

# ARTIFICIAL NEURAL NETWORKS

---

## PROJECT REPORT

---

*Authors:*

Michiel BONGAERTS

Marjolein NANNINGA

Tung PHAN

Maniek SANTOKHI

May 31, 2015

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Concept</b>	<b>3</b>
2.1	Impression . . . . .	3
2.2	MoSCoW . . . . .	4
2.2.1	Must have . . . . .	4
2.2.2	Should have . . . . .	4
2.2.3	Could have . . . . .	4
2.2.4	Would have . . . . .	4
<b>3</b>	<b>Implementation</b>	<b>5</b>
3.1	Frontend and Backend . . . . .	5
3.2	Dataset . . . . .	5
3.3	Convolutional Neural Network . . . . .	6
3.3.1	LeNet-1 . . . . .	6
3.3.2	Back-propagation . . . . .	6
<b>4</b>	<b>Results and Discussion</b>	<b>8</b>
<b>5</b>	<b>Time schedule</b>	<b>10</b>

# 1. Introduction

Mapping the world around us has always been a human endeavour to advance economical output. A better understanding of the places around us makes for more efficient travelling and exploitation of the land. However, it has always been a very slow and tedious process to produce these maps, something technology has not changed just yet.

A new opportunity has arisen with the arrival of satellite imagery and an ever increasing amount of computational power. An opportunity where this mapping can be done automatically so that this tedious and slow job can be processed even more quickly and perhaps more accurately. It is with this in mind we further analyse any possibilities.

This paper proposes an update. Now a more straightforward approach has been chosen in which the emphasis lies on the actual Neural Network rather than the conversion of interpreted images to vector graphic maps. The new approach deals with image patches rather than pixels. This document discusses the newly acquired concept with a list of features and the actual implementation details. Also an updated schedule will be presented.

## 2. Concept

Earlier attempts to conceptualise the idea to automate map making resulted in a proposal that too heavily focussed on the actual map creation rather than the classification. For a Neural Network course this was deemed not befitting enough. The plan was also quite far reaching to start with. Thoughts were put into downscaling this ambitious plan. We played around a bit and came up with a new concept which will be discussed in this chapter. Firstly, an impression is given how the end user interacts with our system. This will lay the groundwork for how the Neural Network will be constructed. Secondly through the principles of MoSCoW a flexible requirements list is established. Actual talk about the classification is done in the subsequent chapter.

### 2.1 Impression

Below in figure 2.1 an impression is given of what the end user will interact with and a possible result that might come about from said interaction.

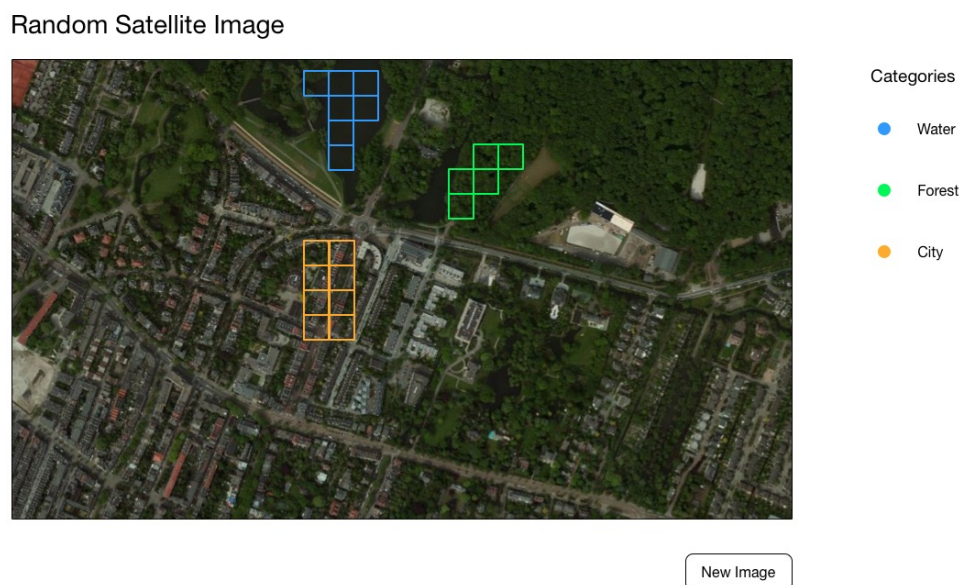


Figure 2.1: Impression of what the user interacts with and a possible outcome.

The most notable attention grabber is the satellite image. This image will be acquired through a random call from a homogenous satellite image database on every new instantiation of the system. Another possibility to acquire new data is by clicking on the 'New image' button. Right of the satellite image one can see labels (Forest, City, Water) which correspond with the labels which are outputs of the classification. Results obtained from our algorithm, given the current satellite image as input, are graphically feed back to the user. The impression above does that by showing a correspondence between image patches and their respective label via color coding laid over the satellite image. This rendition just shows a few islands of results as an example. Normally the entire image will have such arching (which will be a lot more subtle).

Figure 2.2 shows how it works internally. Two grids are maintained. One with patches the size of 25 by 25 pixels. The other by the size of 50 by 50 pixels. The latter is actually fed to the Neural Network from which will be decided for that patch the percentage of type of label it contains. Four 50 by 50 pixels will be grouped together to create a square. In the middle the 25 by 25 pixels patch is placed. This patch will eventually be coloured on the satellite image to indicate the type of label. To decide that, a weighted majority vote of the four 50 by 50 pixels surrounding that smaller patch is computed. Like a convolution filter this construction is shifted over the satellite image as to decide for every area on it. As one can imagine, 25

pixels at the borders are omitted.

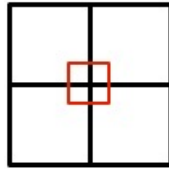


Figure 2.2: Grid structure which facilitates the algorithm.

## 2.2 MoSCoW

Since a limited amount of time is available and we thought of quite some experiments and features that could be added, we used the MoSCoW method to get our priorities straight and focus on the most important requirements. MoSCoW stands for *Must have*, *Should have*, *Could have* and *Would have*. All the requirements are labeled in these four classes.

### 2.2.1 Must have

Must have requirements are critical to project success.

- Create code based on a *Convolutional Neural Network* (CNN) that enables automatic classification of patches from satellite images. Considering images acquired on provincial level and a substantial amount of pixels in one patch (at least 50 x 50 pixels).
- At least the following classes should be recognized: vegetation, city and water.
- Develop a way to visualize the automatic classifications clearly.

### 2.2.2 Should have

Requirements labeled as should have are important to book success, but not necessary for delivery.

- Create a clear interface in which the unlabeled images can be uploaded, and the output consists of labeled images.
- Calculate the uncertainty in the classification and ask user input for very uncertain patches.

### 2.2.3 Could have

It would be very nice if we would be able to reach the Could have features, but they are not critical.

- Experiment with pre-processed images (noise reduction, gradient calculations)
- Analyze images on city level, so with more details present. For this purpose new classes have to be added, such as roadways, cycle paths, buildings, distinct vegetations etc.
- Experiment with other models than the state-of-the art LeNet-1 CNN. For examples, a CNN in which Genetic Algorithms are incorporated, or implementing an Extreme Learning Machine for the training of the weights.

### 2.2.4 Would have

These requirements are implemented only in the most ideal situation. They are considered as the dream project, sometimes serving as a suggestion for further projects.

- Develop a method for high-detailed automatic vector graphics, in which segmentation of the distinct labeled classes is incorporated.
- Use the input of the users to improve the automatic classification.
- Sell the software package to Google.

## 3. Implementation

A plan has been established what kind of application should come about. The previous discussion pressed for a certain kind of structure. One where a Neural Network is central to the problem to be solved. But also a frontend is needed for certain user interaction as well as an infrastructure in the backend which facilitates the communication with said frontend. This chapter discusses these aspects.

### 3.1 Frontend and Backend

The impression given above in figure 2.1 is close to what the end result should look like (although it only displays a certain state). Yet an entire infrastructure outside of the Neural Network is needed to facilitate the application. Below in figure 3.1 an infrastructure is visualised via the communication of the two entities.

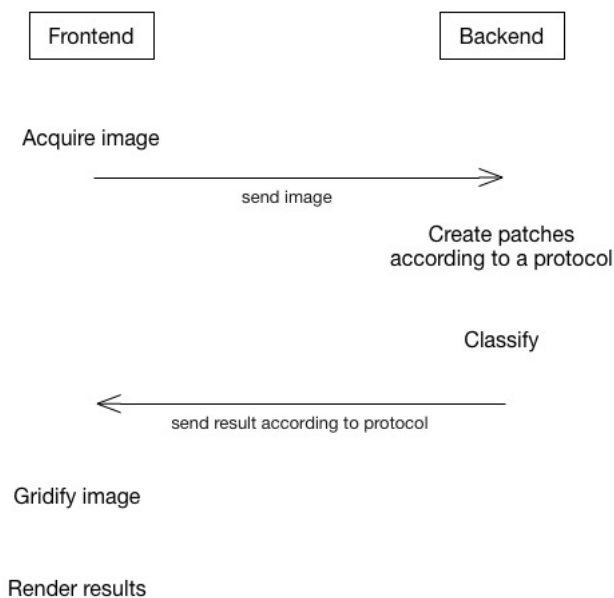


Figure 3.1: Communication visualization between frontend and backend.

First, at the frontend, an image is acquired by the user from a database call. The image is immediately send to the backend. There, patches are extracted according to a certain protocol (specific size and sequence). The patches are classified by our Neural Network. Send back to the frontend are the results in a list according to the sequence defined by the protocol. Also a codification of the results is part of the protocol. Back at the frontend the image is made into a grid which corresponds to the size of the protocol. According to the results of the Neural Network these grid windows are coloured.

### 3.2 Dataset

The image which is being acquired every time comes from the Bing Maps database. A predefined size of the map with satellite view enabled, the labels turned off, no other controls or logos visible and where the hight is 100 meters is put in the frontend. Randomly a location inside the Netherlands is generated. Now the actual image is taken from the map.

For training, satellite images from the same source and practice mentioned above are drawn on a per category/label basis. Patches are acquired code-wise over those images so that for each patch it is known what category/label belongs to it.

### 3.3 Convolutional Neural Network

During our research phase, we encountered a variety of papers that describes the use of image classification based on Neural Networks. Most of these papers have one thing in common which is the use of Convolutional Neural Networks (CNN). Since convolution operations are widely used to extract features from images and since these operations can be represented in terms of a Neural Network these properties lead to the existence of the Convolutional Neural Networks.

In general a CNN architecture consist of multiple Convolution layers and sub-sampling layers. It depends on the architecture how these layers are followed by each other. The convolution layer is the layer which results after the convolution operator is performed by the a convolution kernel. Since the gaol is to extract general features from our input image we want the CNN to generalize. This generalization is partly released by dimension reduction. Since convolution operation reduce the dimension of our input map with  $\frac{M-N+1}{M}$ , where  $N$  is the dimension of the convolution kernel and  $M$  the dimension of the input map, we want the dimension to be reduced more quickly. This is done by dub-sampling kernels which in general 'squeeze' or average the input map with a certain dimension. This operation is equivalent to a convolution operation but with a larger step-size (convolution uses a stepsize equal to 1) and equal weights for each element in the sub-sampling kernel.

#### 3.3.1 LeNet-1

The model we implement is LeCun's LeNet-1. This model uses an alternating sequence of convolution and sub-sampling layers. The architecture of this network is shown in Figure 3.2. The convolution layers act as feature maps, they consist of a window of a certain size, whose pixels are trainable weights. The sub-sampling layers reduce the dimensionality of the outputs.

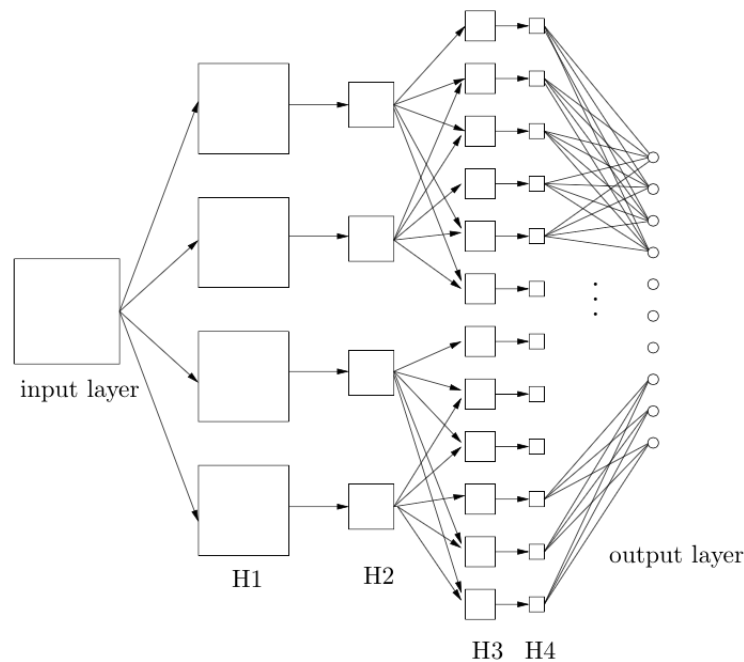


Figure 3.2: The architecture of the CNN proposed by LeNet. H1 and H3 are the convolution layers, H2 and H4 the sub-sampling layers to reduce the dimension. **Source:**

One of the biggest advantages of using CNN, and especially LeCun's LeNet-1 implementation, is the incorporation of backpropagation learning. Meaning that all the weights of the layers are adjusted iteratively, eliminating the need to manually create the convolution masks.

#### 3.3.2 Back-propagation

To train the network we have to adjust several parameters in the network. In CNN

$$\begin{aligned}
f(x) &= \frac{1}{1 + e^{-x}} \\
E &= \sum_k \frac{1}{2} |t_k - F_k|^2 \\
F_k &= f(x_k) \\
x_k &= \sum_{ij} W^{4rbk}[i, j] H^{4rb}[i, j] - b_k \\
H^{4rb}[i, j] &= \sum_{u, v} W^{3rb}[u, v] H^{3rb}[2i + u, 2j + v] \\
H^{3rb}[2i + u, 2j + v] &= f(x^{3rb}[2i + u, 2j + v]) \\
x^{3rb}[2i + u, 2j + v] &= \sum_{nm} W^{2rb}[n, m] H^{2rb}[n + (2i + u), m + (2j + v)] - b^{3rb}[2i + u, 2j + v]
\end{aligned} \tag{3.1}$$

In these equations  $r$  stands for row which is the first branch after the first convolution i.e. the total rows is equal to the amount of different feature maps in H1.  $b$  stands for the amount of branches per row.

$$\Delta W^{4rbk}[i, j] = \sum_k -\eta \frac{dE}{dF_k} \frac{dF_k}{dx_k} \frac{dx_k}{dW^{4rbk}[i, j]} \tag{3.2}$$

$$\begin{aligned}
\Delta W^{2rb}[n, m] &= \\
\sum_k -\eta \frac{dE}{dF_k} \frac{dF_k}{dx_k} \sum_{ij} \frac{dx_k}{dH^{4rb}[i, j]} \sum_{uv} \frac{dH^{4rb}[i, j]}{dH^{3rb}[2i + u, 2j + v]} \frac{dH^{3rb}[2i + u, 2j + v]}{dx^{3rb}[2i + u, 2j + v]} \sum_{nm} \frac{dx^{3rb}[2i + u, 2j + v]}{dW^{2rb}[n, m]} \\
\Delta b^{3rb}[2i + u, 2j + v] &= \\
\sum_k -\eta \frac{dE}{dF_k} \frac{dF_k}{dx_k} \sum_{ij} \frac{dx_k}{dy^{4rb}[i, j]} \sum_{uv} \frac{dy^{4rb}[i, j]}{dy^{3rb}[2i + u, 2j + v]} \frac{dy^{3rb}[2i + u, 2j + v]}{dx^{3rb}[2i + u, 2j + v]} \frac{dx^{3rb}[2i + u, 2j + v]}{db^{3rb}[2i + u, 2j + v]}
\end{aligned} \tag{3.3}$$



## 4. Results and Discussion

In the first phase of our research we made a convolutional neural network based on LeNet-1. We experimented with different convolution kernels in the first layer (H1) of the network and differed the amount of branches in the third layer (H3). We started with manually chosen convolution kernels for both layers H1 and H3. We began with input patches of  $50 \times 50$  followed by a convolution kernel of  $5 \times 5$  resulting in a feature map of dimensions  $46 \times 46$ . Next, sub-sampling with a dimension reduction of 2 results in a feature map of dimension  $23 \times 23$ . Again we did different  $4 \times 4$  convolution kernel operation creating branches in H3 resulting in  $20 \times 20$  feature maps. Sub-sampling with a dimension reduction of 5 resulted in an  $4 \times 4$  output map. These outputmaps are then connected to two different classification neurons for Forest and City.

We used the Sigmoid function as activation function only for our classification-neurons. Thus, no activation functions were used earlier in the network (compared to LeNet-1). We performed back-propagation on the weights connecting the output maps with the classification neurons and biases. The first functioning results were found with 4 feature maps in H1 and 2 feature maps in H3 each. We used the convolution kernels shown in figure 4.1. The Network was trained on structure only so no color features were included.

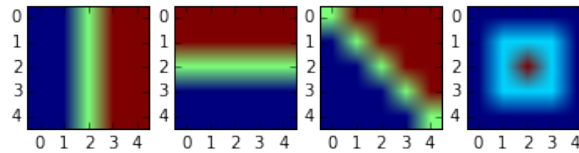


Figure 4.1: The four filters used in the first stage where the CNN was able to distinguish forest and city.

In the second phase we tried to extend our classes with the class *water*. We tried to investigate whether training only the last weights and biases were sufficient to train the network to classify three classes. During this phase we made some adjustments to the network. We realised that the last sub-sampling with a dimension reduction of a factor 5 would probably lead to less accuracy by the network since more information might be lost in this reduction. Although, we did no quantitative research on this hypothesis the underlying reasoning is that in each layer of the network more generalizations are made by the network which should represent the information of your input patch. When you downsize the feature maps at the end of the network this more contained information will be lost in this sub-sampling compared with sub-sampling earlier in the network. Therefore we introduced a slightly different CNN. Now, the inputsize of the patches is  $41 \times 41$ . This is convolved with a  $9 \times 9$  convolution kernel resulting in a  $33 \times 33$  feature map. Now sub-sampling with a factor 3 is introduced reducing it to a  $11 \times 11$  feature map. From this map we can branch off (H3) with a convolution operation with a  $4 \times 4$  convolution kernel resulting in a  $8 \times 8$  feature map. This last map is then sub-sampled with a reduction of 2 resulting in the output map of  $4 \times 4$ . In this way we presume that the network makes better generalization than the earlier network.

Also, in this new architecture the first convolution kernels are increased in dimension. This adjustment has been introduced due to a problem with respect to the classes *forest* and *water*. It turns out that these two classes are very similar in structure (see figure 4.2). Conceivably, for the network it becomes hard to distinguish these classes because of their similarities. Therefore we increased the size of the first convolution kernels. In this way more context per pixel is included leading to larger differences in feature maps between these classes.

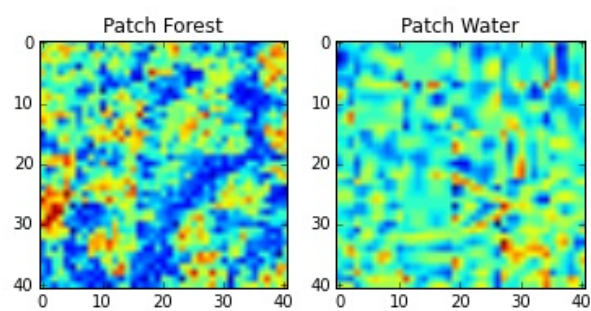


Figure 4.2: A forest and water patch are shown. Their similarity makes it hard for the CNN to distinguish these two classes.

## 5. Time schedule

Week	Description of work	Deadline	Labour (hours)
20	<ul style="list-style-type: none"> <li>• Implement backpropagation output layer</li> <li>• Start examining how to take into account RGB values</li> <li>• Train on more pictures and examine performance on other classes like city and water</li> <li>• Start developing frontend/backend architecture</li> </ul>	Finish updated proposal	40
21	<ul style="list-style-type: none"> <li>• Work on implementing communication protocol</li> <li>• Experiments with connecting the nodes of the different layers of the Neural Networks</li> </ul>	Finish backpropagation output layer, and frontend implementation of the communication protocol	40
22	<ul style="list-style-type: none"> <li>• In-between overall evaluation</li> <li>• Backend and frontend implementation</li> <li>• Experiment with preprocessed images (noise reduction)</li> <li>• Implement backpropagation further, until layer H3, see Figure 3.2</li> </ul>	Finish model with complete RGB values and the frontend/backend architecture must be completely finished	40
23	Start report writing and implement last adjustments	Finished product for demonstration	40
24	Preparation for demo	Demo	30
25	Project finalization	-	30
26	Project finalization	Project report + individual document	40
			Total: 468