

Berliner Hochschule für Technik Berlin
Fachbereich VI – Informatik und Medien

Bachelorarbeit

Erzeugung von Bildern mittels Neuronalen Netzen

Stefan Berger
Medieninformatik
Matrikel-Nr. 854184

Berlin, 6. April 2021

Betreut von: Prof. Dr. F. Gers

Gutachter: Prof. Dr. J. Schimkat

Zusammenfassung

Im Experiment und im Inhalt dieser Bachelorarbeit werden die Fragen

1. ...

2. ...

...

beantwortet.

Inhaltsverzeichnis

1	Einleitung	1
1)	Ziel der Arbeit	1
2)	Bisherige Arbeiten	2
2	Modell	3
1)	Logistische Regression	3
2)	Deep Neural Networks	4
3)	Convolutional Neural Networks	4
4)	U-Net-Architektur	4
5)	Generative Adversarial Networks	4
6)	Image-To-Image-Translation	5
3	Implementierung	6
1)	Entwicklungsumgebung	6
4	Experimente und Resultate	7
1)	Vorbereitung der Eingabedaten	7
2)	Anwendung herkömmlicher Shader	7
3)	Hyperparameter	7
4)	Performancebeobachtungen	8
5	Diskussion	9
	Literatur	11
	Bildnachweis	13
	Abkürzungs- und Symbolverzeichnis	15
	Abbildungsverzeichnis	17
	Tabellenverzeichnis	19

1 Einleitung

Obwohl die Idee für eine Maschine, die anhand eingegebener Daten selbständig Entscheidungen treffen kann und die ersten praktischen Ansätze für künstliche neuronale Netze schon einige Jahrzehnte alt sind, findet der Einsatz derartiger Algorithmen erst seit einigen Jahren statt. Viele Erfindungen, die vor 30 bis 50 Jahren in Filmen und Serien Science Fiction darstellten, sind inzwischen nicht nur Realität, sondern auch alltagstauglich. Zu den wichtigsten Beispielen zählen verbale Schnittstellen an Computersystemen und auch Armbanduhren, autonome Fahrzeuge etwa in Gestalt von Parkassistenten und verschiedene Verfahren zur biometrischen Identitätsprüfung.

Andere rasche technologische Fortschritte aus der jüngeren Vergangenheit haben nicht immer nur die Lebensqualität der beteiligten Personen erhöht, sondern stellten durch Missbrauch gelegentlich sogar Gefahren dar. So wie das Internet auch zur Verbreitung von Falschinformationen und die sichere Verschlüsselung von gespeicherten Daten auch für Erpressungen genutzt werden kann, ist die Generierung von täuschend echten Bildern unter Umständen geeignet, persönlichen, finanziellen oder anders gearteten Schaden zu verursachen.

Weiterhin existieren bei der Auswahl der Trainingsdaten für künstliche neuronale Netze rechtliche Grenzen. Bilddaten sind in mehr als ausreichenden Mengen vorhanden, berücksichtigen aber zum Beispiel nicht immer das Recht am eigenen Bild. Für diese Bachelorarbeit sind die Anforderungen an die Bildqualität außerdem sehr hoch, da möglichst auch Texturen und Lichtreflexionen erlernt werden sollen. Modellgrafiken mit geeigneten Material- und Beleuchtungseigenschaften gibt es zwar auch, aber nur in weitaus geringeren Mengen. Für solche Fremdarbeiten, die sehr arbeitsaufwendig sind, wäre auch die Klärung der Nutzungsrechte erforderlich geworden. Für Trainingsergebnisse, die das mentale Modell einer breiten Nutzergemeinschaft reflektieren, sind grundsätzlich auch Daten aus möglichst vielen verschiedenen Quellen erforderlich.

Es ist deshalb eine Brücke geschlagen worden zwischen Trainingsdaten, die zum einen aus zufällig ausgewählte Benutzereingaben bestehen, und solchen, die bestimmte Qualitätseigenschaften erfüllen und in beliebiger Menge erstellt werden können. Ein künstliches neuronales Netz soll aus Skizzen des "Quick, Draw!"-Datasets von Google hochwertig gestaltete Figuren generieren. TODO: Hierfür kommen als Ein- und Ausgabedaten sowohl Bilddateien als auch die jeweils zugrundeliegenden Dateiformate NDJSON und Wavefront OBJ infrage.

1) Ziel der Arbeit

Die Arbeit verfolgt das Ziel, verschiedene bewährte Architekturen künstlicher neuronaler Netze zur Generierung von Bildern zu untersuchen und in einem praxisorientierten Zusammenhang zu testen. Ich zeige mehrere Möglichkeiten, wirklichkeitsnahe Bilder von Alltagsgegenständen aus Skizzen, die durch Benutzer erstellt wurden, mittels eines künstlichen neuronalen Netzes zu generieren. Die Bilddateien der Skizzen sowie die generierten Bilddatei-

en können dabei aus RGB-Pixelinformationen bestehen oder die Grafik mittels Bildkoordinaten beschreiben, wie es bei Vektorgrafiken und 3D-Modellen der Fall ist. Eine Anwendung der Ergebnisse ist beispielsweise als Feature eines Grafiktablets oder eines Smartboards denkbar.

2) Bisherige Arbeiten

Künstliche neuronale Netze finden erst seit wenigen Jahren breite Aufmerksamkeit, seit auch Heimcomputer in der Lage sind die hohe Anzahl der erforderlichen Rechenoperationen in annehmbarer Zeit auszuführen. Seitdem sind wenige, englischsprachige Einführungen in die Thematik entstanden. Ein häufig genanntes Buch ist "Deep Learning", das online kostenfrei zugänglich ist [1]. Ebenfalls online kostenfrei ist das Buch "Dive into Deep Learning" [2]. Ein weiteres, praxisorientierteres Buch ist "Deep Learning with Python" [3], dessen zweite Auflage bald herausgegeben wird.

Als Architektur für ein künstliches neuronales Netz kommt für diese Arbeit zunächst derselbe Aufbau wie in "Image-To-Image-Translation" [4] zum Einsatz. Dieser setzt seinerseits auf Generative Adversarial Networks [5] und Conditional Generative Adversarial Networks [6] auf.

In "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks" [7] werden ebenfalls realitätsnahe Ergebnisse erzielt. Anders als bei den bisher erwähnten Arbeiten wird Unsupervised Learning verwendet, um das künstliche neuronale Netz zu trainieren.

2 Modell

1) Logistische Regression

Ein vergleichsweise einfacher Algorithmus des Maschinellen Lernens ist Logistic Regression oder Softmax Regression. Es lassen sich damit Klassifizierungen durchführen. Bei binären Klassifizierungen werden die Eingaben in zwei Kategorien unterteilt. Häufige Beispiele sind fehlerfreie oder fehlerhafte Produktionsergebnisse, ärztliche Befunde einer bestimmten Krankheit oder Gesundheit und ob in einem Bild ein bestimmtes Objekt vorhanden ist oder nicht.

Bei der Multiclass Logistic Regression werden die Eingaben in mehr als zwei Klassen unterteilt. Es lassen sich damit beispielsweise in einem Bild verschiedene Arten von Fahrzeugen, Gegenständen oder Tieren unterteilen.

Diagramm: Ein-Schicht-NN (Die Ausgabeschicht wird nicht mitgezählt) Logistic Regression kann mit einer einzelnen oder mehreren neuronalen Schichten implementiert sein. Das künstliche neuronale Netz in Abb. x besitzt eine einzelne aus drei Einheiten bestehende Schicht. Eine Einheit besteht zunächst aus einer differenzierbaren Funktion. Oft ist das eine Multiplikation der Eingaben X mit einer Anzahl lernbaren Parametern W und eine anschließende Addition mit Bias-Werten b . Die Ausgabe einer Schicht wird oft mit \hat{y} bezeichnet:

$$\hat{y} = WX + b \quad (2.1)$$

Dieser folgt eine Aktivierungsfunktion, die dem errechneten Wert einen anderen Wert zuordnet, der entweder nahe 0 oder nahe 1 liegt (gelegentlich nahe -1 oder nahe 1). Beispiele für Aktivierungsfunktionen sind Sigmoid:

$$\text{sig}(t) = \frac{1}{1 + e^{-t}} \quad (2.2)$$

und Rectified Linear Unit, kurz ReLU:

$$f(x) = \max(0, x) \quad (2.3)$$

oder auch Leaky ReLU:

$$f(x) = \max(0.01x, x) \quad (2.4)$$

In einigen künstlichen neuronalen Netzen werden zusätzlich die Eingaben für jede Schicht normalisiert, wodurch die Minimierung der Verlustfunktion (Gradient Descent) optimiert werden kann.

2) Deep Neural Networks

3) Convolutional Neural Networks

Für Bilddaten wendet ein CNN diese Operation typischerweise auf zwei dreidimensionale Matrizen an, nämlich einerseits auf die Eingabedaten und andererseits auf einen sogenannten Filter oder auch Kernel. Die Eingabedaten sind in der ersten Schicht des Netzes die RGB-Pixelinformationen und in allen weiteren konvolutionalen Schichten die Ausgabe der vorherigen Schicht. Ein Filter ist eine Anzahl von trainierbaren Parametern, in diesem Fall auch eine . Beide Matrizen haben also die Form $(H \times W \times C)$, wobei H die Höhe, W die Breite und C die RGB-Farbwerte repräsentiert.

Jede der drei Matrixdimensionen variiert üblicherweise zwischen den verschiedenen Schichten des Netzes. In fast allen CNNs ([1], [8], [9], [4]) nimmt die Kardinalität zunächst ab. Die Reduktion kann durch die Konvolution selbst entstehen oder durch Pooling-Schichten. Beim Pooling werden aus benachbarten Matrixkoeffizienten meist das Maximum, seltener der Durchschnitt oder andere Aggregationen gebildet. Auf diese Weise wird das neuronale Netz darauf trainiert die relevanten Informationen zu extrahieren. [1]

Den konvolutionalen Schichten folgt in einigen Anwendungsfällen eine voll vernetzte Schicht (engl. Fully Connected Layer, FC), in der für jedes Neuron mit jedem Neuron der vorherigen Schicht eine Verbindung besteht. Besonders für die Bilderkennung ist diese Architektur gut geeignet. Die Ausgabe des neuronalen Netzes ist dann ein Vektor, beispielsweise von Wahrscheinlichkeitswerten für das Vorhandensein bestimmter Objekte und gegebenenfalls Bildkoordinaten der erkannten Objekte. Im Fall der Bildgenerierung ist die Ausgabe aber wieder eine Matrix von RGB-Pixelinformationen in der Form $(H \times W \times 3)$.

4) U-Net-Architektur

Die bis hierhin beschriebenen neuronalen Netze besitzen eine gradlinige Struktur, in der die Ausgabe einer Schicht nur an die nächste Schicht übergeben wird. Bei zunehmender Anzahl der Schichten verbessert sich die Performance neuronaler Netze mit diesem Aufbau zunächst, aber verschlechtert sich bei zu vielen Schichten wieder. In einem Residual Neural Network (ResNet) verhindern zusätzliche Verbindungen zwischen nicht direkt aufeinanderfolgenden Schichten diesen Performanceverlust.

Ein U-Net ist eine spezielle Form eines ResNets. Es hat eine annähernd symmetrische Struktur, in der sich die Kardinalitäten der Matrizen zuerst verringern und anschließend wieder erhöhen. U-Nets erzielen selbst mit wenigen Trainingsdaten gute Ergebnisse und benötigen dafür vergleichsweise wenig Rechenleistung. [10]

5) Generative Adversarial Networks

Ein Generative Adversarial Network (GAN) besteht zunächst aus einem Generator und einem Discriminator [5]. Der Generator lernt während des Trainings täuschend echt aussehende Bilddaten zu generieren. Der Discriminator wird dagegen darauf trainiert, echte von generierte

Bildern zu unterscheiden. Anschließend können beide Modelle “gegeneinander antreten”. Deswegen wird es Generative *Adversarial* Network genannt.

GANs unterscheiden sich von anderen Modellen durch ihren Aufbau und in Bezug auf das Trainingsziel. Künstliche neuronale Netze ermitteln häufig einen skalaren Wert wie beispielsweise einen Wahrscheinlichkeitswert und minimieren zu diesem Zweck eine Verlustfunktion. In einem GAN sind zwei CNNs im Einsatz. Das erste, der Generator, erstellt Tensoren n-ter Klasse. In diesem Beispiel sind das RGB-Bildinformationen. Das zweite CNN wird Discriminator genannt und bekommt als Eingabe die Ausgabe des ersten CNNs. Der Discriminator wird darauf trainiert, generierte Bilder von Bildern aus dem Trainingsset zu unterscheiden. Er minimiert also eine Verlustfunktion. Der Generator wird darauf trainiert, diese Verlustfunktion zu maximieren. Dieser Vorgang stammt aus der Spieltheorie und heißt Minimax. [5]

Eine Erweiterung des GANs ist das Conditional Generative Adversarial Net (cGAN). In cGANs erhält der Generator zusätzliche Eingabedaten (y , “Ground Truth”), die Hinweise für die Generierung enthalten. Der Discriminator erhält diese zusätzlichen Daten ebenfalls, um die Erkennung während des Trainings zu optimieren. [6]

6) Image-To-Image-Translation

Der Image-To-Image-Translation-Algorithmus oder kurz Pix2Pix-Algorithmus verwendet ein GAN, um Bilder in Bilder zu übersetzen. Dafür kommen zwei CNNs zum Einsatz, nämlich eins für den Generator und eins für den Discriminator.

3 Implementierung

1) Entwicklungsumgebung

1).1 Ubuntu Linux

1).2 Python

1).3 Tensorflow

1).4 CUDA

4 Experimente und Resultate

1) Vorbereitung der Eingabedaten

Die Trainingsdaten liegen im NDJSON-Format und als Blender-Dateien beziehungsweise Wavefront OBJ-Dateien vor. Die effizienteste Möglichkeit, die Skizzen und 3D-Modelle für die Verarbeitung in einem CNN vorzubereiten, ist das Rendern und Speichern als Bilddateien. Als Dateiformat kommen JPEG oder PNG infrage. Beide Formate können leicht als Trainingsset mit Tensorflow geladen werden.

Bei der Verarbeitung in einem Convolutional Neural Network spielt die Bildgröße in Bezug auf die Verarbeitungszeit eine wichtige Rolle. Bildformate der Größe 256x256 oder kleiner sind üblich und gut geeignet. Für ein GAN ist es zwar nicht erforderlich, aber sinnvoll, für Ein- und Ausgabedaten dieselben Dimensionen festzulegen. Bei der Bildgenerierung aus Skizzen unterscheiden sich Ein- und Ausgabedaten natürlich bei der Farbtiefe. Während die Skizzen Graustufenbilder sind, besitzen die generierten Bilder drei Farbkanäle (RGB).

Sowohl während der Entwicklung als auch zur Laufzeit kann es vorteilhaft sein, Farbwerte statt auf der oft verwendeten Skala von 0 bis 255 als Fließkommazahlen im Bereich 0,0 bis 1,0 darzustellen. Auch für die Ausgaben der einzelnen Schichten eines künstlichen neuronalen Netzes kann diese Normalisierung durchgeführt werden (TODO: BatchNorm).

Ein ebenso wichtiger wie aufwendiger Vorgang ist die Klassifizierung der Trainingsdaten, also die Zuweisung von Eingaben zu den erlernbaren Ergebnissen. Aufgrund der selbst erstellten Ausgabebilder existieren für diesen Zweck keine vorgefertigten Datasets. Die Sortierung und Zuweisung erfolgt deshalb manuell.

2) Anwendung herkömmlicher Shader

3) Hyperparameter

Die Pix2Pix-Referenzimplementierung ist bereits für die Übersetzung von Skizzen in Fotos eingestellt. Für das Training waren anfangs mehrere Tausend Epochs, also Trainingsdurchläufe, erforderlich, um zufriedenstellende Ergebnisse zu sehen. TODO: Diese Zahl konnte durch ... verringert werden.

Die Eingabebilder sind 256x256 Pixel groß und besitzen einen Farbkanal für Graustufen. Sie werden am Anfang des Trainingsprozesses durch sogenanntes Jittering augmentiert. Dabei werden die Bilder zuerst auf 286x268 Pixel vergrößert und anschließend auf einen zufälligen 256x256 Pixel großer Ausschnitt wieder verkleinert. Diese Pixelgrößen können im Experiment geändert werden.

Der Adam Optimierer [11] erhält für die Learning-Rate den Wert 0,0002. Das Momentum ist auf 0,9 voreingestellt. Diese beiden Werte beeinflussen die Lerngeschwindigkeit und sind in begrenztem Maße anpassbar.

4) Performancebeobachtungen

5 Diskussion

Literatur

1. GOODFELLOW, I. u. a.: *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
2. ZHANG, A. u. a.: *Dive into Deep Learning*. 2020. <https://d2l.ai>.
3. CHOLLET, F.: *Deep Learning with Python*. 1st. USA: Manning Publications Co., 2017. ISBN 1617294438.
4. ISOLA, P. u. a.: *Image-to-Image Translation with Conditional Adversarial Networks*. 2018. Abgerufen unter arXiv: 1611.07004 [cs . CV].
5. GOODFELLOW, I. J. u. a.: *Generative Adversarial Networks*. 2014. Abgerufen unter arXiv: 1406.2661 [stat . ML].
6. MIRZA, M.; OSINDERO, S.: *Conditional Generative Adversarial Nets*. 2014. Abgerufen unter arXiv: 1411.1784 [cs . LG].
7. RADFORD, A. u. a.: *Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks*. 2016. Abgerufen unter arXiv: 1511.06434 [cs . LG].
8. LECUN, Y. u. a.: Object Recognition with Gradient-Based Learning. In: *Contour and Grouping in Computer Vision*. Springer, 1999.
9. RONNEBERGER, O. u. a.: U-Net: Convolutional Networks for Biomedical Image Segmentation. In: *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Springer, 2015, Bd. 9351, S. 234–241. LNCS. Auch verfügbar unter: <http://lmb.informatik.uni-freiburg.de/Publications/2015/RFB15a>. (available on arXiv:1505.04597 [cs.CV]).
10. HE, K. u. a.: *Deep Residual Learning for Image Recognition*. 2015. Abgerufen unter arXiv: 1512.03385 [cs . CV].
11. KINGMA, D. P.; BA, J.: *Adam: A Method for Stochastic Optimization*. 2017. Abgerufen unter arXiv: 1412.6980 [cs . LG].

Bildnachweis

Abkürzungs- und Symbolverzeichnis

Abkürzungen

AC	Air Compressor, Luftverdichter
APH	Air Preheater, Luftvorwärmer
CC	Combustion Chamber, Brennkammer
EXP	Expander
HRSG	Heat Recovery Steam Generator, Abhitzekessel

Lateinische Symbole

c	Spezifische Kosten je Exergieeinheit, €/J _{ex}
\dot{C}	Kostenstrom, €/h
CC	Kapitalgebundene Kosten, €
cf	Capacity Factor, Jährliche Auslastung, –
e	Spezifische Exergie, J/kg
\bar{e}	Spezifisch molare Exergie, J/mol
\dot{E}	Exergiestrom, W
f	Exergoökonomischer Faktor, –
fc	Spezifische Brennstoffkosten, €/J
FC	Brennstoffkosten, €
h	Spezifische Enthalpie, J/kg
\dot{H}	Enthalpiestrom, W
HHV	Brennwert, J/kg

Griechische Symbole

Δ	Differenz
ε	Exergetischer Wirkungsgrad, –
η_s	Isentroper Wirkungsgrad, –
κ	Isentropenexponenten, –
λ	Luftzahl, –

Hoch- und tiefgestellte Indizes

0	Referenzpunkt, Thermodynamische Umgebung
a	Average, Mittlere
D	Destruction, Vernichtung
F	Fuel, Brennstoff, Aufwand
net	Netto

Abbildungsverzeichnis

Tabellenverzeichnis

Anhang