



Analyzing the Information Content of Multiple Views of an Object in Object Detection with Neural Networks

Dennis Kraus

Electrical Engineering and Information Technology

10.06.2019

Supervisor:
Sebastian Schrom, M.Sc.

REGELUNGSMETHODEN
UND ROBOTIK 

Prof. Dr.-Ing. J. Adamy

Eidesstattliche Erklärung

Erklärung zur Abschlussarbeit gemäß § 22 Abs. 7 und § 23 Abs. 7 APB TU Darmstadt

Hiermit versichere ich, Dennis Kraus, die vorliegende Master-Thesis gemäß § 22 Abs. 7 APB der TU Darmstadt ohne Hilfe Dritter und nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die Quellen entnommen wurden, sind als solche kenntlich gemacht worden. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen. Mir ist bekannt, dass im Falle eines Plagiats (§38 Abs. 2 APB) ein Täuschungsversuch vorliegt, der dazu führt, dass die Arbeit mit 5,0 bewertet und damit ein Prüfungsversuch verbraucht wird. Abschlussarbeiten dürfen nur einmal wiederholt werden.

Bei der abgegebenen Thesis stimmen die schriftliche und die zur Archivierung eingereichte elektronische Fassung gemäß § 23 Abs. 7 APB überein.

Darmstadt, 10.06.2019

Dennis Kraus

(English translation of above declaration for information purposes only)

Thesis Statement

pursuant to § 22 paragraph 7 and § 23 paragraph 7 of APB TU Darmstadt

I herewith formally declare that I, Dennis Kraus, have written the submitted Master's Thesis independently pursuant to § 22 paragraph 7 of APB TU Darmstadt. I did not use any outside support except for the quoted literature and other sources mentioned in the paper. I clearly marked and separately listed all of the literature and all of the other sources which I employed when producing this academic work, either literally or in content. This thesis has not been handed in or published before in the same or similar form. I am aware, that in case of an attempt at deception based on plagiarism (§38 paragraph 2 APB), the thesis would be graded with 5,0 and counted as one failed examination attempt. The thesis may only be repeated once.

In the submitted thesis the written copies and the electronic version for archiving are pursuant to § 23 paragraph 7 of APB identical in content.

Abstract

english abstract

Zusammenfassung

deutscher abstract

Contents

1	Introduction	1
1.1	Overview	1
1.2	Motivation	1
2	Related Work	3
3	Fundamentals	5
3.1	Artificial Neural Networks	5
3.1.1	Overview	5
3.1.2	Multilayer Perceptron	6
3.1.3	Convolutional Neural Networks	11
3.2	Software	11
4	Methods	13
5	Results	15
6	Discussion	17
6.1	Conclusions	17
6.2	Outlook	17

List of Figures

3.1	Model of a Neuron	6
3.2	Model of a Perceptron	9
3.3	Multilayer Perceptron	10
3.4	Handwritten Digit from the MNIST Digit Dataset	11

List of Tables

3.1 Truth Tables of Logical Operations	7
--	---

Chapter 1

Introduction

1.1 Overview

1.2 Motivation

Chapter 2

Related Work

Chapter 3

Fundamentals

This chapter covers the fundamentals necessary to understand the methods presented and their application. It is divided into a section on neural networks and one on the used software and frameworks. The former starts with the principle of a neural network. It continues with an explanation of the required steps to train the network and ends with an overview of Convolutional Neural Networks, which are networks that are more suited for image processing. The latter explains which software and frameworks support building and training a model.

3.1 Artificial Neural Networks

This section examines the types of neural networks that are important for this work. Furthermore, it explains how these types are build in order to manage the wanted use case.

3.1.1 Overview

Artificial neural networks are vaguely inspired by the biological neural networks that constitute animal brains for recognizing patterns. Its task is being an universal approximator for any unknown function $f(x) = y$ where x is the input and y the output. There are two conditions that need to be fulfilled. One is the relation of x and y and the other is the presence of numerical data. So every data like images, text or time series must be translated. The complexity of the approximated function depends on the use case but usually it is highly non-linear. General use cases for neural networks embrace classification, clustering and regression.

Classification means the network divides given data like images into categories by recognizing patterns. This is the task used in this work. The correct category of each input is given as an additional label. Therefore, the network learns the correlation between data and labels. Kind of an downside here is that every input must be labeled by human knowledge beforehand. This kind of learning is called supervised learning, because each predicted category by the network needs to be compared with the ground truth label. Use cases are for example the classification of cars in images or even the

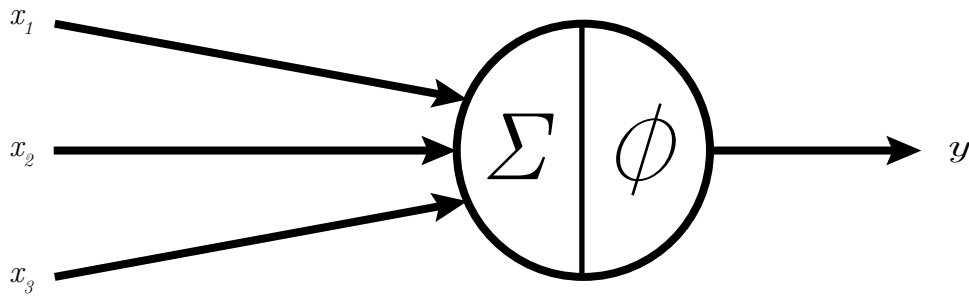


Figure 3.1: Model of a neuron. The inputs x_i are summed up and put into the activation function ϕ whose result is the neuron's output y .

type of car in an image or whether an email is spam. Again, it all depends on the wanted use case and given data.

Clustering divides data into clusters or groups, respectively, but without requiring labels. Therefore, this learning type is called unsupervised learning. So it is kind of an classification task with a dynamic category creation. Use cases are comparing data to each other and finding similarities or anomalies. Because unlabeled data occurs way more often than labeled data in real world examples, a network can train on a broader range of related data and probably gets more accurate than a classification one.

Regression is the prediction of a future event by establishing correlations between past events and future events. A simple use case is the prediction of a price of a house given its size and the size-price data pairs of different houses. A more advanced use case is the prediction of hardware breakdowns by establishing correlations of already known data.

3.1.2 Multilayer Perceptron

This section starts with an explanation of a single computational neuron and its development to become a perceptron and ends with an overview of the multilayer perceptron architecture.

3.1.2.1 Perceptron

McCulloch and Pitts[1] were the first who defined a computational model of a neuron that corresponds to the functionality of one in neurobiology. This neuron has several logical inputs which can either be true or false and a logical output. Therefore, this neuron works as a linear classifier separating two categories where only one is the correct one. This is called binary classification. A schematic of this model can be seen in Fig. 3.2.

Because numerical values are required for later operations, every logical value is transformed to 0 if it is false or 1 if it is true. After summing up the inputs, a threshold

Table 3.1: Truth Tables of Logical Operations

(a) Logical AND					(b) Logical OR				
x_1	x_2	x_3	Thresh	Output	x_1	x_2	x_3	Thresh	Output
0	0		2	0	0	0		1	0
0	1		2	0	0	1		1	1
1	0		2	0	1	0		1	1
1	1		2	1	1	1		1	1
0	1	1	3	0	0	0	0	1	0
1	1	1	3	1	0	0	1	1	1

(c) Logical XOR				(d) Logical XNOR			
x_1	x_2	x_3	Output	x_1	x_2	x_3	Output
0	0		0	0	0		1
0	1		1	0	1		0
1	0		1	1	0		0
1	1		0	1	1		1
0	1	1	0	0	1	1	0
1	1	1	1	1	1	1	1

activation function is applied. In a mathematical sense

$$y = \phi \left(\sum_i^n x_i \right) \quad (3.1)$$

describes this operation where n is the amount of inputs, x_i a single input and ϕ the used activation function, in this case the threshold activation function. Plots of several activation functions including their equations are shown in Fig. ???. That means, if a given threshold is reached the output of the perceptron is 1 and 0 otherwise. This corresponds to the neurobiological spike of a neuron. The truth tables for the logical AND and OR operations are shown in Fig. 3.2a and Fig. 3.2b, respectively. The first one outputs true if all inputs are true. The second one outputs true if at least one input is true. It can be seen, that by changing the threshold value a different logical operation can be represented. For two inputs a threshold of 2 represents the logical AND operation. If the number of inputs are changed, a related change of the threshold still models the same operation. In this example the number of inputs is changed to 3 whereas the threshold must be changed to 3. The same procedure is valid for the logical OR operation. But here the threshold value is always 1 regardless of the amount of inputs and therefore differs to other logical operations. One limitation of this neuron type is its inability to represent exclusive logical operations like XOR and XNOR[2], whose truth tables are shown in Fig. 3.2c and Fig. 3.2d, respectively. The first one outputs true if the inputs are different. The second one outputs true if all inputs are identical which

equals an inverted XOR operation. It can easily be seen, that no threshold value can be found for meeting the requirements. For example, in the first case if the threshold is at 1 the first three combinations of inputs could be covered. If all inputs are false the threshold value is not reached and therefore the neuron outputs a 0. If the inputs differ from each other, the threshold value is reached which outputs a 1. But as soon as the fourth one needs to be classified the threshold value is no longer valid. The sum of the inputs equals two which exceeds the threshold and would output a logical true like in the OR case. But according to the truth table a logical false should be outputted. The reasons for this misbehavior are the neuron's limitation to binary values and a threshold activation function and therefore the classification of only two categories.

Donald Hebb states "When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased." [3] on how neurons learn. This means, that if a neuron repeatedly and persistently stimulates a immediately subsequent neuron, i.e. the more often two wired neurons are active, their synaptic efficacy increases. This is known as the Hebbian Theory. Hebb summarizes this with his famous quote "neurons that fire together, wire together".

Frank Rosenblatt developed the first perceptron [4]. Considering the Hebbian Theory the original McCulloch-Pitts-Neuron needs to be modified by adding associated weights for the inputs in order to simulate the strength of a synapse. Thus, Theorem 3.1 changes to

$$y = \phi \left(\sum_i^n x_i \cdot w_i \right) \quad (3.2)$$

by considering the weights w_i . Furthermore, the perceptron allows the usage of real-valued inputs and weights and uses the Heaviside step function as the activation function. According to Fig. ?? the Heaviside function outputs 0 if its parameter is negative and 1 otherwise. Therefore, its difference to the threshold activation function is just an offset of the threshold or bias, respectively. Adapting Fig. 3.1 to this results in Fig. ?. There is an input with the value 1 which is weighted by the bias. Thus, when multiplied representing the bias. This result is part of the weighted sum of the inputs which is fed to an activation function whose result is the output of the perceptron. Combining this with Theorem 3.2 yields

$$y = \phi(x_1 \cdot w_1 + x_2 \cdot w_2 + \dots + x_n \cdot w_n + b) \quad (3.3)$$

$$= \phi \left(\left(\sum_i^n x_i \cdot w_i \right) + b \right) \quad (3.4)$$

where x_i are the inputs, w_i the weights, b the bias and ϕ the Heaviside activation function.

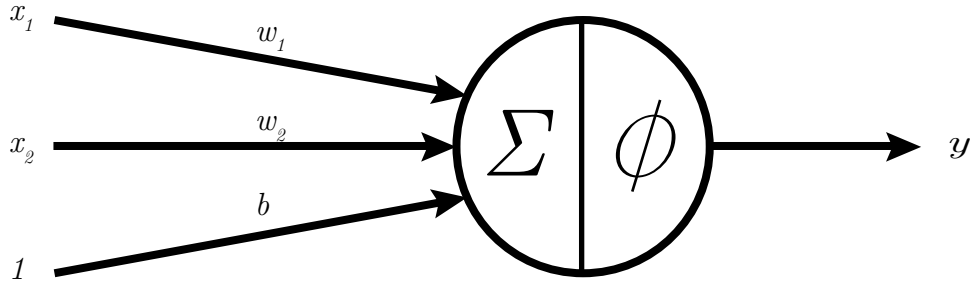


Figure 3.2: Model of a perceptron. The inputs x_i are weighted by w_i and are summed up with the bias b . This sum is put into the activation function ϕ whose result is the perceptron's output y .

Let's write the inputs and weights as vectors of

$$\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T \quad (3.5)$$

$$\mathbf{w} = (w_1 \ w_2 \ \cdots \ w_n)^T \quad (3.6)$$

for simplicity. Inserting this in Theorem 3.3 results in

$$y = \phi(\mathbf{x} \cdot \mathbf{w}^T + b) \quad (3.7)$$

with the same parameters as before. However, this model still works as a linear classifier and thus is unable to represent logical exclusive functions. This can be solved by concatenating multiple perceptrons and building a multilayer artificial network.

3.1.2.2 Multilayer Perceptron

A multilayer perceptron consists of multiple perceptrons divided in layers and solves complex tasks. It is an universal approximator for every function[5] regardless of the activation functions used[6]. Because of the multiple layers and the non-linear activation functions non-linearity is introduced into the network. Thus, it can distinguish data that is not linearly separable as most real world data is.

There are at least three layers. Each layer contains several perceptrons that are not connected to each other. However, every perceptron is connected to every perceptron of its subsequent layer. This type of connection is called fully-connected network. Because the data flow within the network is only in one direction and does not contain circles, the architecture is called feedforward neural network. A visualization of this is shown in Fig. 3.3. Although the weights are not displayed for clarity, they still exist and follow the same principle as with a single perceptron. Like the single perceptrons every node in the networks still holds a single numerical value. In this kind of network architecture perceptrons are often referred to as nodes.

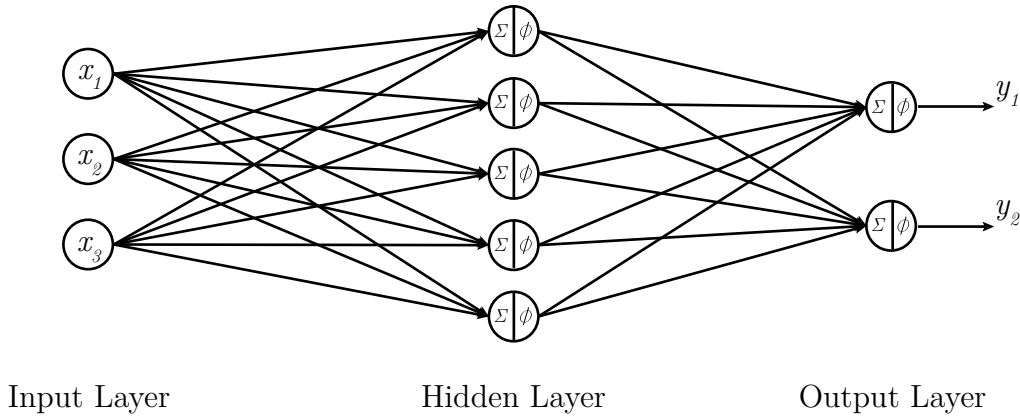


Figure 3.3: Multilayer Perceptron

The input layer serves as an interface for the data. It does not perform any calculations and just passes the data to the next layer. The number of nodes in this layer depends on the data and how it can be divided. If the real world data is an image, for example, the number of nodes should be equal to the number of pixels, so that every node can hold the intensity value of one pixel.

The output layer is responsible for transferring the network data to the outside so that it can be interpreted and worked with. The number of nodes in this layer depends on the expected results. If kinds of animals need to be detected in an image, every output node would represent a single kind or category, respectively. Let's say there are three kinds of animals possible, then there need to be three output nodes. In theory the node representing the correct kind of animal holds a one and every other a zero, if the values are squashed within this range.

Every layer between the input and the output layer is a hidden layer. They have no direct connection to the outside, neither to the input nor the output, hence, their name. Their task is to transfer the input information to the output by performing calculations. With at least one hidden layer every continuous function can be approximated.

Let's take again an image as an example. The task is to classify a handwritten digit from the MNIST dataset[7]. The digit can be seen in Fig. 3.4. Each grid cell represents a pixel. Because every pixel should have an associated node in the input layer of the network, this real world data is transferred into the network by flattening the intensity values of the image matrix to a vector. Therefore, the vector contains $28 \cdot 28 = 784$ elements which equals the number of input nodes. If the correct weights and biases of the network are found, it knows that if in another image more or less the same pixels or nodes, respectively, have high intensities, the same number is written. The downside of the flattening is, that every relation of pixels like position is lost, which means a loss in overall information. The consequence of this is, that if a digit is not centered and, for example, takes only half the space of the image, it cannot be classified, because the

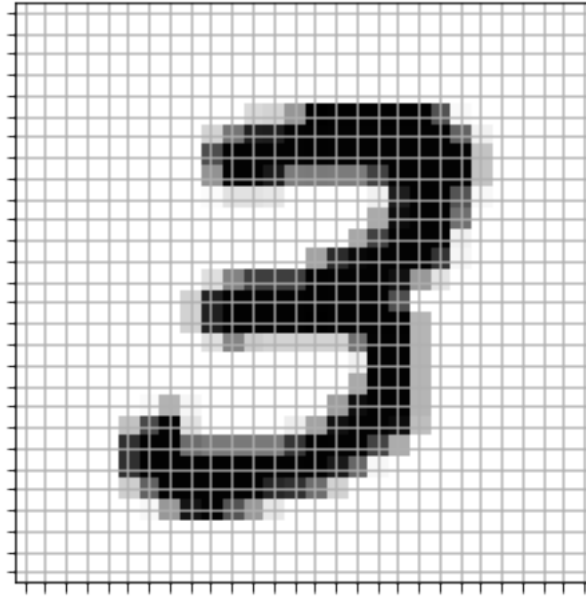


Figure 3.4: Handwritten Digit from the MNIST Digit Dataset. Represented as a 28×28 pixel matrix. Each cell represents a pixel.

network has never seen something similar to these input values.

3.1.2.3 Training

3.1.2.4 Improving Performance

3.1.3 Convolutional Neural Networks

3.1.3.1 Convolutional Layer

3.1.3.2 Pooling Layer

3.1.3.3 Fully-Connected Layer

3.2 Software

This section focuses on explaining which software and frameworks are used for implementing the network and generating the dataset.

Chapter 4

Methods

Chapter 5

Results

Chapter 6

Discussion

6.1 Conclusions

6.2 Outlook

Bibliography

- [1] McCulloch, Warren S. ; Pitts, Walter: Neurocomputing: Foundations of Research. Cambridge, MA, USA : MIT Press, 1988. – ISBN 0–262–01097–6, Kapitel A Logical Calculus of the Ideas Immanent in Nervous Activity, S. 15–27
- [2] Minsky, Marvin ; Papert, Seymour: *Perceptrons: An Introduction to Computational Geometry*. Cambridge, MA, USA : MIT Press, 1969
- [3] Hebb, Donald O.: *The Organization of Behavior: A Neuropsychological Theory*. Wiley, 1949 (A Wiley book in clinical psychology). – ISBN 9780471367277
- [4] Rosenblatt, Frank: The Perceptron: A Probabilistic Model for Information Storage and Organization in The Brain. In: *Psychological Review* (1958), S. 65–386
- [5] Cybenko, G.: Approximation by superpositions of a sigmoidal function. In: *Mathematics of Control, Signals and Systems* 2 (1989), Dec, Nr. 4, 303–314. <http://dx.doi.org/10.1007/BF02551274>. – DOI 10.1007/BF02551274. – ISSN 1435–568X
- [6] Hornik, Kurt: Approximation capabilities of multilayer feedforward networks. In: *Neural Networks* 4 (1991), Nr. 2, S. 251–257
- [7] Lecun, Yann ; Bottou, Léon ; Bengio, Yoshua ; Haffner, Patrick: Gradient-based learning applied to document recognition. In: *Proceedings of the IEEE*, 1998, S. 2278–2324