

TECHNISCHE UNIVERSITÄT DORTMUND

Fakultät Informatik

Design Automation for Embedded Systems Group

Training eines binären Neuronalen-Netzwerkes auf den MNIST Datensatz

FACHPROJEKT

Jack Diep, Florian Köhler, Yannick Naumann

Betreuung:

M.Sc. Mikail Yayla

1. September 2021

Inhaltsverzeichnis

Abbildungsverzeichnis	ii
1 Einleitung	1
1.1 Motivation	1
2 Neuronale Netzwerke	3
2.1 Neuronen	4
2.2 Schichten und Kanten	4
2.3 Aktivierungsfunktionen	5
2.4 Training	7
2.4.1 Ablauf	7
2.4.2 Backpropagation	8
3 Das BNN	9
3.1 Aktivierungsfunktion	9
3.2 Binärer Linear-Layer	10
3.3 BatchNorm	10
3.4 Auswertung der letzten Schicht	11
3.5 Verluste durch binäre Linear-Layer	12
4 Training des BNNs	13
4.1 Lernrate	13
5 Binarisierung	15
6 Export	16
6.1 Export der Kantengewichte	16
6.2 Export der Schwellwerte	17
Literatur	18

Abbildungsverzeichnis

2.1	Struktur eines neuronalen Netzwerks	3
2.2	Berechnung eines Neuronenwertes mit der <i>tanH</i> Aktivierungsfunktion . .	5
2.3	ReLU	6
2.4	Sigmoid	6
2.5	Softmax	6
2.6	tanH	6
3.1	Die Vorzeichenfunktion	9
3.2	Training mit und ohne BatchNorm	11
3.3	Verlustmessung für den <i>Linear-Layer</i>	12
4.1	Training mit verschiedenen Lernraten	14

Kapitel 1

Einleitung

Neuronale Netzwerke bilden eine Unterkategorie des Machine-Learning und erlauben Auswertungen von Eingaben auf Basis von zuvor angelernten, empirischen Ergebnissen. Ein neuronales Netzwerk bildet ein System aus Neuronen ab, welche Schichtweise verbunden sind einen unidirektionalen Datenfluss erzeugen. Dieses System besteht üblicherweise aus einer Eingabeschicht, einer Ausgabeschicht und dazwischen beliebig viele *versteckte* Schichten, welche die eigentliche Arbeit des Netzwerks verrichten. Die Anzahl der Neuronen in den Eingabe- und Ausgabeschichten ist intuitiv wählbar. Die Größe der Eingabeschicht wird häufig durch die Anzahl der möglichen Eingaben bestimmt, die Größe der Ausgabeschicht durch die Anzahl der möglichen Ergebnisse. Die Größe und Anzahl der dazwischen liegenden Schichten hingegen muss je nach Anforderung und Gegebenheiten individuell ermittelt werden. Je größer das Netzwerk desto höher sind die Anforderungen an die benötigte Hardware um dieses zu betreiben. Bei schwächerer Hardware oder Einschränkungen bezüglich der Energieversorgung können kleinere Netzwerke eingesetzt werden, wenn auch häufig mit geringerer Genauigkeit verglichen mit einem größeren Netzwerk.

1.1 Motivation

Die Größe eines Netzwerks kann beim Entwurf dessen direkt beeinflusst werden. In Anwendungsgebieten, bei denen der Fokus auf geringen Hardwareanforderungen liegt, stoßen schnell das Problem der schwindenden Genauigkeit. Auf IoT-Geräten oder mobilen Plattformen finden klassische neuronale Netzwerke daher nur eingeschränkt Nutzen.

2016 veröffentlichten M. Courbariaux und Y. Bengio [CB16] eine wissenschaftliche Arbeit und stellen dort das *binarisierte neuronale Netzwerk* (kurz BNN) vor. Dieses könne im Vergleich zu einem klassischen neuronalen Netzwerk eine theoretische Geschwindigkeitssteigerung auf das 32-fache erreichen. Die Genauigkeit des BNN liege jedoch nur knapp unter derer klassischer Netzwerke. Das BNN ermöglicht dank dieser Eigenschaften den Einsatz von neuronalen Netzen auf vergleichsweise schwacher Hardware und verspricht zugleich nur geringe Genauigkeitseinbuße.

In dieser Ausarbeitung wird die allgemeine Funktionsweise von neuronalen Netzwerken erläutert und anschließend der Entwurf eines binarisierten Netzwerk mit Fokus auf Handschriftenerkennung auf Basis des MNIST Datensatzes dokumentiert. Dabei werden grundlegende Überlegungen, das Vorgehen, Herausforderungen sowie dazu erarbeitete Lösungen vorgestellt.

Kapitel 2

Neuronale Netzwerke

Unter einem neuronalen Netzwerk versteht man ein System aus Neuronen. Diese sind schichtweise organisiert wobei jedes Neuron einer Schicht jeweils zu allen Neuronen der direkt anliegenden Schichten verbunden ist. Abbildung 2.1 zeigt beispielhaft ein neuronales Netzwerk aus 19 Neuronen mit insgesamt 5 Schichten.

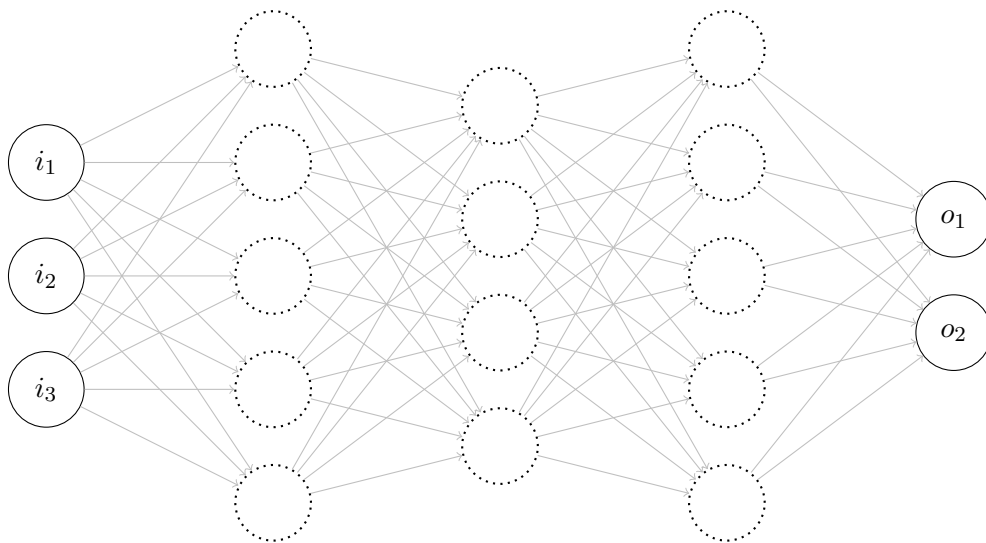


Abbildung 2.1: Struktur eines neuronalen Netzwerks

2.1 Neuronen

Unter einem Neuron versteht sich bei neuronalen Netzen lediglich ein Knoten, in dem üblicherweise ein 32-bit großer Wert hinterlegt ist. Häufig kommen hier Gleitkommazahlen zwischen -1.0 und 1.0 zum Einsatz, da sich dieser Wertebereich besonders gut zum Rechnen eignet und Eigenschaften bezüglich der Multiplikation besitzt, die eine Wertexplosion verhindern. Die Werte aller Neuronen, mit Ausnahme derer in der Eingabeschicht, setzen sich jeweils aus den Werten aller Neuronen der direkt davor liegenden Schicht zusammen. Das genaue Vorgehen bei der Wertermittlung hängt jeweils vom Netzwerk und den dort verwendeten Aktivierungsfunktionen ab.

2.2 Schichten und Kanten

In einem neuronalen Netzwerk ist jedes Neuron einer Schicht mit allen Neuronen der jeweiligen davor liegenden und danach liegenden Schicht über Kanten verbunden. Neuronale Netzwerke sind unidirektionale Graphen, demnach fließen Informationen über Schichten (und folglich Kanten) nur in eine Richtung.

Allen Kanten wird initial eine Gewichtung zugewiesen, welche erneut netzwerkabhängig generiert werden oder durch zuvor angelernte Daten bestimmt werden. Beim Trainieren des Netzwerks werden diese bei jeder Lerniteration (auch *Epoch* genannt) justiert, während sie beim Betrieb für gewöhnlich keine Änderungen mehr erfahren. Die Genauigkeit eines Netzwerks wird überwiegend durch diese Gewichte bestimmt, daher ist das Ziel beim Trainieren eines Netzes die Optimierung jener.

Kantengewichte wirken sich maßgeblich auf die Wertberechnung von Neuronen aus. Diese wird in zwei Schritten ausgeführt. Im ersten Schritt wird aus den Kantengewichten aller eingehenden Kanten und den Werten der darüber verbundenen Neuronen die Produktsumme gebildet. Für den zweiten Schritt sind jeweils Aktivierungsfunktionen notwendig, welche im nachfolgenden Kapitel erläutert werden.

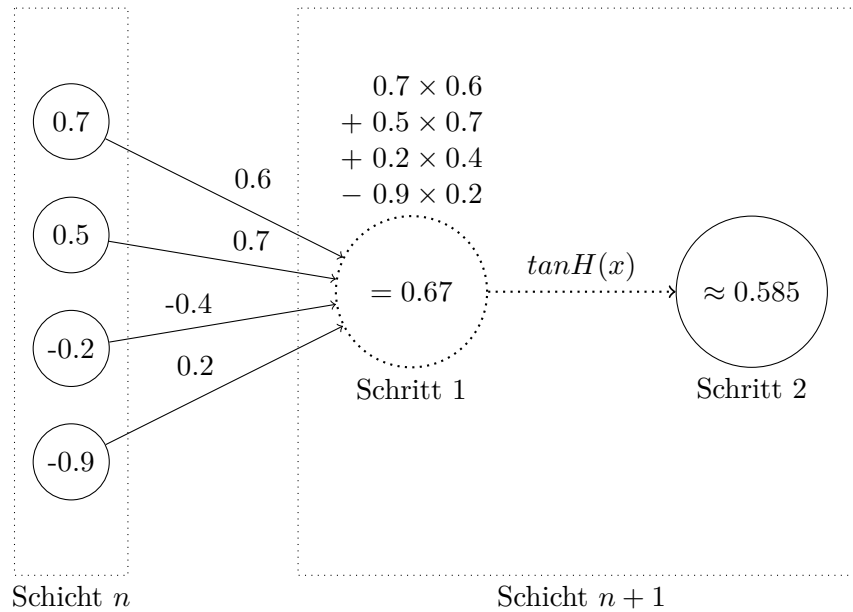


Abbildung 2.2: Berechnung eines Neuronenwertes mit der $\tan H$ Aktivierungsfunktion

2.3 Aktivierungsfunktionen

Im zweiten Schritt wird der zuvor errechnete Wert durch eine weitere Funktion modifiziert. Sogenannte *Aktivierungsfunktionen* können beliebig gewählt werden, müssen jedoch offensichtlich alle möglichen Eingaben auf einen Wert abbilden können. Die verwendete Aktivierungsfunktion kann je nach Schicht variieren. Einmal gewählt, ist diese jedoch für die jeweilige Schicht im Netzwerk für die Laufzeit fest.

Die beiden erläuterten Berechnungsschritte werden schichtweise in Richtung des Datenflusses des Netzwerks für alle Neuronen durchgeführt. Im folgenden Beispiel wird bildhaft dargestellt, wie solch eine Neuronenwertberechnung für ein Netzwerk aussieht, welche in der betrachteten Schicht die $\tan H$ -Aktivierungsfunktion verwendet.

Die Wahl der Aktivierungsfunktion beeinflusst die möglichen Werte, die Neuronen innerhalb einer Schicht annehmen können. Um eine Wertexplosion zu vermeiden (z.B. bei Verwendung von $\text{ReLU}(x) : \max(0, x)$ als Aktivierungsfunktion), können zusätzliche Normalisierungsschichten verwendet werden, welche die Neuronenwerte einer Schicht auf einen erwünschten Zielbereich einschränken. Es existieren jedoch auch Aktivierungs-

funktionen, die diese *Squashing*-Eigenschaft direkt besitzen. Darunter zählt auch die im vorherigen Beispiel verwendete Funktion \tanh , bei der alle Eingaben auf den Zielbereich $[-1,1]$ abgebildet werden.

Im Folgenden sind einige, für neuronale Netzwerke übliche Aktivierungsfunktionen abgebildet.

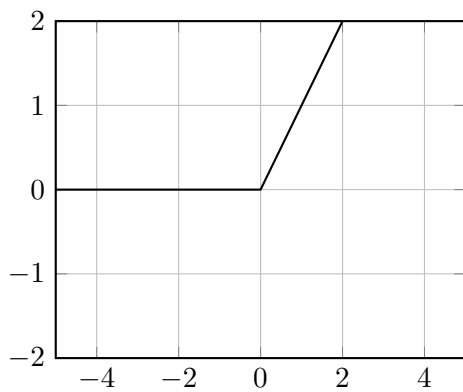


Abbildung 2.3: ReLU

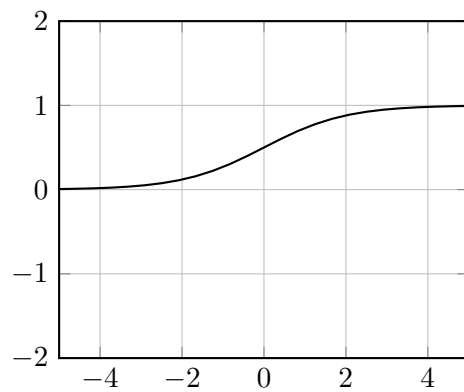


Abbildung 2.4: Sigmoid

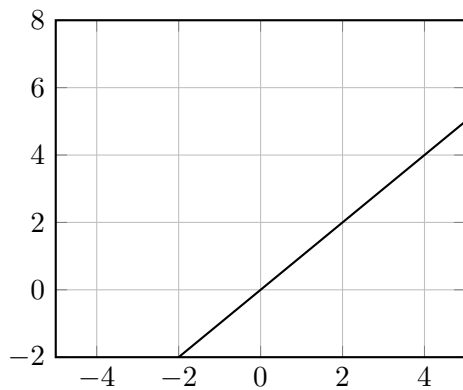


Abbildung 2.5: Softmax

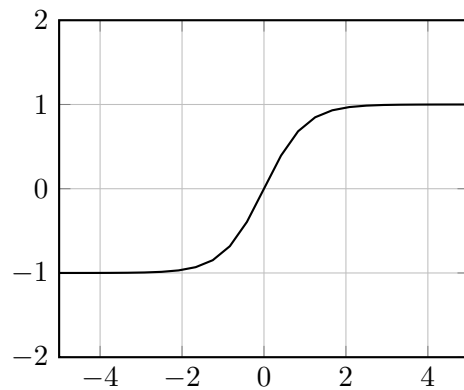


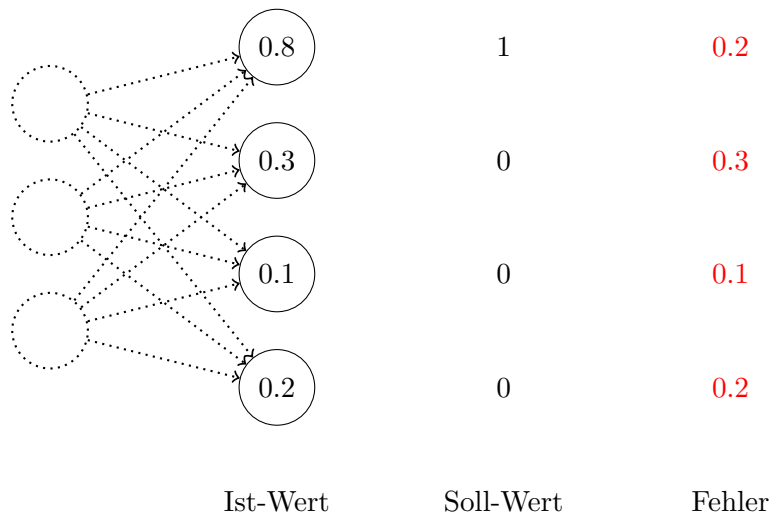
Abbildung 2.6: tanH

2.4 Training

Die Stärken und Schwächen eines neuronalen Netzwerks liegen neben der Struktur der Neuronen, der Schichten und den gewählten Aktivierungsfunktionen vor allem in den verwendeten Kantengewichten. Diese zu optimieren kann ein langwieriger und rechenintensiver Prozess sein. Kleine Änderungen im Aufbau des Netzwerks können große Teile von zuvor erlernten Gewichten unbrauchbar machen, daher werden Netzwerke häufig an sich selbst trainiert. Das bedeutet, dass beim Trainieren eines neuen Netzwerks die initialen Kantengewichte keine empirischen Daten aus anderen (ähnlichen) Netzwerk verwenden, sondern diese zu Beginn pseudozufällig gewählt werden.

2.4.1 Ablauf

Beim Training werden wiederholt verschiedene Eingaben in das Netzwerk getätigt, zu denen das korrekte Ergebnis bekannt ist. In der letzten Schicht des Netzwerks, der Ausgabeschicht, werden die vom Netzwerk bestimmten Werte der jeweiligen Neuronen mit den erwarteten Werten verglichen. Die Differenz dieser wird als Fehler bezeichnet. Das Ziel des Trainings ist es, den Fehler durch Anpassung der Kantengewichte zu verringern oder im Optimalfall ganz zu eliminieren (Fehler = 0). In folgender Abbildung ist ein Fehlerbeispiel dargestellt.



2.4.2 Backpropagation

Es stellt sich nun die Frage, ob und wie errechnet werden kann, welches Kantengewicht wie angepasst werden muss, um den Fehler zu verringern. Der Begriff der *Backpropagation* beschreibt das Durchlaufen des Netzwerks, jedoch entgegen der eigentlichen Laufrichtung und das Rückschließen der notwendigen Gewichts Anpassung, um die Ausgabe in die gewünschte Richtung zu verändern.

Bevor der genaue Wert für eine einzige Anpassung ermittelt werden kann, muss die Entwicklungsrichtung des Fehlers bei Änderung des Kantengewichts errechnet werden.

Kapitel 3

Das BNN

Bei unserem Neuronalen Netz handelt es sich um ein binäres Netzwerk. Im folgenden werden die Anpassungen, beziehungsweise die verwendeten Komponenten beschrieben.

3.1 Aktivierungsfunktion

In Neuronalen-Netzwerken kann, wie in Kapitel 2.3 beschrieben, eine ganze Reihe von Aktivierungsfunktionen verwendet werden. Jede dieser Funktionen hat, je nach Anwendungskontext, verschiedene Vor- und Nachteile.

Für BNNs ist hier die $\text{sign}(x)$ Funktion eine populäre Wahl. Diese überführt die Akti-

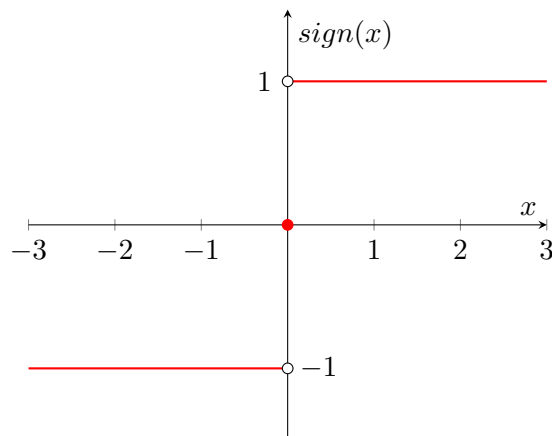


Abbildung 3.1: Die Vorzeichenfunktion

vierungen direkt in binäre Werte. Da für die *backpropagation* allerdings die Ableitung

der Aktivierungsfunktion benötigt wird und die Vorzeichenfunktion nicht stetig differenzierbar ist, wird diese oft approximiert. In unserem Netzwerk wird sie, wie häufig verwendet, durch die *hard tanH* Funktion approximiert [Kim+20]. Diese ist günstiger zu berechnen als die Funktion *tanH* und ist bei Werten $x \leq -1 \vee x \geq 1$ gleich zur Vorzeichenfunktion.

3.2 Binärer Linear-Layer

3.3 BatchNorm

Die Rolle des *BatchNorm-Layers* in nicht-binären Neuronalen-Netzwerken ist eine leicht andere als in unserem BNN. Durch *BatchNorm* werden die Gewichte einer Schicht normalisiert. Hier wird dafür gesorgt, dass *Batches* einen Mittelwert von Null und eine Standardabweichung von Eins haben[IS15] .

Durch diese reduzierte Streuung der Werte, können in Neuronalen Netzen höhere Lernraten verwendet werden.

Bei binären Netzwerken hat dies, aufgrund der diskreten Kantengewichte, kaum einen Effekt und hat somit kaum Auswirkungen auf die verwendbare Lernrate. In BNNs erfüllt die *BatchNorm* Schicht eine zentrale Rolle für das Lernen des Netzwerkes. Primär dient die Normalisierung der Schicht dazu, das *expoding-gradient* Problem zu verhindern. Hier wird beim Training des Netzwerkes, welches die Error-Werte reduzieren soll, ein sehr hoher Error-Wert akkumuliert. Dies führt zu einer zu starken Anpassung der Gewichte, welche folgend zu einer niedrigeren Genauigkeit führt[SBN19].

Da durch wiederholtes auftreten die Genauigkeit stark abgesenkt wird, kann dieses Netzwerk sich, ab einer gewissen Schwelle, nicht mehr verbessern.

Der Versuch in Abbildung 3.2 zeigt einen Vergleich des Netzwerkes mit und ohne *BatchNorm*. Hierfür wurde ein Netzwerk für 50 Epochen trainiert. nach jeder Epoche wurde die Genauigkeit getestet. Anschließend wurde weiter Trainiert. Um eine bessere Vergleichbarkeit zu erzielen, wurden die Bilder über die Schwellwert-Methode binarisiert.

Wie in Abbildung 3.2 zu sehen ist, performt das Netzwerk mit *BatchNorm* Schicht zu jeder Trainingsdauer besser als ohne *BatchNorm* Schicht. Während in den ersten 20 Epochen mit *BatchNorm* ein klarer Aufwärtstrend zu erkennen ist, bleibt die Genauigkeit

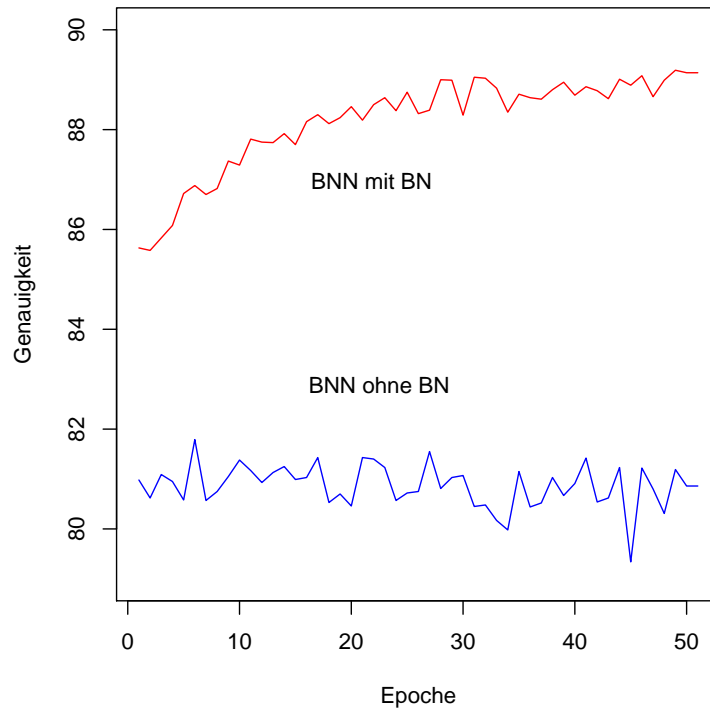


Abbildung 3.2: Training mit und ohne BatchNorm

ohne *BatchNorm* immer zwischen 79% und 82%. Hier sind besonders gut die starken Einbrüche nach Unten zu erkennen, die mit *BN* deutlich schwächer sind. Diese zeigen besonders deutlich den starken Genauigkeitsverlust, der durch *exploding gradients* verursacht wird. Die Konsequenz ist, dass eine längere Trainingsdauer keine Verbesserung des Netzwerkes mehr bedeutet, was bedeutet, dass die Konvergenz-Schwelle bei < 82 erreicht ist. Ab dieser Schwelle verbessert sich das Netzwerk ohne *BN* nicht mehr.

3.4 Auswertung der letzten Schicht

Um am Ende das Ergebnis des Netzwerkes auszuwerten, müssen die Ergebnisse des letzten *Linear-Layers* normalisiert werden. Diese Normalisierung entscheidet, ob ein Neuron feuert oder nicht. Für die Normalisierung der Daten wird hier die $\text{logsoftmax}(x)$ Funktion verwendet. Durch $\text{softmax}(x)$ werden die Aktivierungen der letzten Schicht so norma-

lisiert, dass ihre Summe eins ergibt. Dies ist mit Wahrscheinlichkeiten zu vergleichen, mit der das Bild die jeweilig zugeordnete Zahl widerspiegelt. Um bessere Ergebnisse in Kombination mit der Verlustfunktion, *negative log likelihood loss*, zu erzielen, wird die *logsoftmax*-Funktion verwendet.

Da bei der Anwendung des MNIST-Datensatzes immer nur das Neuron ausgewählt werden sollte, da immer nur eine Zahl auf einem Bild abgebildet ist, wird bei der Auswertung über die *argmax*-Funktion das aktivste Neuron ausgewählt. Dieses entspricht dann der Zahl, die am wahrscheinlichsten abgebildet ist.

3.5 Verluste durch binäre Linear-Layer

Durch die binarisierung der *Linear-Layer* ist zu vermuten, dass diese, im Vergleich zu normalen *Linear-Layer*, etwas schlechter performen. Dies ist der Fall, da die Anzahl der möglichen Kantengewichte stark, auf Null und Eins, eingeschränkt ist.

Durchgang	binär	normal
1	88.29	97.43
2	87.32	96.98
3	87.19	97.2

Abbildung 3.3: Verlustmessung für den *Linear-Layer*

Trainiert wurde hier das gleiche Netzwerk, ein mal mit binären *Linear-Layern*, das andere mal mit normalen *Linear-Layer*. Jedes Netzwerk wurde für 50 Epochen trainiert, bevor die Genauigkeit ausgewertet wurde. Um sicher zu gehen, ob die Genauigkeit gegen diesen Wert konvergiert, wurde jede Messung drei mal wiederholt.

Wie in Abbildung 3.3 zu sehen, leidet die Genauigkeit des Netzwerkes beachtlich unter der Binarisierung der *Linear-Layer*. Der Mittelwert für das Training mit normalem *Linear-Layer* ist hierbei 97,2%, während bei binären Schichten ein Durchschnitt von 87,6% erreicht wird. Es ist klar zu sehen, dass die Einschränkung der Gewichte auf Null und Eins und der damit einhergehende Granularitätsverlust, sich stark auf die Genauigkeit des Netzwerkes, bei gleicher Größe, auswirken.

Kapitel 4

Training des BNNs

Wie bereits in Kapitel 2.4 beschrieben, werden Netzwerke über eine rechenintensive Annäherung der Kantengewichte an das optimale Ergebnis trainiert. Diese Konvergenz in Richtung eines besseren Ergebnisses kann hierbei, selbst bei schlechtem Training, meist erreicht werden, da selbst marginale Verbesserungen eine Auswirkung haben. Da *BNNs* jedoch die stärkste Form von quantisierten Netzwerken sind, können hier keine stetigen Änderungen an den Kantengewichten vorgenommen werden.

4.1 Lernrate

Die Lernrate ist ein zentraler Parameter beim Training von Neuronale Netzwerken. Sie gibt die Rate an, mit der Kantengewichte in einem Durchlauf angepasst werden. Je höher die Lernrate also ist, desto schneller passt sich das Modell an gerade trainierte Daten an. Bei der Wahl der Rate sind die Probleme des *overfitting* und *underfitting* zu beachten. Als *underfitting* bezeichnet man eine zu geringe Anpassungsfähigkeit des Netzwerkes, ausgelöst über zu kleine Netze oder zu niedrige Lernrate. Ist die Lernrate zu klein, kann der Error nicht reduziert werden, die Genauigkeit nimmt nur, wenn überhaupt, sehr langsam zu. Das Netz kann aus Daten keine Lernerfolge ziehen.

Overfitting ist eine zu schnelle Anpassung des Modells. Ist zum Beispiel die Lernrate zu hoch, kann die optimale Genauigkeit des Netzwerkes nicht erreicht werden, da die Anpassungen pro Trainingsdatum zu groß sind. Hier wird das Optimum immer, durch eine Überkorrektur, übersprungen und kann nie erreicht werden[Smi18].

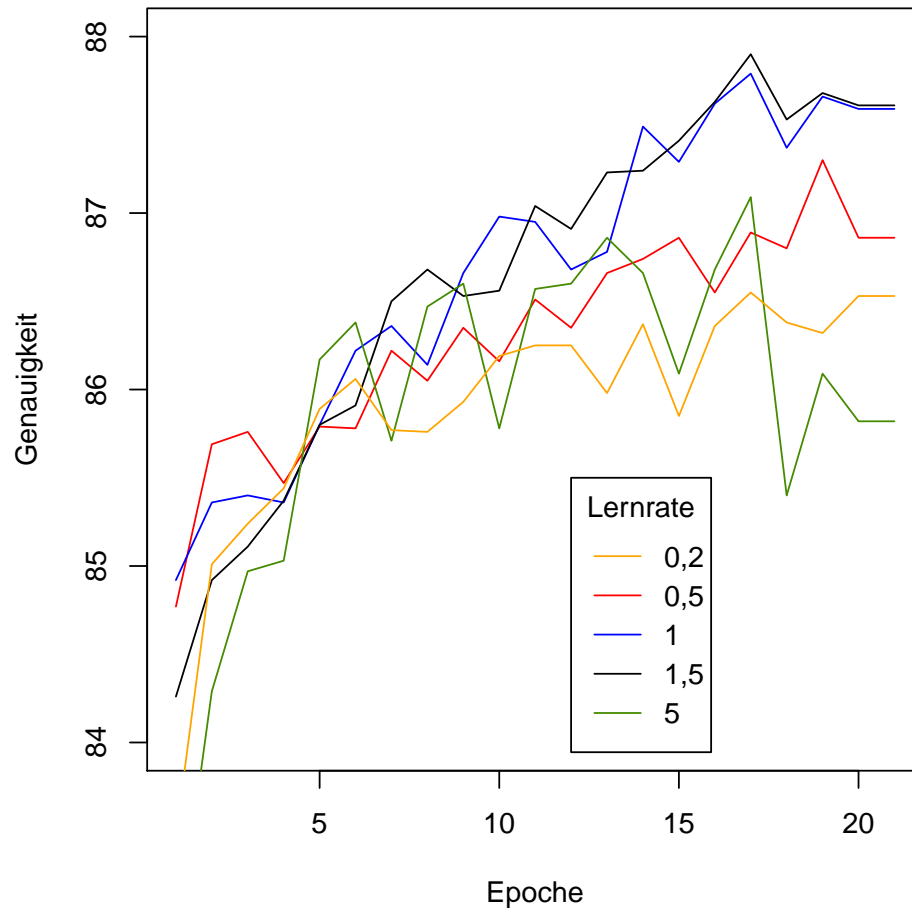


Abbildung 4.1: Training mit verschiedenen Lernraten

Zur Ermittlung der optimalen Lernrate haben wir unser Netzwerk für 20 Epochen mit verschiedenen Lernraten trainieren lassen. In Abbildung 4.1 sind die Ergebnisse mit einer Messung der Genauigkeit nach jeder Epoche dargestellt. Bei einer Lernrate von 5 ist deutlich das *overfitting* ab Generation 17 zu erkennen. Hier werden zu starke Anpassungen vorgenommen, weshalb sich die Genauigkeit vom Optimum entfernt. Lernraten unter Eins hingegen nähern sich dem Optimum erst gar nicht genug an. Die Lernraten 1 und 1,5 liefern hier die besten Ergebnisse. Da das Training mit einer Rate von 1,5 weniger große Ausschläge zeigt, wurde sich für diese Lernrate entschieden.

Kapitel 5

Binarisierung

Kapitel 6

Export

Nachdem das Netzwerk nun trainiert wurde, müssen die Ergebnisse, die Kantengewichte und Schwellwerte der Neuronen, nun exportiert werden. Im Folgenden sollen diese dann in den BNN-Beschleuniger Baustein importiert und verwendet werden. Da der Import in VHDL stattfindet, eignen sich hier simple Formate, sprich eine einfache Textdatei. Diese kann dann, Zeichen nach Zeichen, von dem Import-Buffer eingelesen und in einer Matrix gespeichert werden.

6.1 Export der Kantengewichte

Da es sich bei unserem Netzwerk um ein *FullyConnected Neural Network* handelt, ist insbesondere jedes Neuron mit jedem Neuron der Folgenden Schicht verbunden. Bei unserem BNN ergibt sich also folgende Kantenanzahl

$$784 \cdot 500 + 500 \cdot 1024 + 1024 \cdot 1024 = 1.952.576$$

Diese Gewichte müssen alle, mit möglichst wenig Mehrkosten, in die Datei geschrieben werden. Da es sich bei den Gewichten lediglich um binäre Werte, Einsen und Nullen, handelt, ist kein Trennzeichen zwischen den Gewichten notwendig. Die Gewichte sind außerdem, trivialer Weise, präfixfrei und können fortlaufend in die Datei geschrieben werden.

Um die Gewichte zu extrahieren, wird zuerst über jeden *Layer* iteriert. In jedem *Layer* wird nun jedes Neuron abgelaufen. Jedes dieser Neuronen hat nun jeweils eine Kante

zu jedem Neuron in der nachfolgenden Schicht. Hier wird ebenfalls über alle Kanten iteriert und das jeweilige Gewicht wird hinten an eine Variabel an gehangen. Ist nun ein *Layer* fertig, wird der Inhalt der Variable, welche als Zwischenspeicher dient, in die Datei *weights.txt* geschrieben. So wird für jede Schicht ein IO-Zugriff gemacht.

6.2 Export der Schwellwerte

Literatur

- [CB16] Matthieu Courbariaux und Yoshua Bengio. “BinaryNet: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1”. In: *CoRR* abs/1602.02830 (2016). arXiv: 1602.02830. URL: <http://arxiv.org/abs/1602.02830>.
- [IS15] Sergey Ioffe und Christian Szegedy. “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”. In: *Proceedings of the 32nd International Conference on Machine Learning*. Hrsg. von Francis Bach und David Blei. Bd. 37. Proceedings of Machine Learning Research. Lille, France: PMLR, 2015, S. 448–456. URL: <https://proceedings.mlr.press/v37/ioffe15.html>.
- [Kim+20] Hyungjun Kim u. a. “Improving Accuracy of Binary Neural Networks using Unbalanced Activation Distribution”. In: *CoRR* abs/2012.00938 (2020). arXiv: 2012.00938. URL: <https://arxiv.org/abs/2012.00938>.
- [SBN19] Eyyüb Sari, Mouloud Belbahri und Vahid Partovi Nia. “A Study on Binary Neural Networks Initialization”. In: *CoRR* abs/1909.09139 (2019). arXiv: 1909.09139. URL: <http://arxiv.org/abs/1909.09139>.
- [Smi18] Leslie N. Smith. “A disciplined approach to neural network hyper-parameters: Part 1 - learning rate, batch size, momentum, and weight decay”. In: *CoRR* abs/1803.09820 (2018). arXiv: 1803.09820. URL: <http://arxiv.org/abs/1803.09820>.