
DLM Semesterprojekt: Bestimmung der 3D-Boundingbox eines Roboters

GruppeNR: 1, Namen: Tim Eisenacher (1870755) und Paul Manea()
EMail: 1870755@stud.hs-mannheim.de

Abstract

This Latex and Style file are modified versions from NeurIPS 2018.[1]

1 Einleitung

Objekterkennung ist einer der vielversprechendsten und am stärksten im Forschungsfokus stehenden Bereiche von *Computer Vision* und *Machine Learning* [7]. Besonders im Bereich der medizinischen und industriellen Innovationen konnten dadurch bereits jetzt schon große Fortschritte erzielt werden [9, 6]. So können beispielhaft Operationsroboter stark von einer auf *Deep-Learning*-Algorithmen basierenden Erkennung und Lokalisation der Operationsinstrumente profitieren [8]. Auch in der Tumordiagnostik kann so die Erkennungsgenauigkeit gegenüber einer menschlichen Klassifikation bei gleichzeitig signifikant niedrigeren Kosten verbessert werden [2].

Die Anforderungen an die Bildverarbeitung sind dabei in diesen Bereichen besonders groß. Zum einen sorgen immer hochauflösendere Bilder für enorme Datenmengen und damit Hardwareanforderungen. Zum anderen ist die Anzahl an möglichen Klassen und damit Featurekombinationen bei der Objekterkennung enorm. Aufgrund ihrer hohen Trainingsperformanz und der Fähigkeit herausragend effizient Features in Bildern zu erkennen, eignen sich *Convolutional-Neural-Networks* (CNN) besonders gut als Lösungsansatz der beschriebenen Probleme [7].

In der vorliegenden Arbeit wird ein *Deep-Learning*-Modell zur Bestimmung der *3D-Boundingbox* eines Roboters implementiert und evaluiert. Dafür erfolgt anfangs die Abgrenzung und grobe Erläuterung des *Deep-Learnings* (DL) im Kontext des *Machine-Learnings* (ML). Der Fokus liegt dabei auf dem verwendeten DL-Konzept der CNN's und damit verbundener Probleme und Herausforderungen. Anschließend erfolgt eine Abhandlung der für die Implementierung verwendeten Netzstruktur und Metriken. Die Implementierung wird zunächst zur Lösung eines zweidimensionalen Problems erstellt und dann anschließend auf drei Dimensionen erweitert. Die Entwicklung erfolgt in Python 3.7 unter Verwendung des Tensorflow-Frameworks in der Spider und Eclipse IDE. Die Evaluation der Implementierung geschieht anhand eines Vergleiches der geschätzten *Boundingbox* mit der gelabelten *Boundingbox* der Testdaten. Als Datengrundlage für Training und Evaluation dient ein synthetisch generierter und vorgelabelter Datensatz mit RGB-Bildern.

2 Grundlagen und Stand der Technik

Salvaris beschreibt *Machine Learning* als einen Zweig der Computerwissenschaften, bei dem Computern beigebracht wird anhand von Trainingsdaten Entscheidungen zu treffen. Typische Anwendungsgebiete des ML sind Klassifikation, Regression und Clustering. DL ist ein Teilgebiet des ML bei dem besonders komplexe Neuronale Netze mit vielen Schichten und Neuronen Verwendung finden. Ein weitere wesentliche Abgrenzung stellt die Merkmalsextraktion dar. Also die Extraktion der Eigenschaften eines Objekts, die ausschlaggebend für etwaige Klassenzugehörigkeiten sind. Diese entscheidenden Merkmale müssen dem DL-Modell nicht vorgegeben werden, sondern werden von dem Algorithmus selbst gefunden. Dieser Umstand stellt mit die größten Herausforderungen aber auch Chancen beim DL dar [5, S.32-47]. Im Folgenden werden nur die für die vorliegende Arbeit besonders relevanten und speziell angepassten Methoden und Aspekte des DL erläutert. Für weiterführende grundlegende Informationen zum Beispiel zu Neuronen, Schichttypen, Aktivierungsfunktionen und zum Gradientenabstiegsverfahren wird auf [4] verwiesen.

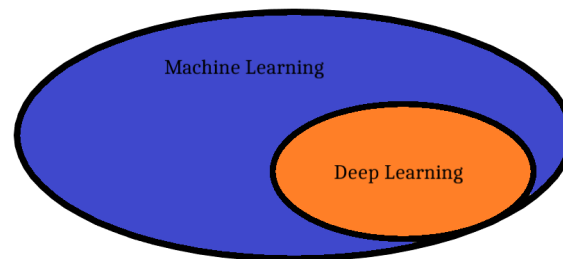


Figure 1: Abgrenzung Deep Learning zu Machine Learning.

2.1 Convolutional Neural Networks

Convolutional Neural Networks (CNN) sind eine spezielle Art von Neuronalen Netzen, die sich für das verarbeiten von gitterartig beschaffenen Daten eignen. Hierzu zählen zum Beispiel auch Bilddaten, deren Pixelraster sich als Gitter oder Matrix interpretieren lassen. Ein typisches CNN besteht dabei aus einem oder mehreren Paaren von *Convolutional*- und *Pooling-Layers*, gefolgt von einem oder mehreren *Fully-Connected-Layers*. Die Folgenden Darstellungen richten sich im wesentlichen nach Goodfellow [3, S.326-366]

Convolutional Layer Bei einem *Convolutional Layer* wird schrittweise ein Filterkernel K über eine Eingabematrix I mit den Dimensionen n und m bewegt (Abb. 2.1). Der Input der folgenden Neuronen $S(i, j)$ berechnet sich dann aus einer Faltungsoperation der jeweils übereinanderliegenden Kernel- und Bildelemente (Gleichung 1).

$$S(i, j) = (I * K)(i, j) = \sum_n \sum_n I(i - m, j - n) K(m, n) \quad (1)$$

Der so berechnete Input eines Neurons wird anschließend abhängig von der verwendeten Aktivierungsfunktion in den Output verwandelt. Zu bemerken ist, dass alle Neuronen eines *Convolutional Layers* die gleichen Gewichte haben (sog. *Parameter Sharing*). Dadurch ist es möglich Speicher gegenüber anderen Netzstrukturen einzusparen, die häufig eine große Gewichtungsmatrix verwenden. Ein weiterer großer Vorteil sind die sog. *Sparse Interactions*. Durch die Verwendung eines Filterkernels der meist nur einen Bruchteil der Größe des zu analysierenden Bildes aufweist, werden nur die Features extrahiert, die wirklich entscheidend sind für die Zugehörigkeit zu einer Klasse. Dies führt ebenso zu einer weiteren Speicher- und Performanzoptimierung.

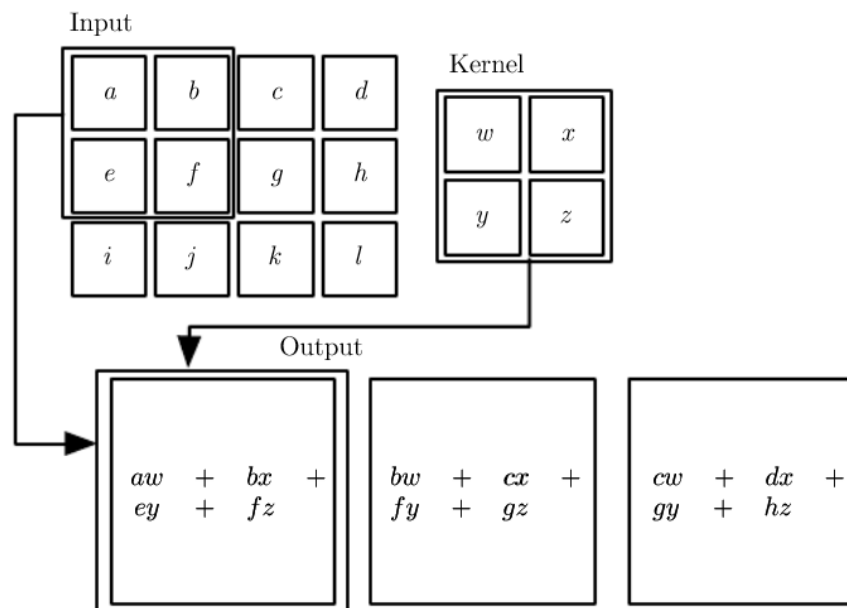


Figure 2: Prinzip eines Convolutional Layers [3, S.330].

Pooling Layer *Pooling Layer* sorgen dafür, dass Features einer Klasse in einem Bild nahezu ortsinvariant gelernt werden können. Ein weitverbreitetes Pooling Verfahren ist das sog. *2X2 Max-Pooling*, bei dem aus jedem 2X2 Quadrat der Neuronen des *Convolutional-Layers* nur das aktivste Neuron an die nächste Schicht weitergeleitet wird. Abbildung 2.1 verdeutlicht dieses Funktionsprinzip. Es werden von den jeweils benachbarten Neuronen nur die mit den höchsten Gewichten an die nächste Schicht durchgeschaltet.

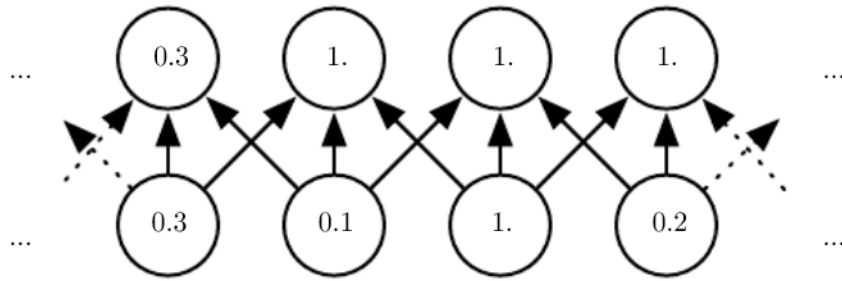


Figure 3: Prinzip eines Pooling Layers [3, S.337].

Fully Connected Layer *Fully-Connected-Layer* oder in *Keras* sog. *Dense-Layer* stellen einen Schichtentyp dar, bei dem jedes Neuron mit jeweils jedem Neuron der vorigen Schicht verschaltet ist. Es ist so möglich, die Ausgaben des letzten *Pooling-Layers* über ein oder mehrere *Fully-Connected-Layer* mithilfe von Aktivierungsfunktionen zum Beispiel in eine Wahrscheinlichkeitsverteilung der Klassenzugehörigkeit zu überführen. Die Anzahl Neuronen in der letzten Schicht entspricht dann der Anzahl zu lernender Klassen oder auch der Anzahl vorherzusagender Features.

2.2 Loss-Funktionen und Metriken

3 Methoden

3.1 Interception over Union

Netzstruktur

4 Ergebnisse

4.1 Headings: second level

4.1.1 Headings: third level

5 Fazit und Ausblick

Vielleicht ganz kurz

Paragraphs

6 Citations, figures, tables, references

6.1 Citations

References

- [1] Md Atiqur, Rahman, Yang, and Wang. Optimizing intersection-over-union in deep neural networks for image segmentation. 2015.
- [2] Angel Cruz-Roa, Hannah Gilmore, Ajay Basavanahally, Michael Feldman, Shridar Ganesan, Natalie N.C. Shih, John Tomaszewski, Fabio A. González, and Anant Madabhushi. Accurate and reproducible invasive breast cancer detection in whole-slide images: A deep learning approach for quantifying tumor extent. *Scientific Reports*, 7(1), apr 2017.
- [3] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [4] Aaron Courville Ian Goodfellow, Yoshua Bengio. *Deep Learning. Das umfassende Handbuch: Grundlagen, aktuelle Verfahren und Algorithmen, neue Forschungsansätze (mitp Professional)*. mitp Professionals, 2018.

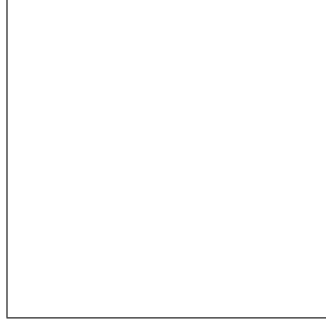


Figure 4: Sample figure caption.

Table 1: Sample table title

Part		
Name	Description	Size (μm)
Dendrite	Input terminal	~ 100
Axon	Output terminal	~ 10
Soma	Cell body	up to 10^6

- [5] Wee Hyong Tok Mathew Salvaris, Danielle Dean. *Deep Learning mit Microsoft Azure*. Rheinwerk Computing, 2019.
- [6] Ahmad Zaib Muhammad Imran Razzak, Saeeda Naz. Deep learning for medical image processing: Overview, challenges and future. 2016.
- [7] Wanli Ouyang, Xiaogang Wang, Xingyu Zeng, Shi Qiu, Ping Luo, Yonglong Tian, Hongsheng Li, Shuo Yang, Zhe Wang, Chen-Change Loy, and Xiaoou Tang. Deepid-net: Deformable deep convolutional neural networks for object detection. 2014.
- [8] Robot-Assisted Surgery, Using Deep, and Learning. Automatic instrument segmentation in. 2018.
- [9] Jonathan Tremblay, Thang To, and Balakumar Sundaralingam. Deep object pose estimation for semantic robotic grasping of household objects. 2018.

6.2 Footnotes

Here are some samples you can use. Please delete.

6.3 Figures

6.4 Tables

All tables must be centered, neat, clean and legible. The table number and title always appear before the table. See Table 1.

Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lower case (except for first word and proper nouns); tables are numbered consecutively.

Note that publication-quality tables *do not contain vertical rules*. We strongly suggest the use of the booktabs package, which allows for typesetting high-quality, professional tables:

<https://www.ctan.org/pkg/booktabs>

This package was used to typeset Table 1.