_{i=0}^{v_K-1} \sum_{j=0}^{h_M-1} \left(K(i, j) \cdot M(i + i, j + j) + b \right) \right) \quad (5)$$

2.1.1 Where is the convolution?

In mathematics, the convolution operation is denoted by $*$ and the convolution s (itself a function) of two functions k and m is defined as

$$s(t) = (k * m)(t) = \int_{-\infty}^{+\infty} k(x) \cdot m(t - x) dx \quad (6)$$

$$s(\textcolor{red}{i}) = (k * m)(\textcolor{red}{i}) = \sum_{i=-\infty}^{+\infty} k(\textcolor{blue}{i}) \cdot m(\textcolor{red}{i} - \textcolor{blue}{i}) \quad (7)$$

For continuous (Eq. (6)) and discrete (Eq. (7)) functions respectively. On a computer, we are in practice always working with discrete convolutions so for us the convolution of interest is Eq. (7). If we make the functions k and m two dimensional—they take two arguments as their input—their convolution also becomes two dimensional.

$$s(\textcolor{red}{i}, \textcolor{red}{j}) = (k * m)(\textcolor{red}{i}, \textcolor{red}{j}) = \sum_{i=-\infty}^{+\infty} \sum_{j=-\infty}^{+\infty} k(\textcolor{blue}{i}, \textcolor{blue}{j}) \cdot m(\textcolor{red}{i} - \textcolor{blue}{i}, \textcolor{red}{j} - \textcolor{blue}{j}) \quad (8)$$

Let us known say that k and m are two functions which index each an array of two dimentions. Meaning that the functions take two indexes as their input arguments. Let us further assume that k and m 's arrays are zero at any index but for those contained in a small area. For k that area is $[0 .. v_k - 1] \times [0 .. h_k - 1]$ and for m it is $[0 .. v_m - 1] \times [0 .. h_m - 1]$. Where h (horisontal length) and v (vertical length)—of both subscripts—are natural non-zero numbers. In this case, the convolution of k and m reduces to

$$s(\textcolor{red}{i}, \textcolor{red}{j}) = (k * m)(\textcolor{red}{i}, \textcolor{red}{j}) = \sum_{i=0}^{v_k-1} \sum_{j=0}^{h_k-1} k(\textcolor{blue}{i}, \textcolor{blue}{j}) \cdot m(\textcolor{red}{i} - \textcolor{blue}{i}, \textcolor{red}{j} - \textcolor{blue}{j}) \quad (9)$$

Where $\textcolor{red}{i} \in [0 .. v_m - 1]$ and $\textcolor{red}{j} \in [0 .. h_m - 1]$.

The resemblance to Eq. (3) is already apparent but there is a subtle difference. Eq. (9) subtracts the blue indices from the red ones while Eq. (3) adds them together. The effect of subtracting instead of adding renders the convolution operation $*$ commutative.

$$(k * m)(\textcolor{red}{i}, \textcolor{red}{j}) = (m * k)(\textcolor{red}{i}, \textcolor{red}{j})$$

$$\sum_{i=0}^{v_k-1} \sum_{j=0}^{h_k-1} k(\textcolor{blue}{i}, \textcolor{blue}{j}) \cdot m(\textcolor{red}{i} - \textcolor{blue}{i}, \textcolor{red}{j} - \textcolor{blue}{j}) = \sum_{i=0}^{v_m-1} \sum_{j=0}^{h_m-1} m(\textcolor{red}{i}, \textcolor{red}{j}) \cdot k(\textcolor{red}{i} - \textcolor{blue}{i}, \textcolor{red}{j} - \textcolor{blue}{j}) \quad (10)$$

Bringing the discussion back to neural networks, kernels and inputs—what is the effect of replacing the pluses in Eq. (3) with minuses? Well at first glance it seems that you cannot do it as you would end up with indexes that are out of bounds for the input image M . But let us assume that M behaves as m meaning that it is zero everywhere where the indices are out of bounds. Suppose we have a kernel that has dimentions $h_K = v_K = 3$ and an input with dmientions $h_M = v_M = 5$, as in Figure 1. Suppose we are computing F at position $i = 5, j = 5$ i.e. $F(5, 5)$. Notice how the top-left coordinate of K pairs up with the bottom-right coordiante of the covered region of M

$$K(0, 0) \cdot M(5 - 0, 5 - 0) = K(0, 0) \cdot M(5, 5)$$

And the top-right coordiante of K pairs up with the bottom-left coordiante of the covered region of M

$$K(0, 2) \cdot M(5 - 0, 5 - 2) = K(0, 0) \cdot M(5, 3)$$

Likewise, the bottom-left pairs with the top-right and the bottom-right pairs with the top-left.

$$K(2, 0) \cdot M(5 - 2, 5 - 0) = K(2, 0) \cdot M(3, 5)$$

$$K(2, 2) \cdot M(5 - 2, 5 - 2) = K(0, 0) \cdot M(3, 3)$$

It appears that what you end with—is a kernel *flipped* about its horizontal and vertical axes. This means that flipping the kernel relative to the input, renders the computation of the feature map commutative. But this is a property that is useful in mathematical contexts and less so for neural networks. It is Eq. (3) that is commonly used in machine learning applications and the associated mathematical operations is actually what is called the *cross-correlation*. Despite this, neural networks that implement these kinds of neurons are called convolutional even if the underlying mathematical operations is more often than not a close relative of the proper convolution operation.

2.1.2 Multiple channels

Let us return to Eq. (5) for a minute and ask: what happens if we add some channels to the input? Let us denote the number of channels by C . The input and kernel are now multi-channel arrays, their individual entries will be denoted as $K(i, j, c)$ with c ranging in $[0..C-1]$. To calculate the feature map when multiple channels are present, we have to take into account the contribution of each individual channel. We do this by simply adding them together. Since we only want to add the bias of the filter once per entry in the feature map, the resulting modified formula becomes

$$F(i, j) = \text{activation} \left(\sum_{i=0}^{v_K-1} \sum_{j=0}^{h_K-1} \left(\sum_{c=0}^{C-1} (K(i, j, c) \cdot M(i + i, j + j, c)) \right) + b \right) \quad (11)$$

2.1.3 The stride

So far we have considered a feature map produced by scanning the kernel over the input one *step* at a time. As is demonstrated in Figure 1 on page 6. However, we could let the kernel take longer steps instead of moving just one row or one column at a time. The ‘length’ of each step that the kernel takes while moving over the input is what is known as the *stride*. Which we can further divide into the horizontal and vertical stride.

What purpose does altering the stride serve? Well the first thing to notice is that increasing the stride reduces the size of the resulting feature map. Meaning that the computational intensity of the neuron is reduced. Which in actuality is a side effect of the fact that the neuron, with a larger stride, studies less information in its input. Which entails that it provides less information in its output.

Processing less information can be both a good and a bad thing. It depends on the characteristics of the input data that the neuron is looking at. If it is judged that the size of the input data’s features are larger than the neuron’s kernel, then increasing the stride is justified. If this is not the case, then increasing the stride will almost certainly leave the network unable to correctly fulfill its intended task.

Let us now further develop Eq. (11) to take into account variable strides and starting positions. We will denote the starting position in the vertical and horizontal positions by i_s and j_s respectively. i_s and j_s are of course numbers in $[0..v_K-1]$ and $[0..h_K-1]$ respectively.

$$F(x, y) = \text{activation} \left(\sum_{i=0}^{v_K-1} \sum_{j=0}^{h_K-1} \left(\sum_{c=0}^{C-1} (K(i, j, c) \cdot M(i + i, j + j, c)) \right) + b \right) \quad (12)$$

2.1.4 Zero padding

2.2 The pooling ‘neuron’

Downsampling and purpose of downsampling. Noise reduction and computational reduction.

2.3 Backward Propagation

Obtain formula for derivative of convolution and pooling.

3 Python Implementation

- The problem the network is going to tackle (classic MNIST image recognition)
- The architecture of the network
- Choice of various functions and their associated code
- Results with time performance
- Possible Modifications

4 Conclusion

The advantages of a CNN over a NN.