

# Machbarkeitsstudie: Smart Warehouse

## Echtzeit-Objektdetektoren im Vergleich

### Studienarbeit

im Rahmen der Prüfung zum  
**Bachelor of Science (B.Sc.)**

des Studienganges Angewandte Informatik  
an der Dualen Hochschule Baden-Württemberg Karlsruhe

von

**Felix Hausberger und Robin Kuck**

Oktober 2019 - Mai 2019

**-Sperrvermerk-**

Abgabedatum:	18. Mai 2020
Bearbeitungszeitraum:	30.09.2019 - 18.05.2020
Matrikelnummer, Kurs:	2773463, 4409176, TINF17B2
Ausbildungsfirma:	SAP SE Dietmar-Hopp-Allee 16 69190 Walldorf, Deutschland
Gutachter an der DHBW:	PD Dr.-Ing. Markus Reischl

# Eidesstattliche Erklärung

Wir versichern hiermit, dass wir unsere Studienarbeit mit dem Thema:

*Machbarkeitsstudie: Smart Warehouse*

gemäß § 5 der „Studien- und Prüfungsordnung DHBW Technik“ vom 29. September 2017 selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt haben. Die Arbeit wurde bisher keiner anderen Prüfungsbehörde vorgelegt und auch nicht veröffentlicht.

Wir versichern zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Karlsruhe, den 26. Oktober 2019

Gez. Felix Hausberger und Robin Kuck

Hausberger, Felix und Kuck, Robin

## Abstract

- English -

In this thesis the object detectors *Regional Convolutional Neural Networks*, *You Only Look Once* and *Single Shot MultiBox Detector* are compared for precision, reactivity and training behaviour and examined for their potential for industrial use. The background scenario of the *Smart Warehouse* offers live video data of a drone with goods in a warehouse, which are to be classified and localized in real time. In the future, this should make it possible to carry out inventories and inventory analyses of a warehouse in a time- and cost-efficient manner conserving resources.

The goal of this feasibility study is to find out whether the *Smart Warehouse* scenario is technically feasible and economically reasonable. In addition, the focus is also on the object detectors themselves and their differences in architecture and behavior in the *Smart Warehouse* environment.

## Abstract

- Deutsch -

In dieser Arbeit werden die Objektdetektoren *Regional Convolutional Neural Networks*, *You Only Look Once* und *Single Shot MultiBox Detector* nach Präzision, Reaktionsvermögen und Trainingsverhalten miteinander verglichen und auf deren Potential zum industriellen Einsatz untersucht. Das Hintergrundscenario des *Smart Warehouses* bietet dabei Live-Video Daten einer Drohne mit Warengegenständen in einem Warenhaus, die in Echtzeit klassifiziert und lokalisiert werden sollen. Dadurch sollen in Zukunft in der Industrie Inventuren und Bestandsanalysen eines Warenhauses zeit- und kostengünstig sowie ressourcenschonend ermöglicht werden können.

Diese Machbarkeitsstudie hat zum Ziel herauszufinden, ob das Szenario des *Smart Warehouse* technisch umsetzbar sowie wirtschaftlich sinnvoll ist. Zusätzlich liegt der Fokus ebenso auf den Objektdetektoren selbst und deren Unterschiede hinsichtlich Architektur und Verhalten im *Smart Warehouse* Umfeld.

# Inhaltsverzeichnis

<b>Abkürzungsverzeichnis</b>	<b>VI</b>
<b>Abbildungsverzeichnis</b>	<b>VII</b>
<b>Formelverzeichnis</b>	<b>VIII</b>
<b>Listenverzeichnis</b>	<b>IX</b>
<b>0 Vorwort</b>	<b>1</b>
<b>1 Einführung</b>	<b>2</b>
1.1 Forschungsumfeld . . . . .	2
1.2 Problemstellung und Motivation . . . . .	3
1.3 Vorgehensweise und Zielsetzung . . . . .	3
<b>2 Grundlagen und Forschungsstand</b>	<b>5</b>
2.1 Neuronale Netze . . . . .	5
2.2 Hyperparameter . . . . .	9
2.3 Objektdetektoren . . . . .	15
2.4 Verwendete Bibliotheken . . . . .	17
2.5 Compute Unified Device Architecture . . . . .	17
<b>3 Konzeption</b>	<b>18</b>
3.1 Bewertungskriterien . . . . .	18
3.2 Initialer Vergleich der Objektdetektoren . . . . .	18
3.3 Echtzeitumgebung . . . . .	18
<b>4 Realisierung</b>	<b>19</b>
4.1 Trainieren der Objektdetektoren . . . . .	19
4.2 Drohnen Anbindung . . . . .	19
4.3 Dashboard Entwicklung . . . . .	19
<b>5 Ergebnisse</b>	<b>20</b>

<b>6</b>	<b>Bewertung</b>	<b>21</b>
<b>7</b>	<b>Zusammenfassung und Ausblick</b>	<b>22</b>
	<b>Literaturverzeichnis</b>	<b>X</b>
<b>A</b>	<b>Anhang</b>	<b>XII</b>

# Abkürzungsverzeichnis

<b>ANN</b>	Artificial Neural Network
<b>CNN</b>	Convolutional Neural Network
<b>COCO</b>	Common Objects in Context
<b>LReLU</b>	Leaky Rectified Linear Unit
<b>PReLU</b>	Parametric Rectified Linear Unit
<b>PascalVOC</b>	Pascal Visual Object Classes
<b>R-CNN</b>	Regional Convolutional Neural Network
<b>ReLU</b>	Rectified Linear Unit
<b>SSD</b>	Single Shot MultiBox Detector
<b>YOLO</b>	You Only Look Once

# Abbildungsverzeichnis

2.1	Linear Threshold Unit . . . . .	6
2.2	Das einschichtige Perzeptron . . . . .	7
2.3	Gradientenverfahren . . . . .	11
2.4	Auswirkung unterschiedlicher Lernraten . . . . .	13
2.5	Sigmoid und Tangens Hyperbolicus . . . . .	14
2.6	ReLU-Aktivierungsfunktionen . . . . .	15
2.7	ELU . . . . .	16



# Formelverzeichnis

2.1	Die Heaviside-Funktion . . . . .	6
2.2	Die Softmax-Funktion . . . . .	6
2.3	Die RMSE-Funktion . . . . .	8
2.4	Neuberechnung der Gewichtungsmatrix durch partielle Differentiation . .	8
2.5	Superpositionsprinzip anhand der Varianz . . . . .	10
2.6	Standardverteilung nach Xavier Initialisierung . . . . .	10

# Listenverzeichnis

# 0. Vorwort

Besonderen Dank ist an unseren Betreuer PD Dr. -Ing. Markus Reischl auszusprechen, ohne den die folgenden Forschungsergebnisse nicht zustande gekommen wären. Auch dem Informatik Labor unter Enrico Hühneborg der DHBW ist für die nötige finanzielle Unterstützung zum Erwerb der Drohne zu danken.

# 1. Einführung

## 1.1. Forschungsumfeld

Einen Teilbereich des maschinellen Lernens (engl.: machine learning) stellt das *Deep Learning* dar, welches auf künstlichen neuronalen Netzen (engl.: artificial neural networks) (ANNs) basiert [1, S. 253]. Unter einer Vielzahl von Typen von ANNs wie Autoencodern, Deep Boltzmann Machines oder rekurrenten neuronalen Netzen befindet sich ebenso die Klasse der *Convolutional Neural Networks* (CNNs), welche hauptsächlich zur Lösung von Klassifikationsproblemen in der Audio-, Text- und Bildverarbeitung genutzt werden [2].

Ein Forschungsfeld im *Deep Learning* stellen Objektdetektoren dar, welche basierend auf CNNs neben Bildklassifikationsproblemen ebenso in der Lage sind, Lokalisationsprobleme zu lösen. Solchen Objektdetektoren werden in der heutigen Zeit immer mehr Bedeutung zugesprochen angesichts neuer Herausforderungen wie autonomen Fahren, automatisierter industrieller Verarbeitung oder aber auch staatlicher Überwachung. Verschiedene Ansätze werden zur Realisierung von Objektdetektoren verwendet, unter anderem Netzarchitekturen wie *You Only Look Once* (YOLO), *Single Shot MultiBox Detector* (SSD) oder *Regional Convolutional Neural Networks* (R-CNN).

Gerade in Zeiten des industriellen Wandels in Richtung *Industrie 4.0* können solche Objektdetektoren ein großes Optimierungspotential für bestehende Industrieszenarien bieten, beispielsweise in der Lagerhaltung und Logistik. Kombiniert mit einer autonomen Drohne können Objektdetektoren es ermöglichen, ohne menschliche Hilfe Inventuren und Bestandsprüfungen in einem Lager- oder Warenhaus durchzuführen. Start-up Unternehmen wie *doks. innovation* werben bereits mit ähnlichen Lösungen, die 80% Zeiteinsparung und 90% Kostensenkung versprechen [3]. Lösungen wie *inventAIRyX* beschränken sich allerdings speziell auf Lagerhäuser, in denen die verpackten Waren mittels Sensoren identifiziert werden, was Großhändler mit Warenhäusern wie *Baumarkt* oder *Selgros* ausschließt. Statt Waren mittels RFID Chips oder Barcodes zu identifizieren, soll in dieser Arbeit der Einsatz von Objektdetektoren für dieses Szenario evaluiert werden.

Wie sich die unterschiedlichen Objektdetektoren unter Echtzeitvoraussetzungen im Betrieb verhalten, soll anhand des Industriebeispiels *Smart Warehouse* innerhalb dieser Arbeit untersucht werden.

## 1.2. Problemstellung und Motivation

Das *Smart Warehouse* beschreibt ein Warenhaus, welches unter Einsatz einer Drohne in der Lage sei soll, Inventuren und Bestandsprüfungen weitgehend ohne menschliche Hilfe durchzuführen. Das Live-Bild der Drohne soll von den Objektdetektoren dazu genutzt werden, Warengegenstände zu lokalisieren und klassifizieren.

Neben der Frage, ob ein solches Industrieszenario überhaupt umsetzbar und wirtschaftlich sinnvoll ist, sollen die Objektdetektoren in diesem Anwendungsszenario nach verschiedenen Kriterien miteinander verglichen und beurteilt werden. Diese Kriterien lassen sich hauptsächlich in die Kategorien Präzision, Reaktionsvermögen und Trainingsverhalten untergliedern und werden später genauer eingeführt. Dadurch lassen sich Aussagen darüber treffen, ob nach dem momentanen Forschungsstand um Objektdetektoren solche das Potential bieten, industriell eingesetzt zu werden.

Falls die Machbarkeitsstudie des *Smart Warehouse* glückt, so kann der Industrie ein kostengünstiges, zeitsparendes und ressourcenschonendes Modell zur Inventurverwaltung eines Warenhauses angeboten werden.

## 1.3. Vorgehensweise und Zielsetzung

Zunächst muss sich mit den theoretischen Grundlagen von CNNs und Objektdetektoren auseinander gesetzt werden. Hierzu ist zunächst eine Einführung in neuronale Netz erforderlich, darunter zu Perzeptronen, dem Gradientenverfahren, dem Backpropagation Algorithmus und Hyperparametern zum Trainieren eines neuronalen Netzes.

Nachdem kurz auf den Grundbaustein moderner Objektdetektoren eingegangen wird, den CNNs, können anschließend die Funktionsweisen und Architekturen der drei miteinander verglichenen Objektdetektoren *YOLO*, *SSD* und Detektoren der *R-CNN* Familie erläutert werden. Hier ist zu bemerken, dass unterschiedliche Evolutionsstufen der drei

Detektoren zu betrachten sind.

Auch bisherige industrielle Einsatzfelder von Objektdetektoren sollen betrachtet werden, bevor technische Grundlagen zum Aufsetzen und Trainieren der drei Modelle näher betrachtet werden.

In der Konzeptionsphase sollen zunächst die Vergleichskriterien eingeführt werden und deren Metriken anschließend für initiale Benchmarkdatensätze für jeden Objektdetektor ermittelt werden. Hierzu wird auf die Datensätze *Pascal Visual Object Classes* (Pascal-VOC), *Common Objects in Context* (COCO) und ImageNet zurückgegriffen. Anschließend wird der Datensatz für das *Smart Warehouse* Szenario eingeführt.

In der Realisierung werden die Herausforderungen zur Steuerung und Anbindung der Drohne betrachtet und zudem die Objektdetektoren auf die realen Datensätze trainiert. Auch die Entwicklung der Webapplikation zur Visualisierung des Live-Bildes und der erkannten Objekte wird Bestandteil dieses Kapitels sein. Die Ergebnisse der Realisierungsphase werden im folgenden Kapitel dargestellt.

Ziel der Arbeit ist es Aussagen über die Fähigkeit von Objektdetektoren zum Einsatz in der Industrie zu treffen, indem eine Bewertung der Verhaltensweisen der Objektdetektoren nach den eingeführten Bewertungskriterien durchgeführt wird. Auch wirtschaftliche Gesichtspunkte werden in diesem Kapitel nicht außer Acht gelassen.

Zuletzt wird das Wesen der Arbeit nochmals kurz zusammengefasst und anschließend auf mögliche Verbesserungen und Ausblicke in die Zukunft aufmerksam gemacht.

## 2. Grundlagen und Forschungsstand

Im folgenden Kapitel soll sich speziell mit den Architekturen der unterschiedlichen Objektdetektoren auseinander gesetzt werden und wie sich diese voneinander abgrenzen. Außerdem wird grundlegendes Wissen über neuronale Netze und wie diese „lernen“ vermittelt, um die späteren Optimierungsverfahren an den Objektdetektoren zu verstehen. Zuletzt werden technische Grundlagen für die Programmierung des *Smart Warehouse* Szenarios vermittelt.

### 2.1. Neuronale Netze

Ein neuronales Netz bildet die Grundlage des *Deep Learnings* [1, S. 253]. Zunächst soll die einfachste Architektur eines neuronalen Netzes, das Perzeptron [1, S. 257], exemplarisch erklärt werden als auch der Lernprozess eines maschinellen Lernmodells an sich, um darauf basierend die Auswirkungen von Hyperparametern auf den Lernprozess des Modells zu erklären.

#### Das Perzeptron

Der Aufbau eines typischen Perzeptrons besteht aus einer oder mehreren Schichten sogenannter *Linear Threshold Units* (LTU) wie in Abbildung 2.1 dargestellt.

Es besteht aus  $n$  Eingängen mit  $x_i \in \mathbb{Q}$ , die im Inputvektor  $x$  zusammengefasst werden. Jeder Eingang wird mit einem Gewicht  $w_i$  aus dem Gewichtsvektor  $w$  versehen [1, S. 257 f.]. Die LTU berechnet das Skalarprodukt  $w^T \circ x$  aller Eingänge  $x$  mit ihren Gewichten  $w$  und wendet anschließend auf das Ergebnis  $z$  eine Aktivierungsfunktion an [1, S. 257 f.]. Das Ergebnis  $h_w(x)$  kann anschließend als Eingabe für ein weiteres Perzeptron dienen. Die einfachste Aktivierungsfunktion für ANNs ist die *Heaviside-Funktion* (siehe Formel 2.1) [1, S. 258 ff.]. Falls eine Klassifizierung mit Wahrscheinlichkeiten vorliegen soll, so ist die letzte Schicht eines Perzeptrons meist mit der *Softmax-Funktion* implementiert (siehe Formel 2.2), die den Wert des  $j$ -ten LTUs einer Schicht mit allen anderen  $n$  Werten

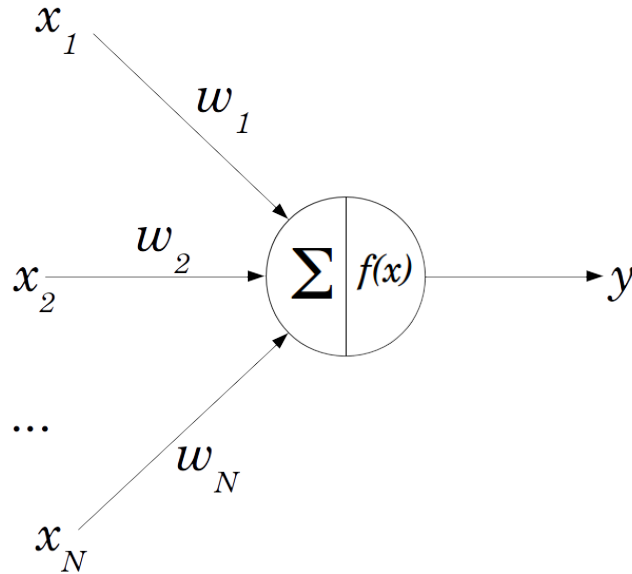


Abbildung 2.1.: Linear Threshold Unit [4]

der LTUs derselben Schicht ins Verhältnis setzt [1, S. 140 ff.]. Es gibt eine Vielzahl an möglichen Aktivierungsfunktionen, die im darauffolgenden Unterkapitel *Hyperparameter* betrachtet werden.

$$h_w(x) = \text{heaviside}(w^T \circ x) = \text{heaviside}(z) = \left( \begin{pmatrix} w_1 & w_2 & \dots & w_n \end{pmatrix} \circ \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \right) = \begin{cases} 1 & \text{wenn } z \geq 0 \\ 0 & \text{wenn } z < 0 \end{cases} \quad (2.1)$$

$$h_w(x) = \sigma(z)_j = \frac{e^{z_j}}{\sum_{i=0}^n e^{z_i}} \quad (2.2)$$

Die Aktivierung einer LTU hängt zusätzlich von einem Schwellwert  $\theta$  ab, der durch



einen sogenannten *Bias* festgelegt wird. Dies ist die Gewichtung des letzten Eingangs, der standardmäßig den Wert 1 liefert. Wählt man die Gewichtung negativ, so ist es schwieriger die LTU zu aktivieren, während eine positive Gewichtung die Aktivierung vereinfacht. [1, S. 258]

Nun bilden ein oder mehrere Schichten solcher LTUs ein Perzeptron. Jede einzelne LTU ist dabei mit allen LTUs der vorherigen Schicht verbunden (siehe Abbildung 2.2) [1, S. 258 f.]. Die beiden LTUs zur Ausgabe können dabei Aussagen über eine Klassifikation von Daten anhand der Eingangsdaten treffen, während die LTUs im Input Layer wesentlich die Daten weiter reichen. Die Verbindungen zur ersten Schicht des Hidden Layer sind stets mit Eins belegt. Existiert keine verborgene Schicht, so bezeichnet man das ANN als einschichtiges Perzeptron, ab einer oder mehr verborgenen Schichten spricht man bereits von einem *Multi-Layer Perzeptron* (MLP), einem mehrschichtigen Perzeptron [1, S. 261 f.]. Ist das neuronale Netz optimal trainiert, so ist am Ende nur eines der LTUs zur Ausgabe aktiviert. Das folgende ANN ist zudem ein Beispiel für ein sogenanntes *Feed Forward Network*, bei dem die Auswertung der Daten von einer Schicht zur nächsten weitergereicht wird, ohne zu bereits besuchten Schichten zurückzukehren [1, S. 263].

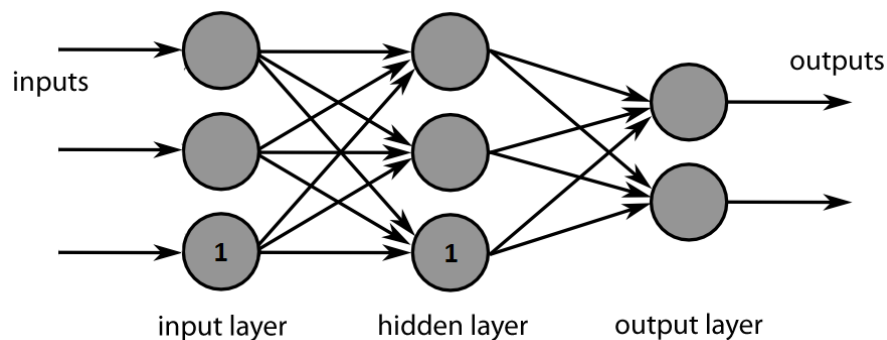


Abbildung 2.2.: Das einschichtige Perzeptron [Wikipedia.20190123]

## Gradientenverfahren und Backpropagation

Um zu verstehen, wie ein neuronales Netz durch Training lernt, muss zunächst der Begriff der Kostenfunktion (engl.: cost function) eingeführt werden. Die Kostenfunktion ist ein Qualitätsmaß dafür, wie weit die Ausgabe einer LTU vom erwarteten Wert abweicht [1,

S. 112 ff.]. Angenommen man übergibt dem neuronalen Netz also einen Datensatz zur Klassifikation, so ist am Ende meist nicht nur eine LTU zu Ausgabe aktiviert, was auf eine eindeutige Klassifikation schließen würde, sondern meist mehrere zu einem frühen Stadium des neuronalen Netzes.

Eine oft genutzte Kostenfunktion ist die *Root Mean Squared Error Funktion* (RMSE) (siehe Formel 2.3) [1, S. 37 f.].

$$E(z, o) = \sqrt{\frac{1}{n} \sum_{k=0}^n \|z_k - o_k\|^2} = \sqrt{\frac{1}{n} \sum_{k=0}^n (z_k - o_k)^2} \quad (2.3)$$

Hierbei ist  $z$  der erwartete Ausgabevektor des Perzeptrons, während  $o$  die momentane Ausgabe darstellt. Da der erwartete Ausgabewert bekannt ist, spricht man auch von sogenanntem überwachtem Lernen [1, S. 8 f.]. Den Fehler der Abweichung dieser beiden Werte gilt es nun schrittweise zu minimieren. Um dies zu erreichen können drei Parameter angepasst werden [1, S. 112 ff.].

1. Die Gewichtung der Verbindungen zum Perzeptron
2. Der Bias zur Aktivierung der LTUs des Perzeptrons
3. Die Stärke der Aktivierung des vorherigen Perzeptrons

Hierbei wird das sogenannte *Gradientenverfahren* eingesetzt. Es berechnet in einem iterativen Prozess über mehrere Testdaten das globale Minimum der Kostenfunktion nach den Gewichtungen der Verbindungen und damit auch nach den Bias Werten, die natürlich ebenso Gewichtungen darstellen. Ergebnis eines Durchlaufs im Gradientenverfahren ist die Gewichtungsmatrix, die die Änderung der Gewichtung jeder einzelnen Verbindung eines Perzeptrons zu jeder LTU des Folgeperzeptrons angibt (siehe Formel 2.4) [1, S. 112 ff.].

$$w_{ijt} = w_{ijt-1} - \eta \frac{\partial E}{\partial w_{ij}} \quad [4] \quad (2.4)$$

Das Gradientenverfahren eignet sich allerdings nur für stetig differenzierbare Funktionen ohne Plateaus. Somit können beispielsweise bei der Heaviside-Funktion als Aktivierungs-

funktion Probleme auftreten, da eine Ableitung der Kostenfunktion stets Null betragen würde, wohingegen bei der später eingeführten Sigmoid-Funktion im gesamten Definitionsbereich immer kleine Änderung der Gewichtungen zu verzeichnen wären. [1, S. 262]

Nun stellt sich auch der Vorteil von MSE als Kostenfunktion gegenüber anderen, durchaus komplexeren Kostenfunktionen heraus. Während MSE genau ein Minimum, das zugleich das globale Minimum der Funktion darstellt, besitzt, haben andere Kostenfunktionen im Gradientenverfahren das Problem, dass anstelle des globalen Minimums auch nur lokale Minima erreicht werden können [1, S. 114 f.]. Dies hat zur Folge, dass mehrere iterative Durchläufe mit mehreren Testdatensätzen nötig werden, um durch unterschiedliche Startkonfigurationen die unterschiedlichen Minima miteinander vergleichen zu können und damit das globale Minimum herauszustellen.

Durch das Gradientenverfahren werden somit nur diejenigen Verbindungen verstärkt, die zum richtigen Ergebnis führen.

Nun bleibt nur noch die dritte Möglichkeit zur Minimierung der Kostenfunktion übrig, die Anpassung der Stärke der Aktivierung des vorherigen Perzeptrons. Zu diesem Problem veröffentlichten David E. Rumelhart, Geoffrey E. Hinton und Ronald J. Williams 1985 den sogenannten *Backpropagation-Algorithmus* [5]. Dieser berechnet mit Hilfe des Gradientenverfahren welchen Anteil am Fehler der Ausgabe jede LTU des letzten Perzeptrons hat und anschließend welcher Anteil davon wiederum auf das vorherige Perzeptron der vergorenen Schicht zurück zu führen ist. Das Gradientenverfahren wird solange wiederholt, bis die Eingangsschicht erreicht wurde, es berechnet also für jede LTU deren Anteil am Fehler des Ergebnisses [1, S. 261 f.].

Mit Hilfe des Gradientenverfahren im Backpropagation Algorithmus wird nun also das neuronale Netz durch mehrere iterative Durchläufe trainiert, wobei das Training als Anpassung der Gewichtungen einzelner Verbindungen zu verstehen ist.

## 2.2. Hyperparameter

Als Hyperparameter versteht man die Parameter, die zur anfänglichen Konfiguration des neuronalen Netzes als auch zur Konfiguration des Lernprozesses heran gezogen werden. Um im Laufe der Arbeit verstehen zu können, wie die Objektdetektoren auf Seiten

der Netzarchitektur und des Lernverhaltens optimiert wurden, ist demnach ein kurzer Einblick in den Themenbereich der Hyperparameter von Nöten.

## Anzahl der LTUs

Die Anzahl der LTUs im ANN ist dafür ausschlaggebend, wie hoch der Komplexitätsanspruch eines Klassifizierungsproblems sein darf, um noch vom ANN gelöst werden zu können. Die Anzahl der LTUs hängt hauptsächlich von den Eingangsdaten ab. Über die optimalste Anzahl an LTUs pro Schicht lässt sich allerdings nur schwer etwas vorhersagen. Generell gilt, dass bei gleicher Anzahl an LTUs tiefere Netze eine weitaus höheren Parametereffizienz aufweisen als breitere Netze, da diese schneller gegen den gewünschten Zustand konvergieren. Zudem lassen sie sich somit schneller und kostengünstiger trainieren. So müssten bei einem 2x32 Netz 1024 Gewichtungen angepasst werden, während es bei einem 32x2 Netz dies nur 128 sind. [1, S. 271 f.]

## Initialisierung der Gewichtungen

Auch stellt die Initialisierung der Gewichte eines ANNs zu Beginn des Trainingsprozesses eine berechtigte Frage dar. Falls keine bereits trainierten ANNs für ein Klassifikationsproblem vorliegen, so werden die Gewichtungen meist zufällig nach einer Normalverteilung gewählt [1, S. 271].

Dies hat allerdings zur Folge, dass nach der Berechnung der gewichteten Summen aller LTUs die Werte der folgenden Schicht nicht mehr normalverteilt sind, da für die Varianz das Superpositionsprinzip gilt (siehe Formel 2.5)

$$\text{Var}(X + Y) = \text{VAR}(X) + \text{VAR}(Y) \quad (2.5)$$

$$W \sim U\left[-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}\right] \quad (2.6)$$

Durch die größer werdende Standardabweichung können demnach Werte entstehen, die

weit über den Mittelwert von Null hinaus gehen. Dies kann wiederum dazu führen, dass der Gradientenabstieg während des Backpropagation-Verfahrens nur langsam vollzogen werden kann, da der Gradient bei bestimmten Aktivierungsfunktionen (siehe Abbildung 2.5) gegen Null konvergiert [1, S. 275 f.].

Eine *Xavier Initialisierung* umgeht das Problem der sogenannten *schwindenden Gradienten*, indem die Gewichte nach 2.6 gleichverteilt werden, wobei  $n_j$  die Anzahl an LTUs der  $j$ -ten Schicht sind. [6, S. 253]

## Auswahl des Gradientenverfahrens

Generell unterscheidet man zwischen drei verschiedenen Arten das Gradientenverfahren durchzuführen (siehe Abbildung 2.3):

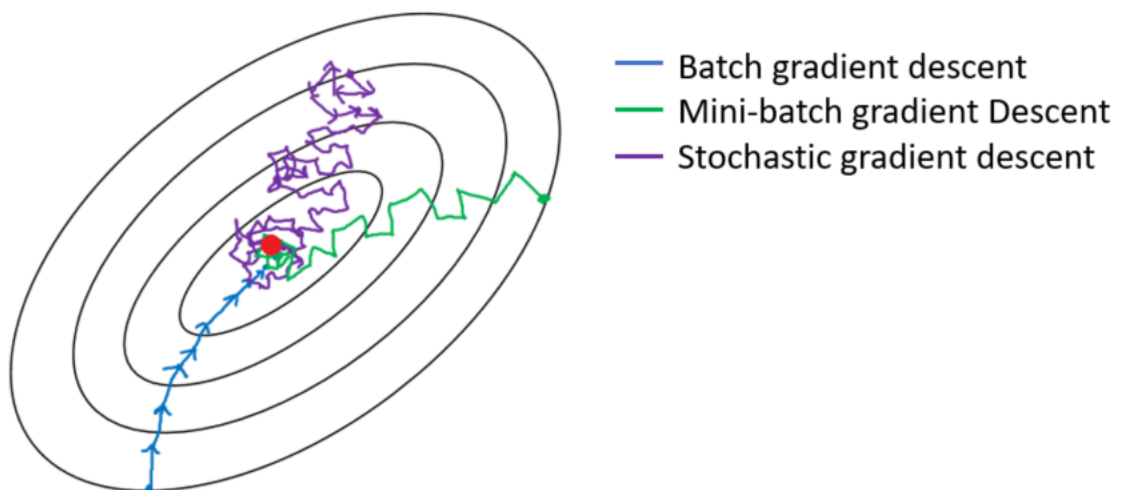


Abbildung 2.3.: Gradientenverfahren [7]

Beim *Batch* Verfahren werden in einem Trainingsdurchlauf, auch *Epoche* genannt, alle vorhandenen Daten des Trainingsdatensatzes herangezogen, um einen Gradientenabstieg zu vollziehen. Dies ist bei großen Trainingsdatensätzen auffällig langsam, dafür aber hinsichtlich der Erreichung des lokalen Minimums sehr zielstrebig. [1, S. 116]

Das *stochastische Gradientenverfahren* führt nach jedem einzelnen Dateneintrag im Trainingsdatensatz einen Gradientenabstieg durch. Da nur wenige Daten des ANNs verändert werden müssen, ist dieses Verfahren deutlich schneller, dafür aber unregelmäßiger hinsichtlich der Erreichung des Minimums. Oft wird das stochastische Gradientenverfahren verwendet, wenn nicht der komplette Trainingsdatensatz in den Hauptspeicher oder Grafikspeicher geladen werden kann. Diese Fähigkeit wird oft als *Out-of-Core* Fähigkeit bezeichnet. Es hat auch den Vorteil, besser das globale Minimum der Kostenfunktion aufzufinden, da bei lokalen Minima die Chance besteht, durch den unregelmäßigen Gradientenabstieg das lokale Minima wieder zu überwinden. [1, S. 118 f.]

Ein Kompromiss der beiden Verfahren bietet das *Mini-Batch* Verfahren, bei dem wiederholt Teilmengen des gesamten Datensatzes für einen Gradientenabstieg verwendet werden. Genauso wie das *Batch* Verfahren bietet das *Mini-Batch* Verfahren den Vorteil, die partiellen Ableitungen als Matrizenoperationen auf die Grafikkarten auszulagern, um die Performanz durch Parallelisierung zu steigern. [1, S. 121]

## Lernrate

Die Lernrate  $\eta$  gibt an, wie groß die Sprünge zum globalen Minimum sein sollen und damit indirekt wie viele Iterationen benötigt werden, um das globale Minimum der Kostenfunktion zu erreichen. Ziel der Anpassung einer Lernrate ist es, mit möglichst wenig Iterationen und Testdaten die optimale Konstellation des neuronalen Netzes zu berechnen. Deshalb wird sie standardmäßig zu Beginn der Iterationen groß gewählt um sich dem Minimum schnell zu nähern während sie am Ende immer kleiner gewählt wird, um nicht über das globale Minimum hinaus zu gehen. Dieses Vorgehen wird als *Simulated Annealing* bezeichnet, während das Funktion zum Festlegen der Lernrate als *Learning Schedule* betitelt wird. [1, S. 113f]

Eine Veranschaulichung der Anpassungen der Lernrate findet sich in Abbildung 2.4.

Die Anzahl der Durchläufe wird zu Beginn des Verfahrens zunächst hoch angesetzt, das Verfahren wird aber genau dann gestoppt, sobald der Gradientenvektor unter eine gewisse Abbruchgrenze fällt. Zwar ist das globale Minimum zu diesem Zeitpunkt noch nicht erreicht, allerdings kann es auch nie vollkommen erreicht werden, da die für das Gradientenverfahren genutzten Aktivierungsfunktionen nie einen partiellen Ableitungs-

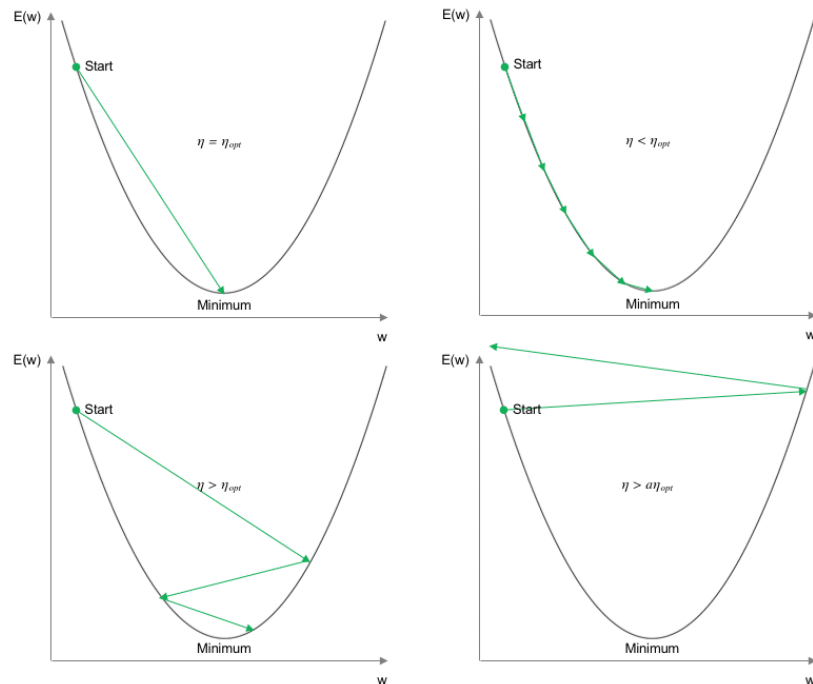


Abbildung 2.4.: Auswirkung unterschiedlicher Lernraten [8]

wert gleich Null zulassen [1, S. 118, S. 272]. In diesem Sinne wird auch von *Toleranz* gesprochen.

## Anzahl an Epochen

Die Anzahl der Epochen beschreibt die Durchläufe durch einen bestimmten Trainingsdatensatz während der Trainingsphase. Ist die Anzahl zu hoch gewählt wird Gefahr gelaufen sogenanntes *Overfitting* des ANNs zu erreichen. Dies bedeutet ein fehlendes Abstraktionsvermögen des ANNs zu erreichen und damit alleinig eine richtige Erkennung der Trainingsdatensätze zu ermöglichen.

## Aktivierungsfunktionen

Zwei bekannte und ähnliche Aktivierungsfunktionen sind die *Sigmoid-Funktion* und die *Tangens Hyperbolicus* Funktion (siehe Abbildung 2.5).

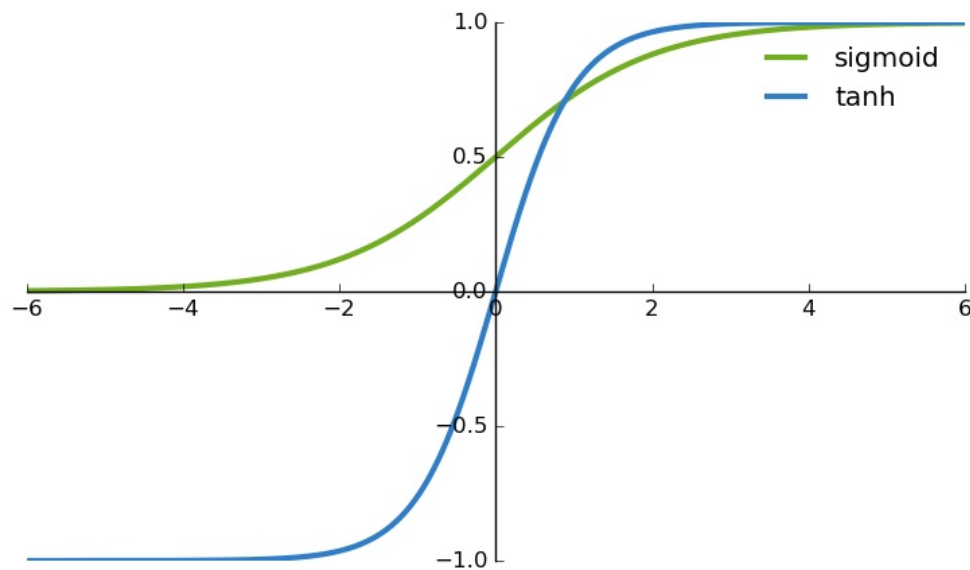


Abbildung 2.5.: Sigmoid und Tangens Hyperbolicus [9]

Da diese allerdings anfällig für das Problem *schwindender Gradienten* sind [1, S. 276], wird die *Rectified Linear Unit* (ReLU) bzw. *Parametric/Leaky Rectified Linear Unit* (PReLU/LReLU) Aktivierungsfunktion bevorzugt (siehe Abbildung 2.6).

Bei ReLU kann es während des Trainingsprozesses dazu kommen, dass LTUs nach dem Gradientenabstieg einen negativen Wert aufweisen, weshalb sie nicht weiter aktiviert werden und für den Rest der Trainingsdauer „tot“ sind. Um dies zu verhindern wurde *LReLU* dazu genutzt, um eine Reaktivierung zu ermöglichen, da auch für negative LTU Werte ein Gradient der Aktivierungsfunktion bestimmt werden kann. Bei *LReLU* ist die Steigung der Funktion im zweiten Quadranten statisch gewählt, während sie bei *PReLU* dynamisch von neuronalen Netz während des Trainingsprozesses selbst gelernt werden kann. [1, S. 280 f.]

Eine letzte Variante der Aktivierungsfunktionen beschreibt die *ELU* Funktion (siehe Abbildung 2.7).

Sie besitzt nicht nur die Eigenschaft schwindende Gradienten und tote LTUs zu ver-



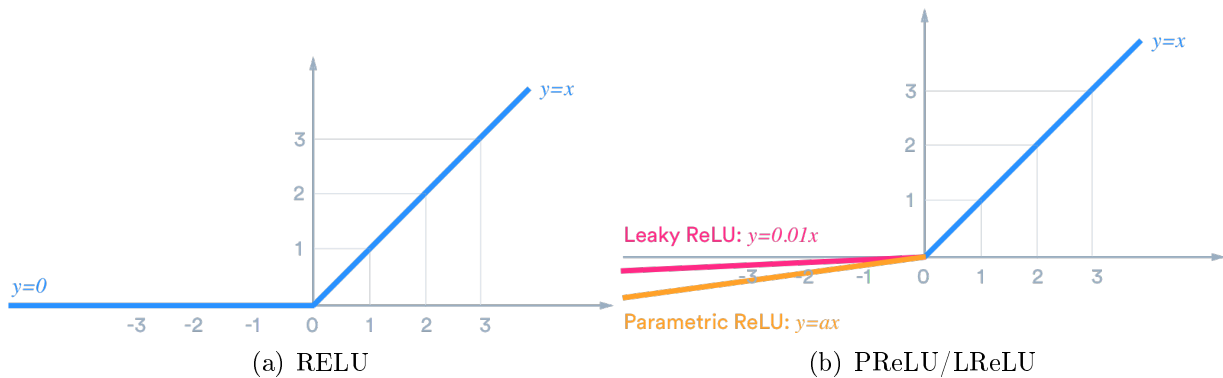


Abbildung 2.6.: ReLU-Aktivierungsfunktionen [10]

hindern, sondern ist im gesamten Definitionsbereich ebenso eine stetig differenzierbare Funktion, was das Gradientenverfahren beschleunigt. Als Standardwert für  $\alpha$  wird oft Eins verwendet. Nachteil der *ELU* Funktion ist der erhöhte Rechenaufwand, was aber durch die schnellere Konvergenz kompensiert wird. [1, S. 280 f.]

## 2.3. Objektdetektoren

### 2.3.1. Convolutional Neural Networks

Auch kann bei tiefen ANNs eine gewisse Klassifikationshierarchie für jede Schicht zugeordnet werden. Bei der Bilderkennung sind beispielsweise die ersten verborgenen Schichten dafür zuständig kleinere Muster zu erkennen, während mit fortschreitenden Schichten diese Muster zu immer größeren Mustern zusammengefasst werden können. [1, S. 271 f.]

Bei einem 28x28 Pixel großen schwarz-weiß Bild sind beispielsweise für die Modellierung des Grauwertes jedes Pixels insgesamt 784 Input LTUs notwendig. Da die erkannten Strukturen von vielen detaillierten nach und nach zu wenigen verallgemeinerten zusammengefasst werden können, sinkt die Anzahl an benötigten LTUs pro Schichten in einem Feed Forward Netz zur Ausgabeschicht. Ein trichterförmiger Aufbau des ANNs ist somit nicht unüblich. [1, S. 271 f.]

Diese Architektur ermöglicht ebenso die Wiederverwendbarkeit einzelner Schichten und Gewichtungen für ähnliche Klassifikationsprobleme, bei denen gleiche Muster vorzufinden

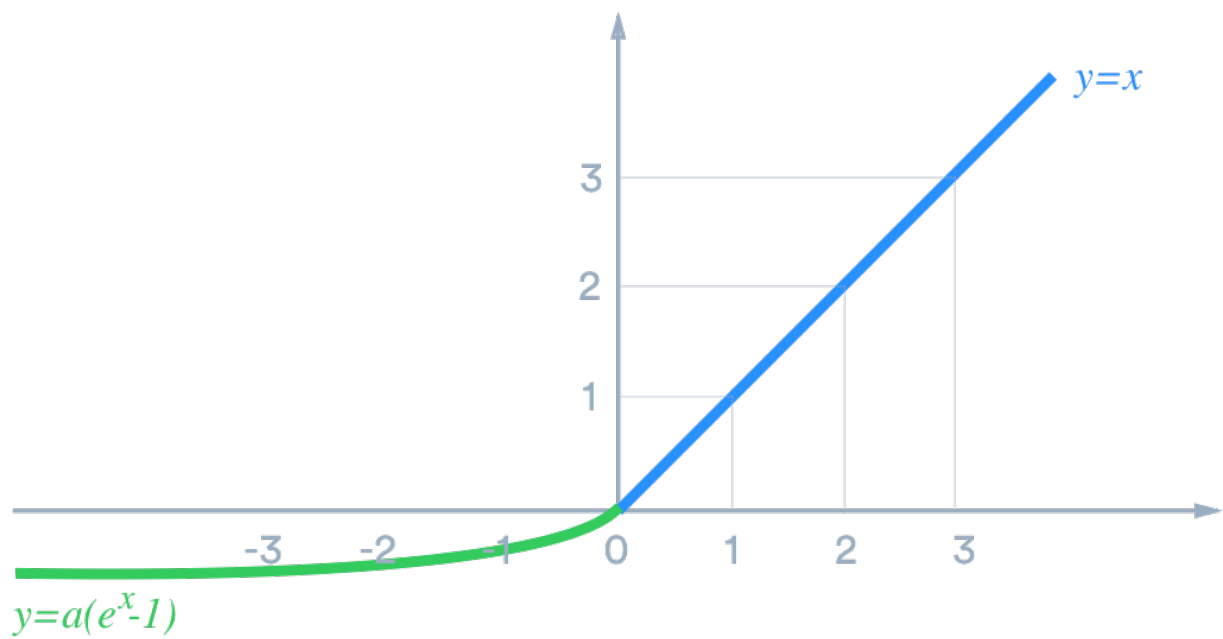


Abbildung 2.7.: ELU [10]

sind. [1, S. 271]

### 2.3.2. Mask Regional Convolutional Neural Network

### 2.3.3. You Only Look Once

### 2.3.4. Single Shot MultiBox Detector

## 2.4. Verwendete Bibliotheken

### 2.4.1. PyTorch

### 2.4.2. OpenCV

### 2.4.3. Django

## 2.5. Compute Unified Device Architecture

- Methoden
- Literaturrecherche
- Herausarbeitung des Neuheitswertes

## 3. Konzeption

### 3.1. Bewertungskriterien

- Kriterien für spätere Bewertung einführen

### 3.2. Initialer Vergleich der Objektdetektoren

- Verfahren incl. Variantendiskussionen
- Funktionsweise an Benchmarkdatensatz erklären
- Parameter d. Methode diskutieren und u.U. bestimmen

### 3.3. Echtzeitumgebung

- Richtigen Datensatz beschreiben

## 4. Realisierung

### 4.1. Trainieren der Objektdetektoren

### 4.2. Drohnen Anbindung

### 4.3. Dashboard Entwicklung

- Herausforderungen bei der Drohne
- Trainieren auf die richtigen Datensätze
- Webapplikation

## 5. Ergebnisse

- Ergebnisse auswerten

## 6. Bewertung

- Ergebnisse bewerten und diskutieren
- Was liefert die Arbeit, was bisher noch nicht bekannt war?

## 7. Zusammenfassung und Ausblick

- Klare Darstellung, was die Arbeit geliefert hat
- Ca. 2-4 Anpunkte: Zukünftige Ziele



# Literaturverzeichnis

- [1] Aurélien Géron. *Machine Learning mit Scikit-Learn & TensorFlow: Konzepte, Tools und Techniken für intelligente Systeme: Übersetzung von Kristian Rother*. 1. Aufl. Heidelberg: dpunkt.verlag GmbH, 2018.
- [2] MathWorks. *Deep Learning: Drei Dinge, die Sie wissen sollten*. MathWorks, 2019. URL: <https://de.mathworks.com/discovery/deep-learning.html> (Einsichtnahme: 12.10.2019).
- [3] doks.innovation. *inventAIRyX*. doks.innovation Homepage, 2019. URL: <https://www.doks-innovation.com/inventairy-x/> (Einsichtnahme: 26.10.2019).
- [4] Philippe Lucidarme. *Simplest perceptron update rules demonstration*. Homepage Blog, 2017. URL: <https://www.lucidarme.me/simplest-perceptron-update-rules-demonstration/> (Einsichtnahme: 26.10.2019).
- [5] David E. Rumelhart/ Geoffrey E. Hinton/ Ronald J. Williams. *Learning Internal Representations by Error Propagation*. Hrsg. von University of California, San Diego. 09/1985. URL: <https://apps.dtic.mil/dtic/tr/fulltext/u2/a164453.pdf> (Einsichtnahme: 26.10.2019).
- [6] Xavier Glorot, Y. B. „Understanding the difficulty of training deep feedforward neural networks“. Diss. Montréal: Université de Montréal, 2010. URL: <http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf> (Einsichtnahme: 26.10.2019).
- [7] Imad Dabbura. *Gradient Descent Algorithm and Its Variants*. Towards Data Science, 2017. URL: <https://towardsdatascience.com/gradient-descent-algorithm-and-its-variants-10f652806a3>.
- [8] Sebastian Heinz. *Wie lernen neuronale Netze?* STATWORX Blog, 2018. URL: <https://www.statworx.com/de/blog/wie-lernen-neuronale-netze/> (Einsichtnahme: 26.10.2019).
- [9] Ronny Restrepo. *Calculating the derivative of Tanh: step by step*. Ronny Restrepo Homepage, 2017. URL: [http://ronny.rest/media/blog/2017/2017\\_08\\_16\\_tanh/tanh\\_v\\_sigmoid.jpg](http://ronny.rest/media/blog/2017/2017_08_16_tanh/tanh_v_sigmoid.jpg).

- [10] Danqing Liu. *A Practical Guide to ReLU*. Medium, 2017. URL: <https://medium.com/@danqing/a-practical-guide-to-relu-b83ca804f1f7> (Einsichtnahme: 26.10.2019).

# A. Anhang

## Abbildungen