

Bachelorarbeit

Titel der Bachelorarbeit: Image Analysis for Food Safety and Health

Name des Autors: TEMKENG Thibaut

Datum der Abgabe: Ende Oktober

Betreuung: Name der Betreuer/ des Betreuers: Shou Liu

Fakultät für Embedded Intelligence for Health Care and Wellbeing

Erklärung

Ich versichere wahrheitsgemäß, die Arbeit selbstständig verfasst, alle benutzten Hilfsmittel vollständig und genau angegeben und alles kenntlich gemacht zu haben, was aus Arbeiten anderer unverändert oder mit Abänderungen entnommen wurde, sowie die Satzung der Universität Augsburg zur Sicherung guter wissenschaftlicher Praxis in der jeweils gültigen Fassung beachtet zu haben.

Ort, den Datum

Inhaltsverzeichnis

1	Einleitung	6
2	Grundlagen	7
2.1	Neuron	7
2.2	Merkmalskarten(Feature-Maps)	7
2.3	Filters	8
2.4	Entwicklung von Künstlichen Neuronalen Netzen	8
3	Convolutional Neural Network	8
3.1	Feedforward	8
3.1.1	Input Layer	8
3.1.2	Faltungsschicht	8
3.1.3	Aktivierungsfunktion	11
3.1.4	Pooling Layer	14
3.1.5	Multi-layer Perzeptron (Fully Connected Layer)	15
3.2	Backforward	15
3.2.1	Fehlerfunktion	15
3.2.2	Gradient	17
3.2.3	Lernrate	17
3.2.4	Gradientenabstiegsverfahren	17
4	Kompression von Tiefe neuronale Netze (DNN)	20
4.1	Pruning Network	21
4.2	Quantisierung von neuronalen Netzwerken	25
4.3	Huffman Codierung	27
5	Experiment	28
5.1	Analyse der Ergebnisse mit Hilfe von Metriken	28
5.1.1	Metriken nennen und vergleichen	28
5.1.2	Klassifizierungsgenauigkeit	28
5.2	Einfluss der Lernrate	28
5.3	Convolution Layers	28
5.4	Besondere CNNs: Bitte nicht lesen	32
5.4.1	AlexNet	32
5.4.2	SqueezeNet	32
5.4.3	Xception	34
5.4.4	MobileNet	35
5.5	Algorithmen zur Optimierung des Gradientenabstiegsverfahren: Optimizer	35
5.5.1	Adaptive Gradient Algorithm (AdaGrad)	35
5.5.2	Root Mean Square Propagation(RMSProp)	36
5.5.3	Adaptive Moment Estimation(Adam)	37

5.6	Problem beim Training von Convolutional neuronale Netzwerke	38
5.6.1	Overfitting	38
5.6.2	Dataset	43
6	Literatur	43

KNN	Künstliches neuronales Netz
CNN	Convolutional Neural Network
KI	Künstliche Intelligenz
NN	neuronales Netz
ConvL	Convolutional Layer
FCL	Fully Connected Layer
Pool	Pooling Layer
ILSVRC	Large Scale Visual Recognition Challenge
DNN	Tiefe neuronale Netze
ML	Maschine Lernen
DSC	Depthwise Separable Convolution

1 Einleitung

2 Grundlagen

2.1 Neuron

Ein künstliches Neuron[18] ist eine mathematische Funktion, die das Verhalten vom biologischen Neuron nachbildet. Künstliche Neuronen sind elementare Einheiten in einem Künstliches neuronales Netz (KNN). Das künstliche Neuron empfängt einen oder mehrere Inputs und bildet sie auf einen Output ab. Normalerweise wird jeder Eingabe x_i separat mit einem Gewicht w_i multipliziert und danach aufsummiert und zum Schluss wird die Summe durch eine Funktion geleitet, die als Aktivierungs- oder Übertragungsfunktion bekannt ist. Eine schematische Darstellung eines künstlichen Neurons ist in Abbildung 1 zu sehen.

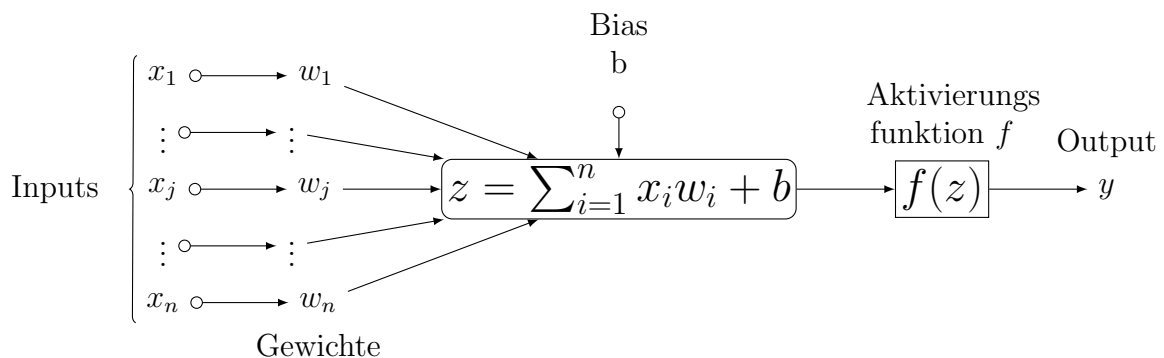


Abbildung 1: Funktionsweise eines künstlichen Neurons

Künstliche Neurone können aufgestapelt werden, um eine Schicht (*Layer*) zu bilden. Ein KNN besteht aus einer oder mehreren Schichten und je nach seiner Position in einem neuronalen Netz (NN) wird eine Schicht anders genannt: Eingangsschicht (*Input Layer*) bzw. Ausgabeschicht (*Output Layer*), wenn das Layer die Eingangsdaten bzw. Ausgabedaten des neuronalen Netzes darstellt und versteckte Schicht (*Hidden Layer*), wenn es keine Eingangs- oder Ausgabeschicht ist. Ein kurzer Überblick über die Darstellung von KNNs kann sich in Abbildung 18 verschafft werden.

2.2 Merkmalskarten (Feature-Maps)

Die Ausgabe einer Schicht wird als Aktivierungskarten oder Feature-Map(s) bezeichnet. Die Anzahl von Feature-Maps in einer Schicht ist gleich die Anzahl der Ausgabekanäle (*output channels*) bzw. Tiefe (*depth*) dieser Schicht. In einer Faltungsschicht ist ein Feature-Map immer zweidimensional, während die Dimension eines Feature-Maps in einem Fully Connected Layer (FCL) nur von der des Inputs abhängt, genauer gesagt, für ein n -dimensionales Inputs ist ein Feature-Map $(n - 1)$ -dimensional. Wie die Feature-Maps berechnet werden, hängt sehr vom Schichttyp ab. Es wird beispielsweise in einer Faltungsschicht die Faltungsoperation mehrmals auf die Eingabe angewendet und jedes Mal wird ein neues Feature-Map erhalten.

2.3 Filters

Ein Filter ist eine kleine Matrix, die die Extraktion von Features ermöglicht. Die modernen Architekturen von Convolutional Neural Network (CNN) verwendet immer mehr Filter (mindestens 2000), um bessere Ergebnisse zu erzielen. Wegen der zufällige Initialisierung von Filtern und ihrer großen Anzahl im CNN kann es vorkommen, dass mehrere Filter das gleiche Feature extrahieren, solche Filter werden als Redundant bezeichnet, oder dass Filter genau das Gegenteil von dem filtern, was wir wollen, nämlich unwichtige Informationen, deshalb werden während des Trainings die Filterparameter ständig geändert. Zwei nichtlineare Filter, d. h. Filter, die nicht proportional zueinander sind, lernen unterschiedliche Dinge. Man kann also daraus schließen, dass je größer die Anzahl an Filter ist, desto mehr Features gefiltert werden können oder besser sind die Ergebnisse, aber das ist in Allgemein falsch, da die guten Ergebnisse sehr stark davon abhängen, wie die Filter im CNN aufgestapelt sind. Zu viele Filtern in einem CNN können dazu führen, dass das CNN sowohl gute Features als auch schlechte schneller lernt und zu wenige Filter schränken die Kapazität des CNN ein, wichtige Features zu extrahieren, denn irgendwann wird es vorkommen, dass jedes Filter etwas Wichtiges gelernt hat, aber noch mehr Informationen werden benötigt, um die Trainingsdaten zu erklären. Es ist also wichtig die richtige Anzahl von Filtern zu finden.

•

2.4 Entwicklung von Künstlichen Neuronalen Netzen

3 Convolutional Neural Network

3.1 Feedforward

3.1.1 Input Layer

Die Eingangsschicht stellt die Eingangsdaten dar. Hier müssen die Eingangsdaten dreidimensional sein. Also die Eingangsdaten von CNN haben immer die folgende Form $W \times H \times D$ wobei (W, H) der räumlichen Dimension und D die Tiefe der Daten entspricht. Z.B $100 \times 100 \times 3$ für ein RGB-Bild und $224 \times 224 \times 1$ für ein Graustufenbild.

3.1.2 Faltungsschicht

Die wichtigste Sicht bzw. der Hauptschicht in einem CNN ist die Faltungsschicht (*Convolution Layer*). Die Eingabedaten eines CNN besteht auf jedenfalls aus wichtigen und unwichtigen Informationen. Als wichtige und unwichtige Informationen haben wir z.B die starke Präsenz der weißen Farbe in einer Muschelsuppe bzw. die Präsenz eines Menschen, wenn man verschiedenen Ernährungsklasse klassifizieren möchte. Während des Trainings eines CNN wird versucht, diese relevanten Informationen aus den Daten zu entnehmen und die irrelevanten auszuschließen. Alles was ein CNN aus den Eingabedaten nutzt, um

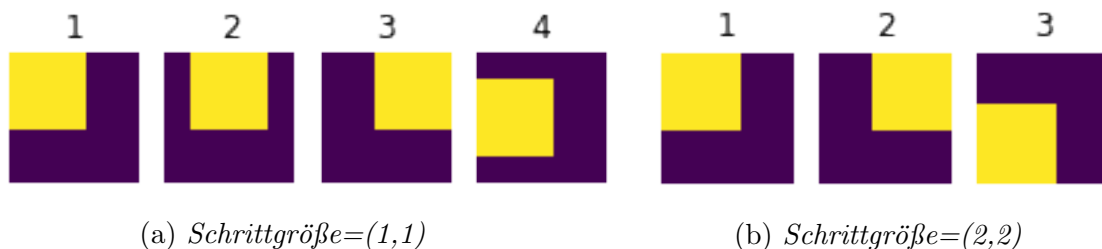
die Daten zu bestimmter Klassen zuzuordnen, wird als Feature bezeichnet. Das einzige bzw. das Hauptziel einer Faltungsschicht besteht darin, die Features aus seinen Inputdaten herauszuziehen. Dass eine Faltungsschicht in der Lage ist, Features selbst zu extrahieren, ohne dass man sie hinweist, was wichtig ist und was nicht, ist wirklich beeindruckend, Aber die Art und Weise, wie sie es tut, ist noch beeindruckender.

In einem Convolutional Layer (ConvL) wird die sogenannte Faltungsoperation (*convolution operation*) durchgeführt, dabei wird das komponentenweises Produkt (*Hadamard-Product*) zwischen einem kleinen Bereich der Eingabedaten und einem Kernel durchgeführt und dann die ganze aufsummiert. Eine Illustration der Faltungsoperation ist in der Abbildung 4a zu sehen. Die Resultante der Faltungsoperation wird die Aktivierung des Neurons genannt. Sollte diese Aktivierung des Neurons null sein, dann sagt man, dass das Neuron "nicht aktiv" ist und sonst ist es *aktiv*. Die null Aktivierung bedeutet, dass das vom Filter gesuchte Feature nicht gefunden wurde. Je stärker oder größer die Aktivierung eines Neurons ist, desto höher ist die Wahrscheinlichkeit, dass das vom Filter gesuchte Feature gefunden ist (siehe Abbildung 4b). Um diese Faltungsoperation durchzuführen, muss einige Hyperparameter vordefiniert sein, und zwar:

- **Anzahl und Größe von Filtern**

Die Anzahl an Filters gibt nicht nur an, wie viele Filter in dem ConvL verwendet werden, sondern auch, wie oft die Faltungsoperation auf die Eingabe durchgeführt wird, d.h die Anzahl der Feature-Maps oder die Tiefe der Schicht. Anstatt das ganze Bild zu betrachten, wenn man auf der Suche nach einem Feature ist, was mit der Tatsache gleichbedeutend ist, dass jedes Neuron der Schicht mit allen Neuronen der vorherigen verbunden ist, wird nur einen lokalen Bereich des Bildes betrachtet, wir verbinden also jedes Neuron nur mit einem lokalen Bereich der Bild. Dieser lokaler Bereich wird als Empfangsfeld (*receptive field*) des Neurons bezeichnet und entspricht der Filtergröße. Es ist zu beachten, dass die Filtergröße nur der räumlichen Dimension des Filters entspricht. Die Tiefe des Filters ist gleich die vom Input, also für ein $(100, 100, 3)$ Bild haben alle Filter die Tiefe 3. Die Verwendung von solchen kleinen Filtern ist die Hauptidee hinter einer Faltungsschicht. Filter haben im Allgemeinen eine kleine räumliche Dimension wie z.B 2×2 , 3×3 oder 5×5 , sonst verliert man einen großen Vorteil von ConvL, der darin besteht, die Speicheranforderung deutlich zu reduzieren, indem es die Gewichte verteilt.

Abbildung 2: Einfluss der Schrittgröße auf die Größe der Feature-Maps



- **Schrittgröße**(*Stride*)

Da ein Filter nur einen kleinen Bereich des Bilds wahrnehmen kann, wird eine Schrittgröße verwendet, um die Bewegung des Filters auf dem Bild zu steuern. Das Filter wird über das Bild von links nach rechts, von oben nach unten bewegt. Sei $S := (n, m)$ die Schrittgröße, dann wird das Filter von n Pixeln nach rechts für jede horizontale Bewegung des Filters und m Pixeln nach unten für jede vertikale Bewegung des Filters bewegt (siehe Abbildung 3a und 3b). Wie die folgende Tabelle zeigt, hängt die räumliche Dimension des Outputs sehr von der Schrittgröße ab.

Imagegröße	<i>Stride</i>	Filtergröße	Output (räumliche Dim.)	Mit Padding
(100, 100, 3)	(1, 1)	(1, 1)	(100, 100)	(100, 100)
(100, 100, 3)	(1, 1)	(3, 3)	(98, 98)	(100, 100)
(100, 100, 3)	(2, 2)	(3, 3)	(49, 49)	(50, 50)
(100, 100, 3)	(1, 1)	(4, 4)	(97, 97)	(100, 100)
(100, 100, 3)	(1, 1)	(10, 10)	(91, 91)	(100, 100)

Tabelle 1: Auswirkung von Schrittgröße, Filtergröße und Padding auf Output eines (100, 100, 3) Bild.

- **Padding**

Mit einem großen Filter lernt man in der Regel "mehr" als mit einem kleinen. Aber wie es in der Tabelle 1 zu sehen ist, habe man eine Reduktion der Dimension, wenn ein Filter mit Filtergröße > 1 angewendet wird und um die Dimension zu behalten, wird vor der Anwendung der Filter auf das Bild eine Füllung (*padding*) an den Rändern des Bilds gemacht. Diese Füllung ermöglicht erstens den Entwurf immer tiefer Netzwerke, denn es gibt nach jedem ConvL einen kleinen Dimensionverlust und zweitens, dass die Information an Rändern nicht zu schnell verschwunden werden. Um die Ränder einer Eingabe zu füllen, werden sehr oft Nullen (*zero padding*) verwendet, oder die Pixel, die an der Grenze liegen, werden wiederholt.

In einem CNN mit mehreren ConvLs beschäftigen sich die ersten ConvLs mit dem Erlernen einfacher Merkmale wie Winkel, Kanten oder Linien und je tiefer das CNN ist, desto komplexer sind die extrahierten Merkmale. Das liegt daran, dass jede nachfolgende Schicht einen größeren Bereich des Originalbildes betrachten oder „sehen“ kann. Nehmen wir an, dass die zwei ersten Layers eines CNNs Filter von Größe 3×3 und ein *Stride* = (1, 1) verwenden. Die erste Schicht betrachte immer 3×3 benachbarte Pixel des Originalbildes und speichert seine Aktivierung auf einem Pixel. Jetzt wenn die zweite ConvL 3×3 benachbarte Pixel betrachtet, betrachtet sie eigentlich 5×5 benachbarte Pixel des Originalbildes. Wir können uns deshalb vorstellen, dass irgendwo in späteren Schichten Filter das gesamte Originalbild betrachten.

(a) Matrixdarstellung

0	0	3	3	3	0	0	0	0
0	2	0	3	2	0	0	2	0
2	0	2	2	0	2	2	0	2
1	1	1	1	1	1	1	1	1
1	0	1	0	2	0	1	0	1
0	0	1	2	0	2	0	0	1
1	1	1	1	1	1	1	3	1
3	0	3	1	2	0	3	0	3
0	3	0	2	0	2	0	0	0

(9,9)Bild

 \times

0	0	0
0	1	0
1	0	1

Filter

 $=$

6	6	6
1	6	1
0	6	0

(3,3)Feature-Map

(b) Pixeldarstellung

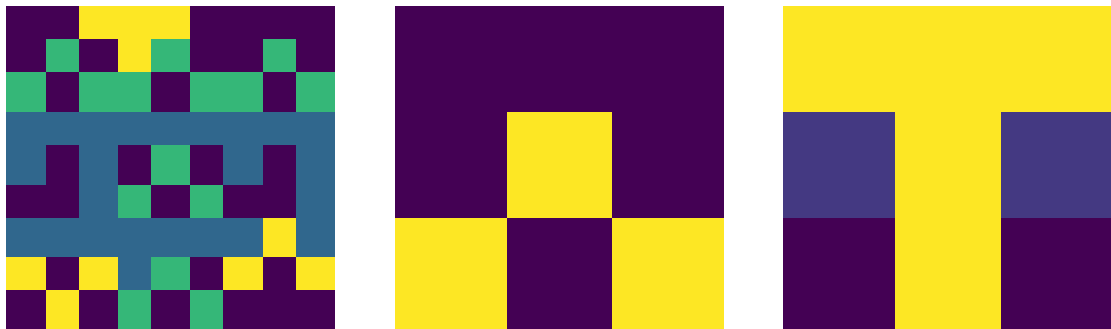


Abbildung 4: Faltungsoperation mit einem 3×3 -Filter und Schrittgröße = 3

3.1.3 Aktivierungsfunktion

Das neuronale Netzwerk wird während dem Training mit sehr vielen Daten gespeist und das sollte in der Lage sein, aus diesen Daten zwischen relevanten und irrelevanten Informationen Unterschied zu machen. Die Aktivierungsfunktion auch Transferfunktion oder Aktivitätsfunktion genannt, hilft dem NN bei der Durchführung dieser Trennung. Es gibt sehr viele Aktivierungsfunktionen und in folgenden werden wir sehen, dass eine Aktivierungsfunktion je nach zu lösende Aufgaben vorzuziehen ist.

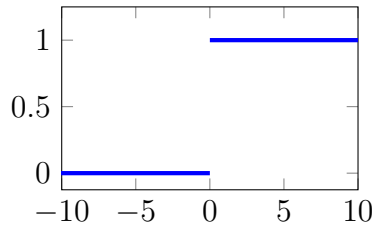
$$\begin{cases} Y = f(\Sigma(\text{Gewicht} * \text{Input} + \text{Bias})) \\ f := \text{Aktivierungsfunktion} \end{cases}$$

Binäre Treppenfunktion ist extrem einfach, siehe Abbildung 5, definiert als $f(x) =$

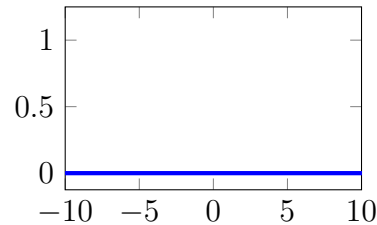
$$\begin{cases} 1, & \text{if } x \geq a \text{ (a:= Schwellenwert) } \\ 0, & \text{sonst} \end{cases}$$

Sie ist für binäre Probleme geeignet, also Probleme wo man mit *ja* oder *nein* antworten sollte. Sie kann leider nicht mehr angewendet werden, wenn es mehr als zwei Klassen klassifiziert werden soll oder wenn das Optimierungsverfahren gradientenbasierend ist, denn Gradient immer null.

Abbildung 5: Binäre Treppenfunktion



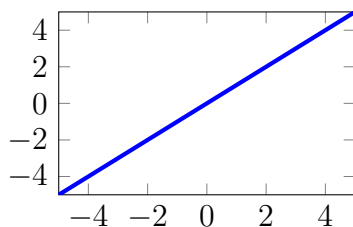
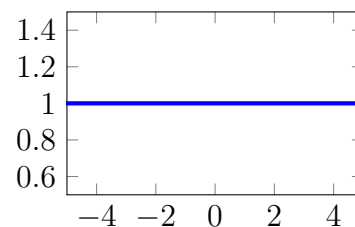
(a) Binäre Treppenfunktion



(b) Ableitung Binäre Treppenfunktion

Lineare Funktion ist definiert als $f(x) = ax$, $f'(x) = a$, siehe Abbildung 7. Sie ist monoton, null zentriert und differenzierbar. Es ist jetzt möglich, nicht mehr nur binäre Probleme zu lösen und mit gradientenbasierenden Optimierungsverfahren während der Backpropagation Parameter anzupassen, denn Gradient nicht mehr null, also sie ist besser als binäre Funktion. Nutzt ein mehrschichtiges Netz die lineare Aktivierungsfunktion, so kann es auf ein einschichtiges Netz überführt werden und mit einem einschichtigen Netz können komplexe Probleme nicht gelöst werden. Außerdem ist der Gradient konstant. Der Netzfehler wird also nach einigen Epochen nicht mehr minimiert und das Netz wird immer das Gleiche vorhersagen.

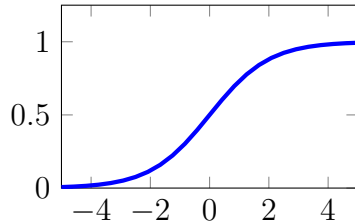
Abbildung 7: Lineare Funktion


 (a) Lineare Funktion: $f(x) = x$

 (b) Ableitung Lineare Funktion: $f'(x) = 1$

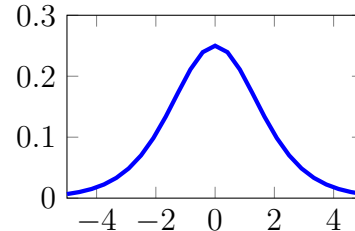
Logistische Funktion ist definiert als $f(x) = \frac{1}{1+\exp(-x)}$, $f'(x) = \frac{\exp(x)}{(1+\exp(x))^2}$, siehe Abbildung 9. Sie ist differenzierbar, monoton, nicht linear und nicht null zentriert (hier nur positive Werte). Zwischen $[-3, +3]$ ist der Gradient sehr hoch. Kleine Änderung in der Netzininput führt also zu einer großen Änderung der Netzausgabe. Diese Eigenschaft ist bei Klassifikationsproblemen sehr erwünscht. Die Ableitung ist glatt und von Netzininput abhängig. Parameter werden während der Backpropagation je nach Netzininput angepasst. Außerhalb von $[-3, 3]$ ist der Gradient fast gleich null, daher ist dort eine Verbesserung der Netzleistung fast nicht mehr möglich. Dieses Problem wird Verschwinden des

Gradienten (*vanishing gradient problem*) genannt. Außerdem konvergiert das Optimierungsverfahren sehr langsam und ist wegen der exponentiellen (e^x) Berechnung rechenintensiv.

Abbildung 9: Logistische Aktivierungsfunktion: $\text{sigmoid}(x)$.



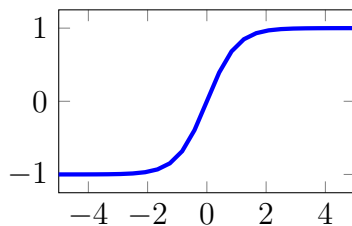
(a) Logistische Aktivierungsfunktion.



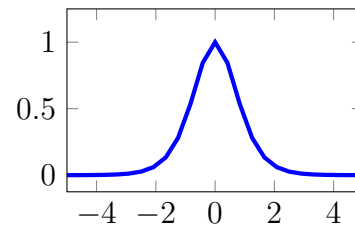
(b) Ableitung der Logistische Funktion.

Tangens Hyperbolicus ist definiert als $\tanh := 2\text{sigmoid}(x) - 1$, siehe Abbildung 11. Außerdem dass sie null zentriert ist, hat sie die gleichen Vor- und Nachteile wie die Sigmoid Funktion. **Sättigung fehlt noch**

Abbildung 11: Tangens Hyperbolicus.



(a) Tangens Hyperbolicus.

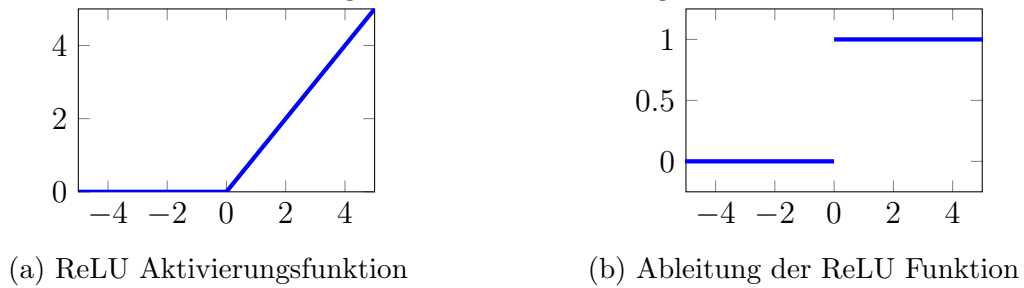


(b) Ableitung der Tangens Hyperbolicus.

Rectified Linear Unit (ReLU) ist definiert als $f(x) = \max(x, 0)$, siehe Abbildung 13. Sie ist sehr leicht zu berechnen. Es gibt keine Sättigung wie bei *Sigmoid* und *tanh*. Sie ist nicht linear, deshalb kann der Fehler schneller propagiert werden. Ein größter Vorteil der ReLU-Funktion ist, dass nicht alle Neurone gleichzeitig aktiviert sind, negative Eingangswerte werden zu null, daher hat die Ausgabe von Neuronen mit negativen Eingangswerten keinen Einfluss auf die Schichtausgabe, diese Neurone sind einfach nicht aktiv. Das Netz wird also spärlich und effizienter und wir haben eine Verbesserung der Rechenleistung. Es gibt keine Parameteranpassungen, wenn die Eingangswerte negativ sind, denn der Gradient ist dort null. Je nachdem wie die Bias initialisiert sind, werden mehrere Neurone getötet, also nie aktiviert und ReLU ist leider nicht null zentriert. **ReLU haben die wünschenswerte Eigenschaft, dass sie keine Eingangsnormalisierung benötigen, um eine Sättigung zu verhindern.**

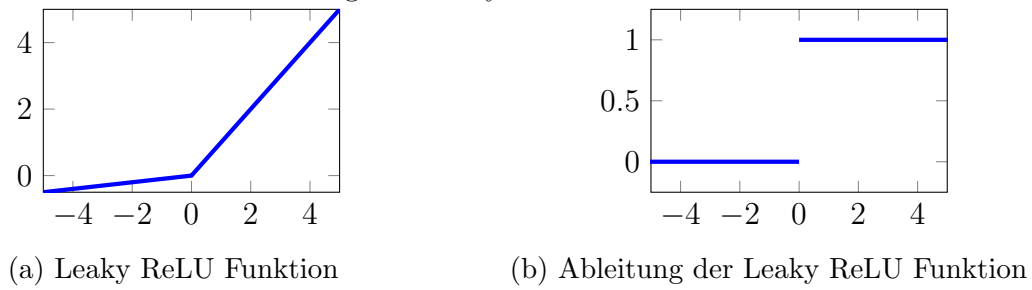
Leaky ReLU Funktion ist definiert als $f(x) = \begin{cases} x, & \text{if } x > 0 \\ 0.01x, & \text{sonst} \end{cases}$, siehe Abbildung 15. Sie funktioniert genauso wie die ReLU-Funktion, außer dass sie das Problem des

Abbildung 13: ReLU Aktivierungsfunktion



toten Neurons und sie ist null zentriert. Es gibt somit immer eine Verbesserung der Netzleistung, solange das Netz trainiert wird. Wenn das Problem von Leaky ReLU nicht gut gelöst wird, wird empfohlen, die *Parametric ReLU* (PReLU) Aktivierungsfunktion zu verwenden, die während der Training selber lernt, Problem der toten Neurone zu lösen.

Abbildung 15: Leaky ReLU Funktion



Softmax ist definiert als $f(x_1, x_2, \dots, x_n) = \frac{(e^{x_1}, e^{x_2}, \dots, e^{x_n})}{\sum_{i=1}^n e^{x_i}}$. Die Softmax-Funktion würde die Ausgänge für jede Klasse zwischen null und eins zusammendrücken und auch durch die Summe der Ausgänge teilen. Dies gibt im Wesentlichen die Wahrscheinlichkeit an, dass sich der Input in einer bestimmten Klasse befindet.

In allgemein wird die ReLU aufgrund des Problems der toten Neurone nur in versteckte Schichten und die Softmax-Funktion bei Klassifikationsproblemen und Sigmoid-Funktion bei Regressionsproblemen in Ausgabeschicht verwenden.

3.1.4 Pooling Layer

Die Funktionsweise von Pooling-Schichten ist sehr ähnlich zu der von ConvLs. Das Filter wird über die Inputdaten bewegt und dabei anstatt die Faltungsoperation durchzuführen, werden die Inputdaten Blockweise zusammengefasst. Ein Pooling-Layer besitzt nur eine Schrittgröße und eine Filtergröße. Das Filter in Pool ist in Gegensatz zu Filtern in ConvL nicht lernbar ist, es gibt nur an, wie groß der Block, der zusammengefasst wird, sein muss. Als Standard werden ein 2×2 Filter und eine 2×2 Schrittgröße verwendet, was die Dimension der Inputdaten um mindestens die Hälfte reduziert. Interessanter dabei ist, dass die wichtigen Informationen oder Muster nach der Pooling-Layer vorhanden bleiben und damit haben wir nicht sowohl eine Erhöhung der Rechengeschwindigkeit als

auch eine Reduzierung der Netzparameter. Noch dazu sind Pools invariant gegenüber kleiner Veränderung wie Parallelverschiebung.

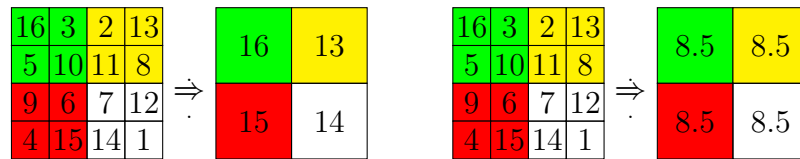


Abbildung 17: Funktionsweise der Pooling-Sicht mit Pooling_size= (2, 2) und Stride = 2

Je nachdem wie die Blöcke in Pooling Layer (Pool) zusammengefasst werden, haben die Pools unterschiedliche Namen. Werden die Werte eines Blocks durch den Maximalwert des Blocks, dann sprechen wir von Max-Pooling-Layer (siehe Abbildung 17 links), wenn sie durch den Mittelwert des Blocks ersetzt wird, wird von Average-Pooling-Layer (siehe Abbildung 17 rechts) und wenn die Filtergröße gleich die räumliche Dimension der Eingangsdaten ist, sprechen wir von Global-Max-Pooling-Layer und Global-Average-Pooling-Layer, es wird also alle Neurone in einem Kanal zu einem Neuron, die Ausgabedimension solche Schicht entspricht der Anzahl der Kanäle bzw. Tiefe der Inputdaten.

Das Global-Pooling-Layer wird sehr oft angewendet, um das Vorhandensein von Merkmale in Daten aggressiv zusammenzufassen. Es wird auch manchmal in Modellen als Alternative zur Flatten-Schicht, die mehrdimensionale Daten zu eindimensionale umwandelt, beim Übergang von ConvLs zu einem FCL verwendet.

3.1.5 Multi-layer Perzeptron (Fully Connected Layer)

Nachdem die relevanten lokalen Merkmale durch die Wiederholung von Conv und Pooling-Schichten extrahiert werden sind, wenden sie in einem FCL kombiniert, um das Ergebnis jeder Klasse zu berechnen. Die FCLs des CNN ermöglichen, Informationssignale zwischen jeder Eingangsdimension und jeder Ausgangsklasse zu mischen, so dass die Entscheidung auf dem gesamten Bild basieren kann und ihm eine Klasse zugewiesen werden kann. Die FCLs funktionieren eigentlich genau wie ConvLs, außer dass jedes Neuron in FCL mit allen Neuronen und nicht mit einem kleinen Bereich von Neuronen im vorherigen Layer verbunden ist. Ein NN mit nur FCLs sieht wie in Abbildung 18 aus.

Aufgrund der hohen Anzahl von Verbindungen zwischen Neuronen in einem FCL wird viel Speicher und Rechenleistung benötigt und verlangsamt auch das Training, es ist auch einer der Gründe, weshalb die FCLs meist nur in der letzten Schicht von CNN zur Klassifizierung verwendet werden und die Anzahl der Neuronen in letzter Schicht entspricht der Anzahl von Klassen.

3.2 Backforward

3.2.1 Fehlerfunktion

Das Training von CNNs besteht darin, den vom CNN begangenen Fehler zu korrigieren bzw. zu minimieren, daher wird es sehr oft als ein Optimierungsverfahren betrachtet.

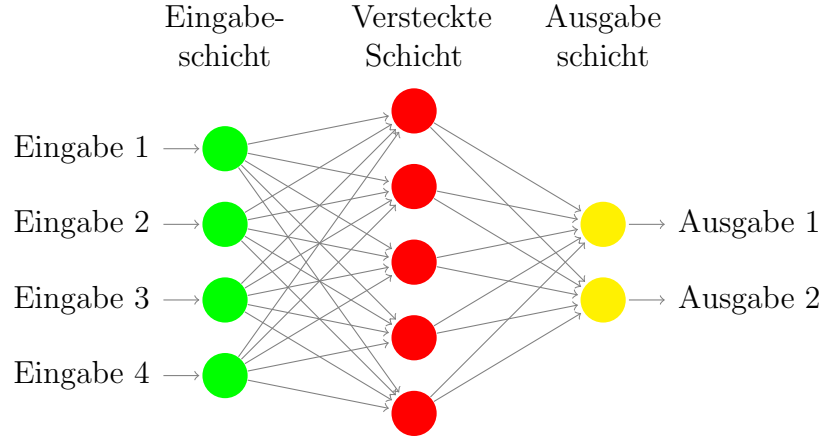


Abbildung 18: Darstellung eines neuronalen Netzes

Wie gut die Vorhersage des neuronalen Netzes gerade ist, wird durch eine Fehlerfunktion auch Kostenfunktion genannt quantifiziert oder angegeben. Die Fehlerfunktion bringt die Ausgabewerte des NNs mit den gewünschten Werten in Zusammenhang. Sie ist ein nicht-negativer Wert und je kleiner dieser Wert wird, desto besser ist die Übereinstimmung des CNNs. Es wird in Laufe des Trainings von CNNs versucht, diese Kostenfunktion mit Gradientenbasierten Verfahren zu minimieren (siehe Absatz 3.2.4).

Die meisten benutzten Kostenfunktionen sind die Kreuzentropie (*cross-entropy*, Gleichung 3.2) (CE) und die mittlere quadratische Fehler (*mean squared error*, Gleichung 3.1) (MSE).

$$MSE(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (3.1)$$

$$CE(Y, \hat{Y}) = -\frac{1}{n} \sum_{i=1}^n (Y_i \log(\hat{Y}_i) + (1 - Y_i) \log(1 - \hat{Y}_i)) \quad (3.2)$$

$Y := \{Y_1, \dots, Y_n\}$: Die tatsächlichen Werte

$\hat{Y} := \{\hat{Y}_1, \dots, \hat{Y}_n\}$: Die Ausgabewerte des neuronalen Netzes

Im Gegenteil zu CE Fehlerfunktionen, die sich nur auf Wahrscheinlichkeitsverteilungen anwenden lassen, können die MSE auf beliebige Werte angewendet werden. Nach [10, Pavel et al] ermöglicht die CE-Verlustfunktion ein besseres Finden lokaler Optima als die MSE-Verlustfunktion und das soll daran liegen, dass das Training des MSE Systems schneller in einem schlechten lokalen Optimum stecken bleibt, in dem der Gradient fast null ist und damit keine weitere Reduzierung der Netzfehler ermöglicht. Im Allgemeinen ist die CE Kostenfunktion für die Klassifikationsprobleme und die MSE Fehlerfunktion für die lineare Regression-Probleme besser.

3.2.2 Gradient

Der Gradient einer Funktion ist die erste Ableitung einer Funktion und in mehrdimensionalem Raum ist der Gradient einer Funktion der Vektor, dessen Einträge die ersten partiellen Ableitungen der Funktion sind. Der Gradient an einem Punkt gibt die Richtung der steilsten Anstieg der Funktion an diesem Punkt. Also da wir die Kostenfunktion minimieren möchten, sollen wir lieber immer in die Gegenrichtung des Gradienten gehen. In auf dem Gradienten basierten Optimierungsverfahren wird der Gradient benutzt, um die lokalen oder globalen Extremwerte (hier das Minimum) zu erreichen. Da wir jetzt die Richtung des Minimums herausgefunden haben, bleibt noch zu bestimmen, wie wir in diese Richtung gehen sollen.

3.2.3 Lernrate

Die Lernrate oder Schrittweite beim maschinellen Lernen ist ein Hyperparameter, der bestimmt, inwieweit neue gewonnene Informationen alte Informationen überschreiben sollen [13], in anderen Worten wie schnell wir ans Ziel kommen. Je nachdem wie die Lernrate gesetzt wird, werden bestimmte Verhalten beobachtet (Siehe Absatz 5.2) und sie nimmt sehr oft Werte zwischen 0.0001 und 0.5: Die Lernrate muss allerdings im Intervall $]0, 1[$ Werte annehmen, sonst ist das Verhalten des NN nicht vorhersehbar bzw. konvergiert das Verfahren einfach nicht. Für jeden Punkt x aus dem Parameterraum gibt es eine optimale Lernrate $\eta_{opt}(x)$, sodass das globale oder lokale Minimum sofort nach der Parameteranpassung erreicht wird. Da $\eta_{opt}(x)$ am Trainingsanfang leider nicht bekannt ist, wird die Lernrate in die Praxis vom Programmierer basiert auf seine Kenntnisse mit NNs oder einfach zufällig gesetzt.

3.2.4 Gradientenabstiegsverfahren

Aktuelle leistungsfähige DNN bestehen fast immer aus Million Variable (lernbarer Parameter). Wir können uns ein DNN als eine Gleichung mit Millionen von Variablen vorstellen, die wir lösen möchten. Mit Hilfe der Daten wollen wir uns in einem Raum, dessen Dimension größer als eine Million ist, bewegen, um die optimalen Parameter (Parameter, die die Trainingsdaten korrekt abbilden) zu finden. Aufgrund der unendlichen Anzahl von Punkten in solchen Räume wäre es nicht sinnvoll, einen Punkt zufällig auszuwählen, dann überprüfen, ob er optimal ist und, wenn nicht, nochmals einen anderen Punkt zufällig auszuwählen. Genauer zu diesem Zeitpunkt kommen Gradientenabstiegsverfahren zum Einsatz. Die Gradientenabstiegsverfahren sind Verfahren, die auf dem Gradienten basieren, um Optimierungsprobleme zu lösen. Hier wird Gradientenabstiegsverfahren verwendet, um sich der optimalen Parametern anzunähern oder sie zu finden.

Ablauf eines Gradientenverfahrens im DNN

Das Gradientenabstiegsverfahren kann in drei Hauptschritte aufgeteilt werden. Die Abbildung 19 stellt das Backpropagation-Verfahren bildlich dar. Beim ersten Schritt wird ein zufälliger Punkt aus dem Parameterraum ausgewählt und davon ausgehend wird

der Parameterraum exploriert. Das entspricht der Netzparameterinitialisierung am Trainingsanfang. Der zweite Schritt besteht darin, die Abstiegsrichtung zu bestimmen. Dazu werden zuerst die Eingangsdaten in das NN eingespeist (*Forwardpropagation*), danach wird der Fehler zwischen den Netzvorhersagen und den korrekten Werten berechnet und ein Fehler gibt es (fast) immer, denn die Initialisierung wird zufällig gemacht und die Wahrscheinlichkeit, dass wir von Anfang an die optimalen Werte finden, ist verschwindend klein und zuletzt wird der Gradient (3.2.2) der Kostenfunktion in abhängig von den gegebenen Eingangsdaten und den erwarteten Werten berechnet. Beim letzten Schritt wird die Schrittweite bestimmen. Die Lernrate (3.2.3) oder die Schrittweite wird vor Trainingsbeginn festgelegt oder während des Trainings abhängig von aktuellem Netzzustand allmählich adaptiert.

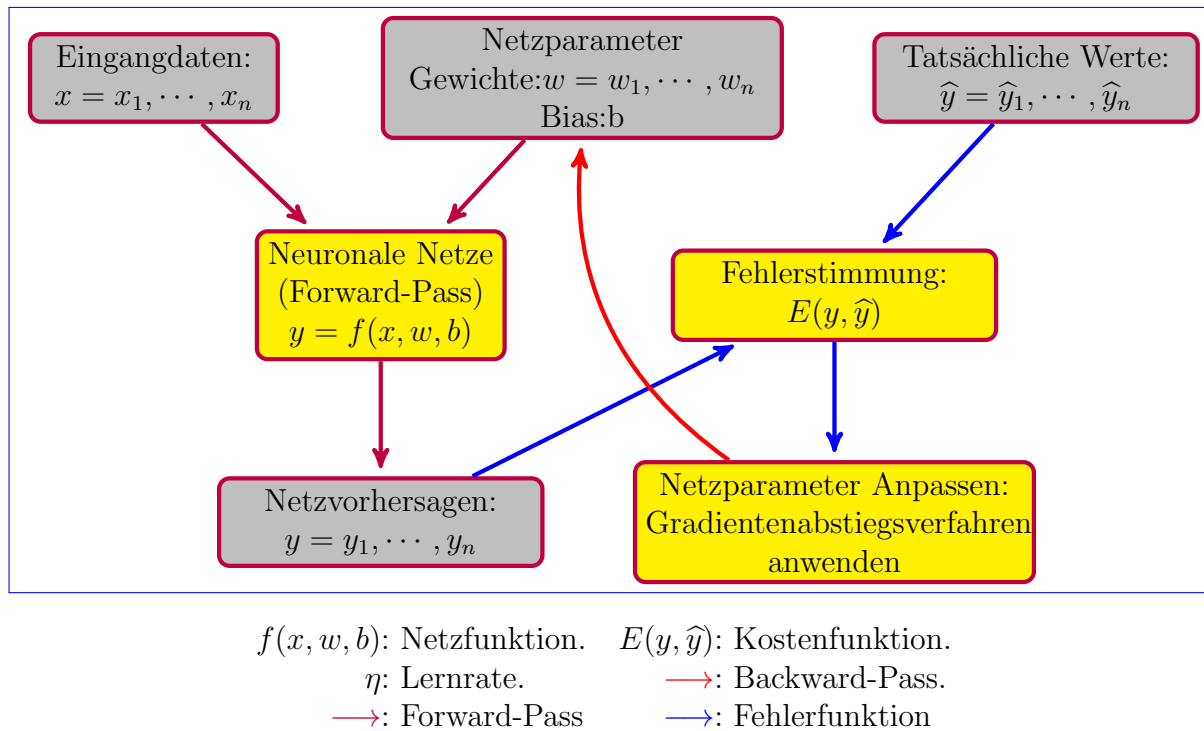


Abbildung 19: Ablauf der Backpropagation

Variante des Gradientenverfahrens

Bisher existiert drei Variante des Gradientenabstiegsverfahren, die sich nur durch die Größe der Daten, die sie verwendet, um den Gradienten der Kostenfunktion berechnet, unterscheidet.

Stochastic Gradient Descent (SGD) : Bei SGD wird jeweils ein Element bzw. Sample aus der Trainingsmenge durch das NN durchlaufen und den jeweiligen Gradienten berechnen, um die Netzwerkparameter zu aktualisieren. Diese Methode wird sehr oft

online Training genannt, denn jedes Sample aktualisiert das Netzwerk. SGD verwendet geringer Speicherplatz und die Iterationen sind schnell durchführbar. Zusätzlich kann die Konvergenz für großen Datensatz wegen der ständigen Aktualisierung der Netzwerkparameter beschleunigt werden. Diese ständigen Aktualisierung hat die Schwankung der Schritte in Richtung der Minima zur Folge, was die Anzahl der Iteration bis zum Erreichen des Minimums deutlich ansteigt und dabei helfen kann, aus einem unerwünschten lokalen Minimum zu entkommen. Ein großer Nachteil dieses Verfahren ist der Verlust der parallelen Ausführung, es kann jeweils nur ein Sample ins NN eingespeist werden. Der Algorithmus 1 zeigt den Ablauf von SGD.

Input: loss function E , learning rate η , dataset X, y und das Modell $F(\theta, x)$
Output: Optimum θ which minimizes E

```

1 while converge do
2   | Shuffle  $X, y$ 
3   | for  $x_i, y_i$  in  $X, y$  do
4   |   |  $\tilde{y} = F(\theta, x_i)$ 
5   |   |  $\theta = \theta - \eta \cdot \frac{1}{N} \sum_{i=1}^N \frac{\partial E(y_i, \tilde{y}_i)}{\partial \theta}$ 
6   | end
7 end

```

Algorithm 1: Stochastic Gradient descent(SGD) [2].

Batch Gradient Descent (BGD) : BGD funktioniert genauso wie SGD, außer dass der ganze Datensatz statt jeweils ein Element aus dem Datensatz zur Netzwerkparameteraktualisierung genutzt wird. Jetzt kann das Verfahren einfach parallel ausgeführt werden, was den Verarbeitungsprozess des Datensatzes stark beschleunigt. BGD weist weniger Schwankungen in Richtung des Minimums der Kostenfunktion als SGD auf, was das Gradientenabstiegverfahren stabiler macht. Außerdem ist das BGD rechenefizienter als das SGD, denn nicht alle Ressourcen werden für die Verarbeitung eines Samples, sondern für den ganzen Datensatz verwendet. BGD ist leider sehr langsam, denn die Verarbeitung des ganzen Datensatz kann lange dauern und es ist nicht immer anwendbar, denn sehr große Datensätze lassen sich nicht im Speicher einspeichern. Der Algorithmus 2 zeigt den Ablauf von BGD.

Input: loss function E , learning rate η , dataset X, y und das Modell $F(\theta, x)$
Output: Optimum θ which minimizes ϵ

```

1 while converge do
2   |  $\tilde{y} = F(\theta, x)$ 
3   |  $\theta = \theta - \eta \cdot \frac{1}{N} \sum_{i=1}^N \frac{\partial \epsilon(y, \tilde{y})}{\partial \theta}$ 
4 end

```

Algorithm 2: Batch Gradient descent [2].

Mini-batch Stochastic Gradient Descent(MSGD): MSGD ist eine Mischung aus SGD und BGD. Dabei wird der Datensatz in kleine Mengen (*Mini-Batch oder Batch*) möglicherweise gleicher Größe aufgeteilt. Je nachdem wie man die Batch-Größe setzt, enthalten wir SGD oder BGD wieder. Das Training wird Batchweise durchgeführt, d.h. es wird jeweils ein Batch durch das NN propagiert, der Verlust jedes Sample im Batch wird berechnet und dann deren Durchschnitt benutzt, um die Netzwerkparameter zu anzupassen. MSGD verwendet den Speicherplatz effizienter und kann von Parallelen Ausführung profitieren. Noch dazu konvergiert MSGD schneller und ist stabiler. In die Praxis wird fast immer das MSGD Verfahren bevorzugt. Der Algorithmus 3 zeigt den Ablauf von MSGD

<p>Input: loss function E, learning rate η, dataset X, y und das Modell $F(\theta, x)$ Output: Optimum θ which minimizes E</p> <pre> 1 while converge do 2 Shuffle X, y 3 for each batch of x_i, y_i in X, y do 4 $\tilde{y} = F(\theta, x_i)$ 5 $\theta = \theta - \eta \cdot \frac{1}{N} \sum_{i=1}^N \frac{\partial E(y_i, \tilde{y}_i)}{\partial \theta}$ 6 end 7 end </pre>

Algorithm 3: Mini-Batch Stochastic Gradient descent(MSGD) [2].

4 Kompression von DNN

Die neueren maschinellen Lernmethoden verwenden immer tiefer neuronale Netze wie z.B. *Xception(134 Layers)*, *MobileNetV2(157 Layers)*, *InceptionResNetV2(782 Layers)*, um Ergebnisse auf dem neuesten Stand der Technik zu erzielen. Aber die Verwendung von sehr tiefer NNs bringt mit sich nicht nur eine deutliche Verbesserung der Modellleistung, sondern auch einen bedeutenden Bedarf an Rechenleistung und an Speicherplatz, was der Einsatz solcher Modelle auf Echtzeitsystemen mit begrenzten Hardware-Ressourcen schwierig macht. Es wurden bisher mehrere Ansätze untersucht, um die dem NN zugewiesenen Ressourcen effizienter zu nutzen:

- Die Modellbeschneidung (*Network pruning*), die die redundanten und die nicht relevanten Verbindungen zwischen Neuronen entfernt.
- Die Destillation von NNs, die ermöglicht, die großen Modellen in kleinen zu komprimieren.
- Die Quantisierung von NN, die für die Darstellung von einzelner Netzparameter weniger als 32 Bits nutzt.
- Huffman-Codierung, die eine komprimierte Darstellung des Netzwerks ermöglicht.

Im folgenden werden nur die Beschneidung, die Quantisierung von NN und die Anwendung von Huffman auf NN mehr eingegangen werden.

4.1 Pruning Network

Wie oben schon erwähnt, wird beim *Pruning* neuronaler Netzwerke versucht, unwichtige oder redundante Verbindungen oder komplette Neuronen aus dem Netzwerk zu entfernen, um ein Netz mit möglichst geringer Komplexität zu erhalten. Mit unwichtigen Verbindungen werden die Parameter (Gewichte und Bias) gemeint, die fast null sind, denn Parameter mit Nullwert haben keinen Einfluss auf das Output des Neurons, sie sind einfach überflüssig. Während oder nach dem Training gibt es mehrere Parameter, die nicht wirklich oder nicht zu viel zum Neuronenergebnis beitragen, obwohl sie keinen Nullwert haben, deshalb ist es zum Reduzieren der Netzwerkdicke notwendig, anderen Maßstäbe als den Nullwert anzulegen, um Verbindungen zu entfernen.

Das Pruning-Verfahren bietet einige Vorteile wie Reduzierung der Speicher- und Hardwarekosten, die Trainingsbeschleunigung, die schnellere Antwortzeit und das Verringern der Wahrscheinlichkeit der Overfitting (5.6.1). Es ist sehr wichtig zu beachten, dass die Anwendung vom Pruning-Verfahren auf ein NN nur Sinn macht, wenn das NN teilweise oder komplett trainiert ist, sonst macht das Pruning nur eine Reduktion der Anzahl der Netzwerkparameter.

Es gibt zwei Hauptszenarien für das Pruning von NN. Die erste besteht darin, die irrelevanten Verbindungen in einem komplett trainierten NN zu entfernen. Mit komplett trainierten NN wird gemeint, dass die erwünschte Genauigkeit schon erreicht ist. Im zweiten Szenario wird Pruning während des Trainings durchgeführt, das wird sehr oft als *iteratives Pruning* bezeichnet. Dabei wird vor Trainingsbeginn bestimmte Dinge festgelegt, wie z.B. ab wann wird das Netzwerk beschnitten und wie oft es durchgeführt werden soll. Das erste Szenario ist einfacher anzuwenden, denn man muss nur darauf warten, bis es keine Verbesserung der Genauigkeit des NN mehr gibt und dann Pruning auf das NN anwenden, aber damit verliert man großen Vorteile des Pruning, die die Beschleunigung des Trainings und Reduzierung der Overfitting-Wahrscheinlichkeit sind. Aus diesen Gründen wird in der Praxis das zweite Szenario bevorzugt, aber die richtigen Hyperparameter zu finden, ist kompliziert. Eine zu früh Beschneidung kann z.B. die Genauigkeit des NN zu sehr verschlechtern, denn es kann sein, dass eine Verbindung, die nur nach einer späteren Gewichts Anpassung zum Ergebnis eines Neurons hätte beigetragen können, entfernt würde, aber es bietet eine Beschleunigung des Trainings. Was eine zu spät Beschneidung angeht, haben wir fast die gleichen Nachteile wie im ersten Szenario und eine sehr häufiges Pruning kann das Training auch verlangsamen. Das ganze Prozess muss sehr oft leider mehrmals angewendet, um die richtigen Hyperparameter zu finden, sodass es sich wirklich nicht mehr lohnt. Ein graphischer Ablauf der Beschneidung ist in der Abbildung 20 zu sehen. Wenn wir zu viel auf einmal schneiden, könnte das Netzwerk so sehr beschädigt werden, dass es nicht mehr wiederhergestellt werden kann.

Die Hauptarbeit bei Pruning-Verfahren ist sicherlich, die guten Kriterien für die Bewertung der Wichtigkeit von Parametern zu finden und es ist vielleicht einer der Gründe, warum Pruning-Verfahren bisher nicht so populär ist, obwohl es immer noch zu schwierig

ist, tiefe NNs beispielsweise auf mobile Geräten mit eingeschränkten Ressourcen durchzuführen.

Bisher gibt es viele Kriterien für die Bewertung der Wichtigkeit eines Parameters, die sich miteinander unterscheiden und je nach Anwendung Vor- und Nachteile aufweisen. Im Folgenden werden einige Kriterien vorgestellt.

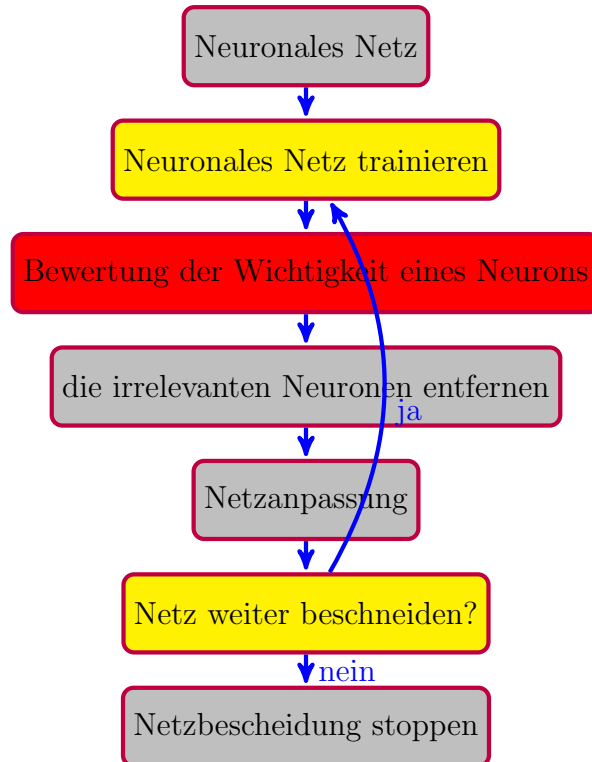


Abbildung 20: Netzbeschneidung während des Trainings

- **Schwellenwert**

Für die Bewertung der Wichtigkeit eines Parameters verwenden [6, Han et al] einen Schwellenwert θ , also alle Parameter mit einem Wert in $[-\theta, \theta]$ werden eliminiert. So ein Kriterium macht das Pruning schneller und effizienter, weil es sehr einfach ist, die zu eliminierenden Parameter zu finden. Da nach dem Pruning immer ein Genauigkeitsverlust auftritt, muss das NN erneuert trainiert werden, um seine Genauigkeit wiederherzustellen. Aber dieser Ansatz hat viele Nachteile. Erstens muss der optimale Schwellenwert gefunden werden und dafür muss das ganze Prozess (Schwellenwert auswählen, NN beschneiden, Verlust vergleichen) mehrmals wiederholt werden, um nur den optimalen Schwellenwert zu finden, deshalb können wir nie sicher sein, dass wir den optimalen Schwellenwert gefunden haben. Zweitens wird eine große Reduktion der Rechenkosten nur in FCLs und nicht in ConvLs beobachtet [7]. Die modernen Architekturen von CNNs können also bezüglich der Rechenkosten aus dem Verfahren keinen Vorteil ziehen, denn sie bestehen am meisten aus Faltungsschichten. Drittens muss spärliche Bibliotheken oder spezielle

Hardware verwendet werden, damit die Bewertung des beschnitten NN effektiv wird.

- **Filter und Feature-Map Pruning**

Anstatt die Gewichte in einem Filter elementweise zu entfernen, schlagen [7, Li et al] vor, das ganze Filter und die entsprechende Feature-Map zu entfernen. Dieses Vorgehen bringt mit sich im Vergleich mit dem oben vorgestellten Ansatz mehrere Vorteile. Es werden z.B. keine Bibliotheken, die eine Beschleunigung durch spärliche Operationen über CNNs ermöglichen, mehr benötigt. Noch dazu werden die Rechenkosten deutlich reduziert. Bei diesem Ansatz werden die weniger nützlichen Filter aus einem vollständig trainierten CNN entfernt. Zur Auswahl der zu entfernenden Filter wird zuerst die l_1 -Norm (4.1) jedes Filter im NN berechnet und dann werden die m Filter mit der kleinsten l_1 -Norm in jeder Schicht entfernt, wobei m ein Hyperparameter ist, der die Anzahl der zu löschenden Filter angibt.

$$\|F\|_1 = \sum_{i=1}^h \sum_{j=1}^w |F_{i,j}| \quad (4.1)$$

$(w, h) :=$ Breite und Höhe des Filters

Der Grund, warum nur die Filter mit einer kleineren l_1 -Norm entfernt werden, liegt daran, dass sie neigen dazu, Feature-Maps mit geringen Aktivierungen im Vergleich zu den anderen Filtern in der selben Schicht zu erzeugen und das kann gut in der Abbildung 21 festgestellt werden. Wobei das Feature-Map links aus der Abbildung 21 ist das selbe mit dem aus Abbildung 4b und das Feature-Map rechts erhalten wir, indem wir das Filter aus Abb. 4a mit drei multiplizieren.



Abbildung 21: Einfluss der Intensität des Filters auf Feature-Map

Noch interessanter an diesem Verfahren ist, dass das Pruning und das neue Training des Netzwerks auf einmal auf mehrere Schichten durchgeführt werden, was die Beschneidungszeit noch weiter beschleunigt und noch effizienter ist, wenn es um sehr tiefe Netzwerke wie *InceptionResNet* oder *GoogLeNet* angeht. Zur Beschneidung von Filter auf mehrere Schichten kann man die zu entfernenden Filter auf jeder Schicht entweder unabhängig von anderen Schicht oder nicht auswählen. Sollte man die zu löschenden Filter auf jeder Schicht unabhängig von anderen Schichten bestimmen, so muss man mit höheren Rechenkosten rechnen, denn es

werden Filter, die in vorherigen Schichten schon ausgenommen wurden, in der Berechnung der Summe der absoluten Gewichte noch miteinbezogen. Nach Li et al [7] ist die zweite Strategie nicht global optimal und kann trotzdem unter bestimmten Umständen zu besseren Ergebnissen führen. Die Abbildung 22 zeigt, wie die Entfernung eines Filters die Sichtausgabe beeinflusst.

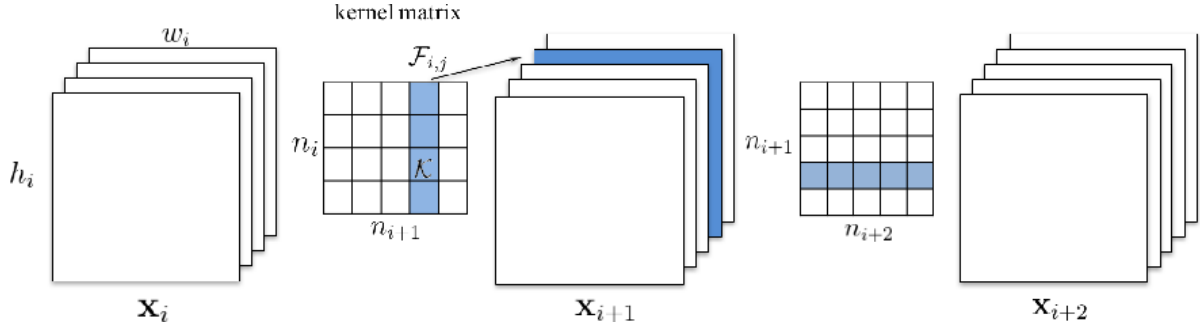


Abbildung 22: Das Beschneiden eines Filters führt zum Entfernen der entsprechenden Feature-Map und der zugehörigen Kernel in der nächsten Ebene[7].

• Automatische Beschneidung

Die bisher beschriebenen Kriterien oder die meisten von ihnen fordern, dass das beschnittene Netzwerk noch trainiert wird, um die Endgewichte für die verbleibenden spärlichen Verbindungen zu lernen. Das ganze macht das Pruning-Verfahren noch schwieriger, denn es wird nur dann unterbrochen, wenn der Genauigkeitsverlust wirklich groß ist, was zu viel Zeit in Anspruch nimmt. Um dieses Problem zu lösen, schlagen *Manessi et al*[8] eine differenzierbare Technik vor, die es ermöglicht, das Pruning durchzuführen, wenn die Gewichte des Netzwerks während der Trainingsphase angepasst werden und die optimalen Schwellwerte automatisch zu suchen. Dieses Verfahren ist inspiriert von dem von Han et al[6] und löst viele Probleme, die das Verfahren von *Han et al* nicht lösen könnten oder wegen der Art und Weise, wie sie das Pruning durchführen, entstanden. Erstens werden die Schwellwerte während des Trainings wie die Gewichte auch gelernt, es kann also eine große Menge von unterschiedlichen Schwellwerten ausprobiert werden und anstatt der selbe Schwellwert auf alle Layers anzuwenden, wird jeder Schicht eine Schwelle zugewiesen. Da Nicht alle Layers gleich empfindlich gegenüber dem Pruning sind und die Verwendung eine selbe Schwelle auf jede Schicht davon ausgeht, dass sie sind, muss dieses Verfahren optimaler. Zweitens, anstatt die Trainingszeit zu erhöhen, weil auch neue Parameter gelernt werden müssen, wird sie reduziert, das liegt daran, dass die Netzwerkparameter irgendwann mehr oder weniger spärlich werden, was Feedforward- und Backpropagationzeit erheblich reduziert. Es ist wichtig zu betonen, dass das Netzwerk nur einmal trainiert werden muss, um eine selbe bzw. bessere Genauigkeit wie bzw. als die anderen Methoden zu erzielen.

Aufgrund seiner Automatisierung und Effizienz wird das Pruning in Zukunft sicherlich immer häufiger eingesetzt werden.

4.2 Quantisierung von neuronalen Netzwerken

Im Allgemein wird von Quantisierung gesprochen, wenn es versucht wird, sich den Werten einer großen Menge mit denen einer kleineren Menge zu nähern. Im Maschine Lernen (ML) ist die Quantisierung der Prozess der Transformation eines ML-Programms in eine approximierte Darstellung mit verfügbaren Operationen mit geringerer Genauigkeit.

Methoden

- Matrixfaktorisierung
- Vektorquantisierung
- low-bit Quantisierung

Matrixfaktorisierung

Die Matrixfaktorisierung versucht, eine Matrix durch die Singulärwertzerlegung (SVD) zu approximieren. Die Singulärwertzerlegung zur Annäherung an eine Matrix $A \in \mathbb{R}^{m \times n}$ verwendet eine Matrix $B \in \mathbb{R}^{m \times n}$ vom Rang $k \leq \min(m, n)$, siehe Gleichung(4.2), wobei U und V orthogonale Matrizen und S eine Diagonalmatrix, deren Diagonaleinträge Singulärwerte heißen und bei Betrachtung von oben nach unten monoton abnehmen. Sollte man B direkt speichern, so hat man keine Kompression gemacht, sondern nur Informationen und Rechenzeiten verschwendet, denn die Matrizen A und B haben die gleiche Dimension, also liegt der Hauptvorteil von der Matrixfaktorisierung darin, wie die Matrix B gespeichert wird.

$$\begin{aligned} A &= USV^T \\ B &= \hat{U}\hat{S}\hat{V}^T = U[:, :k]S[:, :k]V[:, :k]^T \\ \|A - B\|_F &= \sqrt{\sum_{i=1}^m \sum_{j=1}^n (A - B)_{i,j}^2} := \text{Frobenius Norm} \end{aligned} \quad (4.2)$$

Eigentlich werden \hat{U} , \hat{S} und \hat{V} statt B gespeichert, damit wird $mk + kk + kn$ statt mn Speicherplatz verwendet und da \hat{S} Diagonalmatrix ist, kann auch einfach nur die Diagonaleinträge gespeichert werden, was zu einer Kompressionsrate von $\frac{mn}{k(m+n+1)}$ führt. Die Verwendung von Matrizen von kleinen Rang hebt die Tatsache hervor, dass es redundante Parameter in NN-Parametern gibt. Diese Redundanzeigenschaft wird ausgenutzt, um die Geschwindigkeit zu erhöhen [9, Denton et al]. Im Sinne der Frobeniusnorm ist B die bestmögliche Approximation von A mit Singulärwertzerlegung durch eine Matrix von Rang k . Angenommen bildet die zu approximierenden NN-Parameter ein Bild, dann zeigt die Abbildung 23, wie man die NN-Parameter approximieren kann.

Wie man in der Abbildung 23 sieht, reichen relativ wenige Singulärwerte (ca. 5) schon aus, um die groben Konturen zu erkennen. Details erkennt man bei ca. 10 Singulärwerten schon. Generell wird das Bild mit mehr Singulärwerten schärfer.

Quantisierung mit weniger Bits (Low-bit Quantization)

Im Bereich des *Deep Learning* ist das Standard numerische Format für Forschung und Einsatz bisher 32-Bit Fließkommazahlen oder FP32, denn es bietet eine bessere Genauigkeit, aber die anderen Formate wie 8-, 4-, 2- oder 1-Bits werden auch verwendet, obwohl



SW: Singulärwerte

Abbildung 23: Parameter Approximieren durch Matrixfaktorisierung

sie mehr oder weniger ein Verlust an Genauigkeit aufweisen. Die Verwendung von weniger genauen numerischen Formaten hat nicht nur einen kleinen Verlust der Netzleistung zur Folge, sondern auch die Verwendung von deutlich reduzierter Bandbreite und Speicherplatz. Noch dazu beschleunigt die Quantisierung die Berechnungen, denn die ganzzahlige Berechnung zum Beispiel ist schneller als die Fließkommaberechnung.

Bei Quantisierung mit wenigen Bits geht es mehr oder weniger um die Abbildung eines großen Bereiches auf einen kleinen und dazu werden zwei Hauptwerte benötigt: der dynamische Bereich des Tensors und ein Skalierungsfaktor. Angenommen haben wir einen dynamischen Bereich $[0, 500]$ und einen Skalierungsfaktor 5, dann ergibt sich der neue Bereich $[0, 100]$, es werden also Werte zwischen $[5k, 5(k+1)]$ oder $[5k - 0.5, 5k + 0.5]$ auf $5k$ abgebildet. Es ist sinnvoller, die Skalierungsfaktoren unter Berücksichtigung der Anzahl und der Verteilung der Werte im dynamischen Bereich des Tensors auszuwählen, sonst können zu viel Informationen verloren gehen.

Sobald ein anderes Format mit weniger Bits als FP32 verwendet, spricht man in ML von Quantisierung. Als Hauptvorteile von Quantisierung haben wir: Erstens wenn wir ein Modell mit FP32-Format durch das gleiche Modell mit 16-Bit Float (FP16), 8-Bit Integer (INT8) oder 4-Bit Integer (INT4) ersetzen, reduzieren wir den Speicherbedarf um die Hälfte, um ein Viertel bzw. um ein Achtel. Zweitens sind die Hardwares so programmiert, dass die Integer-Operationen im Vergleich zu FP32-Operationen schneller und energieeffizienter sind und drittens wird die Bandbreite durch die Verwendung von kleineren Modellen und dynamischen Werten stark reduziert.

Zur effizienteren Auswahl der Skalierungsfaktoren pro Schicht bzw. pro Filter ist es notwendig, Statistiken über das NN zu sammeln und das kann *offline* oder *online* gemacht werden. Bei der *Offline* Berechnung werden vor der Bereitstellung des Modells einige Statistiken gesammelt, entweder während des Trainings oder durch die Ausführung einiger Epochen auf dem trainierten FP32-Modell und basierend auf diesen Statistiken werden die verschiedenen Skalierungsfaktoren berechnet und nach der Bereitstellung des Modells festgelegt. Bei dieser Methode besteht die Gefahr, dass zur Laufzeit die Werte, die außerhalb des zuvor definierten Bereichs auftreten, abgeschnitten werden, was zu einer Verschlechterung der Genauigkeit führen kann. Bei der *online* werden die *Min/Max*-

Werte für jeden Tensor dynamisch zur Laufzeit berechnet. Bei dieser Methode kann es nicht zu einer Beschneidung kommen, jedoch können die zusätzlichen Rechenressourcen, die zur Berechnung der Min/Max-Werte zur Laufzeit benötigt werden, unerschwinglich sein.[16]

Es gibt zwei Quantifizierungsszenarien. Die erste ist das vollständige Training eines Modells mit einer gewünschten niedrigeren Bit-Genauigkeit (*bit precision*) (kleiner als 32 Bits). Das Training mit sehr geringer Genauigkeit ermöglicht ein potenziell schnelles Training und Inferenz, aber der Hauptproblem mit diesem Ansatz ist, dass Netzparameter nur bestimmte Werte annehmen können, so ist die Aktualisierung der Netzparameter bzw. das Backpropagation nicht mehr wohldefiniert [17]. Das zweite Szenario quantisiert ein trainiertes FP32-Netzwerks mit einer geringeren Bit-Genauigkeit ohne vollständiges Training. Im Allgemeinen gilt: Je geringer die Bitgenauigkeit, desto größer ist der Verlust der Genauigkeit und um diesen Leistungsabfall zu überwinden, wird sehr oft das NN erneuert trainiert oder wird auf eine hybride Quantisierung, die zur Quantisierung verschiedene Formaten z.B. INT8 für die Gewichte und FP32 für die Aktivierung gleichzeitig verwendet, zurückgegriffen.

Vor kurzem haben [17, Yoni et al] zur Quantisierung das lineare Quantisierungsproblem als ein *Minimum Mean Squared Error* (MMSE) Problem formuliert und gelöst. Sie sind in der Lage, die trainierten Modelle zu quantisieren, ohne das NN erneuert trainieren zu haben um seine Genauigkeit wiederherzustellen und benutzen dabei nur das INT4-Format. Obwohl diese Methode minimalen Verlust der Genauigkeit (*accuracy*) aufweist, liefert sie Ergebnisse auf dem neuesten Stand der Technik und nach [17] weist dieser Ansatz den geringeren Genauigkeitsverlust als alle Quantisierungsverfahren auf.

Wie oben erwähnt, eine Quantisierung mit zu weniger Bit führt zu einem großen Genauigkeitsverlust, aber [17, Yoni et al] haben gezeigt, dass Die Quantisierung mit INT4-Format zu einem besseren Ergebnis führen kann. Das sollte daran liegen, dass sie fürs Approximieren eines Tensors, der gegenüber der Quantisierung sehr empfindlich ist, statt nur einen Tensor mehrere Tensors im INT4-Format benutzen.

4.3 Huffman Codierung

Die Huffman-Codierung ermöglicht eine verlustfreie Datenkompression, indem sie jeder einzelnen Dateneinheit eine unterschiedlich lange Folge von Bits zuordnet. Daraus folgt, dass eine gute Möglichkeit zur besseren Verwaltung der dem Modell zugeordneten Ressourcen ist: Erstmal das Netzwerk zu beschneiden, dann zu quantisieren und am Ende die Huffman-Codierung durchzuführen.

5 Experiment

5.1 Analyse der Ergebnisse mit Hilfe von Metriken

5.1.1 Metriken nennen und vergleichen

5.1.2 Klassifizierungsgenauigkeit

Dies funktioniert nur dann gut, wenn der Datensatz nahezu gleichmäßig auf die Klassen verteilt ist.

Datensatz nicht gleichmäßig verteilen

- Die Bewertung Ihres maschinellen Lernalgorithmus ist ein wesentlicher Bestandteil eines jeden Projekts. Ihr Modell kann Ihnen zufriedenstellende Ergebnisse liefern, wenn es mit einer Metrik bewertet wird, sagen wir `accuracy_score`, aber es kann schlechte Ergebnisse liefern, wenn es mit anderen Metriken wie `logarithmic_loss` oder einer anderen solchen Metrik verglichen wird.

•

5.2 Einfluss der Lernrate

Wie in 3.2.3 gesehen, die Lernrate sagt uns, wie schnell wir ans Ziel kommen. sei η_{opt} bzw. η die optimale bzw. die ausgewählte Lernrate für die Lösung unserer Aufgabe.

Gilt $\eta < \eta_{opt}$, so sind wir sicher, ein lokales oder globales Minimum zu erreichen. Aber die Anzahl der benötigten Iterationen bis zum Minimum steigt offensichtlich an und das Verfahren kann leichter in einem unerwünschten lokalen Minimum stecken bleiben.

$\eta > \eta_{opt}$: Hier wird die Anzahl der Iterationen entweder reduziert oder erhöht. Das Verfahren ist im Allgemeinen nicht stabil, denn es wird über das Minimum ständig hinausgegangen und es ist nicht mehr sicher, zu einem lokalen oder globalen Minimum zu gelangen. Wenn eine hohe Lernrate angewendet wird, kann man versuchen, das Verhalten des CNN zu kontrollieren, indem man die Lernrate nach einer Reihe von Iterationen reduziert, bei denen sich die Metrik (Netzwerkgenauigkeit oder Netzwerkfehler) nicht mehr verbessert hat, oder indem man Optimierer(5.5) verwendet, die im Laufe des Trainings die Lernrate entsprechend der Eingabe und der aktuellen Lernrate anpassen. Der Grund, warum die Lernrate immer reduziert und nicht erhöht wird, sollte daran liegen, dass das CNN irgendwann gut trainiert ist und sollte daher neuen Informationen nicht mehr Gewicht beimessen als dem bereits Gelernten.

[hier kommen noch Bild und Erklärung](#)

5.3 Convolution Layers

Heutzutage werden nur Filter der Größe 1×1 und 3×3 und nicht größer verwendet, obwohl größere Filter globale Informationen effizienter entnehmen und so zu besserer Ergebnisse kommen können. Dafür gibt es viele Gründe. Erstens sind die Berechnungen

mit kleinerer Filter viel schneller als mit größerer Filter. Zweitens kann ein ConvL mit einer großen Filtergröße durch mehrere aufeinanderfolgende ConvLs mit einer kleineren Filtergröße ersetzt werden und dabei werden sogar die Anzahl der benötigten Parameter deutlich reduziert. Sollte man z.B. ein 5×5 Filter auf eine Eingabe von Tiefe d angewendet, so braucht man $5 \times 5 \times d$ Gewichte, um ein Feature-Map zu erzeugen. Nimmt man hingegen zwei aufeinanderfolgende ConvLs mit 3×3 Filtern, so braucht man nur $2(3 \times 3 \times d)$ Gewichte, um ein Feature-Map zu erzeugen. Drittens hat man bei kleinen Filtern nur eine langsame Reduzierung der Bilddimension, was den Entwurf von sehr tiefer CNNs wie *InceptionResNet* ermöglicht. Viertens können kleine Filter sehr lokale Merkmale extrahieren, so dass einfache und komplexe Informationen gleichzeitig erfasst werden. Die Anzahl der extrahierten Features ist enorm und werden gefiltert, wenn man tiefer im Netzwerk geht.

Der Absatz(3.1.2) behandelte mehr über die Faltungsoperation bei einem Eingang mit einem Kanal. Im Folgenden werden verschiedene Arten beschrieben, wie die Faltungsoperation an einem Eingang mit mehreren Kanälen durchgeführt werden kann. Für die Erklärung von Konzepten wird folgendes als Eingabe der Faltungsschicht betrachtet.

$$I = \begin{cases} \text{Filtergröße: } F := (F_w, F_h) & n_f: \text{Anzahl von Filtern} \\ \text{Input_size: } W_{in} \times H_{in} \times D_{in} & P: \text{die Anzahl der Nullauffüllung(Padding)} \end{cases}$$

Standard Convolution

Die Standardfaltung funktioniert genau wie in der Abbildung(24). Für eine Eingabe I erzeugt die Standardfaltung Feature-Maps der Größe $(W_{out}, H_{out}, 1)$, wobei W_{out} und H_{out} wie in der Gleichung 5.1 berechnet werden. Für die Berechnung eines Feature-Map wird ein $(F_w \times F_h \times D_{in})$ -Filter auf den Input angewendet. Dabei wird zuerst für jeden Kanal $I_i := I[:, :, i]$ die Faltungsoperation mit dem Filter $F_i := F[:, :, i]$ durchgeführt und dann werden die Ergebnisse jeder Faltungsoperation aufsummiert, um das Feature-Map zu bilden.

$$\begin{aligned} W_{out} &= \frac{W_{in} - F_w + 2P}{S} + 1 \\ H_{out} &= \frac{H_{in} - F_h + 2P}{S} + 1 \\ D_{out} &= n_f \end{aligned} \tag{5.1}$$

Die Berechnungskosten der Standardfaltung für eine Eingabe I unter der Annahme, dass die Schrittgröße eins und Padding berechnet werden, entspricht:

$$F_w \times F_h \times D_{in} \times n_f \times W_{in} \times H_{in} \tag{5.2}$$

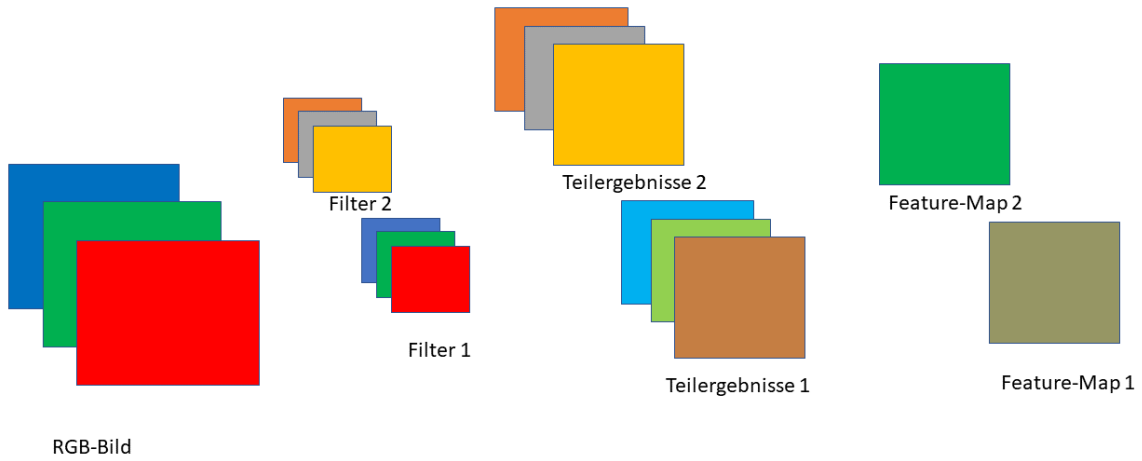


Abbildung 24: Anwendung einer Standardfaltung auf ein Farbbild.

Depthwise Convolution

Genau wie die Standardfaltung funktioniert die *Depthwise Convolution*, außer dass die Teilergebnisse von Standardfaltung schon der Feature-Maps von der *Depthwise Convolution* entsprechen. Festzustellen ist, dass die Anzahl von Feature-Maps immer proportional zur Anzahl der Kanäle des Inputs sein muss. Damit die Anzahl an Feature-Maps der *Depthwise Convolution* mehr als die Anzahl der Kanäle des Inputs ist, muss ein so genannte *Depth_multiplier* α definiert werden. Der *Depth_multiplier* $\in \mathbb{N}_{\geq 1}$ gibt an, wie viele Filter pro Kanal verwendet werden müssen. Die Abbildung(25) zeigt die Anwendung der *Depthwise Convolution* mit einem $\alpha = 2$ auf ein Farbbild. Die *Depthwise Convolution* hat für eine Eingabe I unter der Annahme, dass Schrittgröße eins und Padding berechnet werden, einen rechnerischen Aufwand von:

$$F_w \times F_h \times D_{in} \times W_{in} \times H_{in} \quad (5.3)$$

Wenn man die Gleichungen 5.2 und 5.3 betrachtet, stellt man schnell fest, dass *Depthwise Convolution* im Vergleich zur Standardfaltung äußerst effizient ist. Der Nachteil von *Depthwise Convolution* ist, dass jede generierte Feature-Map nur einen Kanal in Betracht zieht, was bedeutet, dass Kanäle, die nicht so viele wichtige Informationen enthalten, die gleiche Bedeutung erhalten wie Kanäle, die so viele Informationen enthalten.

Pointwise Convolution

Die *Pointwise Convolution* funktioniert genau wie die Standardfaltung, nur dass die Größe der Filter und die Schrittgröße in der Standardfaltung variabel sind und in *Pointwise Convolution* jeweils immer $(1, 1)$. Also jedes erhaltene Teilergebnis ist nur die Skalierung

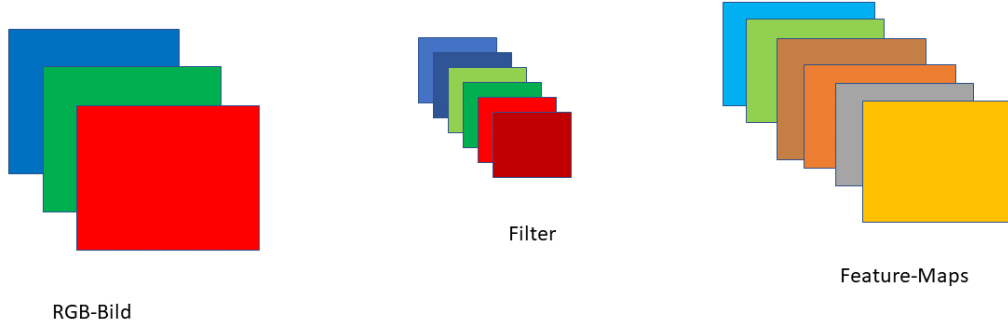


Abbildung 25: Anwendung einer Depthwise Convolution auf ein Farbbild.

eines Eingangskanals. Die *Pointwise Convolution* hat für eine Eingabe I einen rechnerischen Aufwand von:

$$D_{in} \times W_{in} \times H_{in} \times n_f \quad (5.4)$$

Die Hauptidee bei der Verwendung der *Pointwise Convolution* ist es, die Korrelation zwischen den Eingangskanälen zu finden, oder lineare Kombinationen zu lernen und erleichtert die Reduzierung und Erhöhung der Anzahl der Kanäle, wenn es zu viele oder zu wenige Feature-Maps gibt. Es kann jedoch mit *pointwise Convolution* kein Feature extrahieren. **changement de dimension**

Depthwise Separable Convolution Die *Depthwise Separable Convolution* (DSC) besteht aus einer *Depthwise Convolution*, gefolgt von einer *Pointwise Convolution*. Sie kann also Features schneller filtern und kombinieren als die Standardfaltung. Die DSC hat einen rechnerischen Aufwand von:

$$F_w \times F_h \times D_{in} \times W_{in} \times H_{in} + D_{in} \times W_{in} \times H_{in} \times n_f \quad (5.5)$$

Von 5.2 und 5.5 aus sieht man, die DSC weniger Berechnung benötigt als Standardfaltung, und zwar $\frac{1}{n_f} + \frac{1}{F_w \times F_h}$. Da die DSC zu wenig Parameter verwendet, werden die CNN, die die DSC verwenden, als quantisiert bezeichnet. Die DSC ersetzt langsam und sicher die Standardfaltung und das kann man mit *Xception* und *MobileNet* beobachten.

compare accuracy von Depthwise Separable Convolution and Standard Convolution

5.4 Besondere CNNs: Bitte nicht lesen

5.4.1 AlexNet

Wie in 2.4 schon erwähnt, wurde AlexNet von *Krizhevsky et al* [19] im Jahr 2012 entwickelt. AlexNet hat 60 Millionen Parametern und etwas 650.000 Neuronen, besteht aus fünf ConvLs, von denen einige von Max-Pools gefolgt sind, und drei FCLs (siehe Abbildung 26) [19]. Vor AlexNet gab es keine CNNs mit so vielen Schichten, noch interessanter sind die Techniken zur Verbesserung der Leistung von AlexNet :

- **ReLU:** Vor AlexNet waren *tanh* und *sigmoid* die Standard Aktivierungsfunktionen, aber wegen ihrer Sättigung bei hohen oder sehr niedrigen Werten und der Tatsache, dass sie in diesen Bereichen eine Steigung nahe bei null haben, verlangsamten sie stark die Gewichtsanpassungen, was nicht der Fall bei *ReLU* ist, das eine Steigung gleich null nur bei negativen Werten und bei positiven höheren Werte eine Steigung nicht nahe bei null hat. Das Training von CNN wird durch *ReLU* um ein Vielfaches schneller als ein Äquivalent mit *tanh* beschleunigt [19].
- **Überlappendes Pooling:** Die normalen Pools funktioniert wie in 3.1.4 ($pool_size = stride$), aber die Überlappenden verwenden einfach ein Schrittgröße kleiner als das $pool_size$. Nach [19] verbessern die überlappenden Pool die Netzgenauigkeit und macht das Netz gegenüber Overfitting robuster.
- **Dropout und Data Augmentation :** Ein anderer Vorteil von AlexNet ist die Verwendung von Dropout (5.6.1) in FCLs, das sich heute als die beste oder eine der besten Regulierungsmethoden erweist, und der Vermehrung der Daten durch Spiegelung, die die Wahrscheinlichkeit von Overfitting und die Anzahl von Epochen reduziert.

AlexNet kann aufgrund seiner Größe ($\sim 240\text{MB}$) nicht immer auf Systemen mit begrenztem Speicherplatz eingesetzt werden, deshalb ist man seit AlexNet ständig auf der Suche nach neuen Architekturen, die weniger Parameter und Speicherplatz brauchen und gleichzeitig die Ergebnisse auf der Stand der Technik erreichen. Es ist sicherlich in diesem Zusammenhang, dass viele neue effiziente Modelle wie *SqueezeNet*, *Xception* und *MobileNet* entstanden sind.

5.4.2 SqueezeNet

SqueezeNet von [20, Han et al] erreicht die Genauigkeit von AlexNet mit 50x weniger Parametern und einer Größe kleiner als 0.5MB. Die *SqueezeNet* Architektur und einige Variante davon können in Abbildung 27 entnommen werden. Zum Erreichen einer solchen Architektur werden drei Hauptstrategien und ein spezieller Block verwendet. Bezüglich dieser Strategien wurde zuerst die meisten 3×3 Filter durch 1×1 Filter ersetzt, was die Parameteranzahl stark reduziert, denn ein 1×1 Filter nimmt 9x weniger Parameter als ein 3×3 Filter in Anspruch, dann wird die Anzahl der Eingangskanäle auf 3×3 Filter verringern, denn die Tiefe eines Filters in einem ConvL ist gleich die Tiefe vom

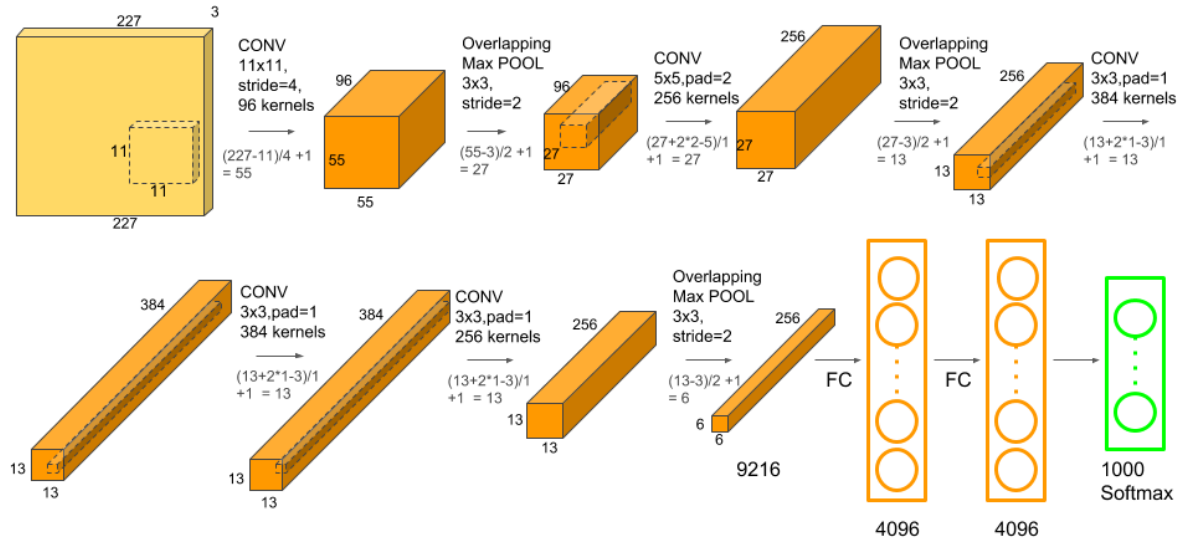


Abbildung 26: AlexNet Architektur source

ConvL Input und ein Input mit einer großen Tiefe wird die Parameteranzahl deutlich zunehmen, obwohl Filter mit kleiner räumlicher Dimension benutzt werden. Die letzte Strategie besteht darin, die Reduzierung der räumlichen Dimension im Netzwerk zu verzögern, so werden Feature-Maps mit größerer Raumaufösung am Ende enthalten. [20, Song Han et al] sind der Meinung, dass größere Feature-Maps zu einer höheren Klassifizierungsgenauigkeit führen.

Um die Ziele der ersten und zweiten Strategien zu erreichen, wird ein neuer Baustein(*fire module*)(FM) eingeführt. Das FM besteht aus einem *Squeeze*-(mit 1×1 Filtern) und einem *Expand*-(mit 1×1 und 3×3 Filtern)-Block(siehe Abbildung28). Der *Squeeze* Block ermöglicht der Reduzierung der Dimensionalität, indem die Anzahl Feature-Maps verringert wird, während die wichtigsten Features erhalten bleiben und die Anwendung von *ReLU* nach dem *Squeeze* Block ermöglicht dem Netzwerk, weitere Features zu lernen.

Was der *Expand* Block angeht, funktioniert der Block mit 1×1 Filtern genauso wie der *Squeeze* Block und der Block mit 3×3 Filtern ermöglicht der Erlangung von Empfangsfelder(*receptive field*), die in der Lage sind, lokale räumliche Informationen zu erfassen, was 1×1 Filter nicht schaffen können.

Wie es in Abbildung 27 zu sehen ist, kann eine Verbindung zwischen zwei Layers über eine oder mehrere Schichte überspringen.Solche Verbindungen werden *Residual connections*(RC) genannt und Netze mit RCs gehören zu der Klasse von Residual Netzwerke (mehr dazu siehe [21]). Wie die Abbildung ?? zeigt, liefert *SqueezeNet* mit Simple RCs als die anderen Varianten

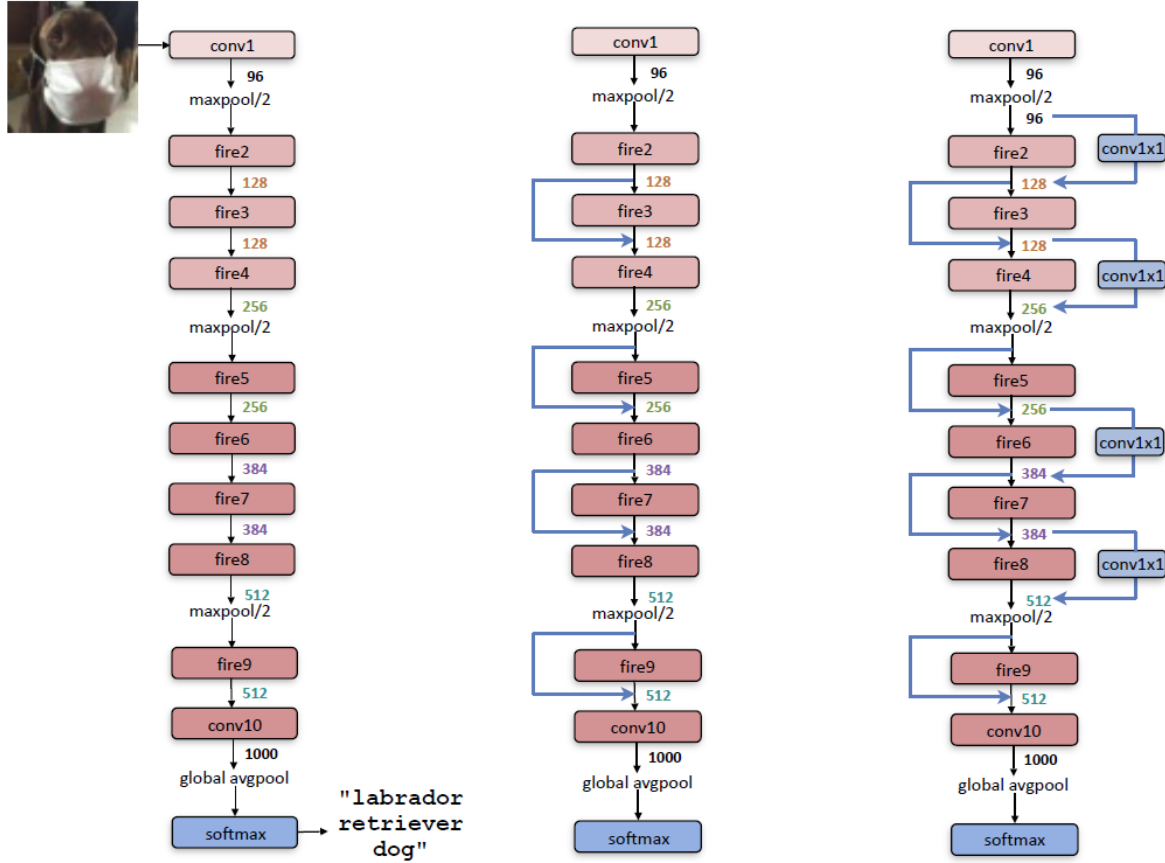


Abbildung 27: SqueezeNet Architekturen [20]

5.4.3 Xception

Xception ([22]) basiert auf *Inception V3* ([23]). Die Inception-Modelle nutzen ein Inception-Modul als Baustein. In einem Inception-Modul werden mehrere Filter unterschiedlicher Größe separat auf eine Eingabe angewendet und die Ergebnisse jeder Faltungsoperation werden zusammengeführt. Die Verwendung von mehreren Filtern unterschiedlicher Größe ermöglicht eine schnellere Feature-Extraktion, da die kleinen Filter sich mit der Extraktion von lokalen Features beschäftigen, während die großen Filter sich mit der Extraktion von globalen Features beschäftigen. Alle Inception-Module werden in *Xception* durch ein *Depthwise Separable Convolution* (DSC)-Layer ersetzt.

Ein DSC besteht aus einer tiefen Faltung (*depthwise convolution*), gefolgt von einer punkweisen Faltung (*pointwise convolution*). Bei der *depthwise convolution* wird die Faltungsoperation für jeden Eingangskanal unabhängig von anderen Kanälen durchgeführt. Die *pointwise convolution* entspricht der 1×1 -Faltung, die die Kombination der Ausgabe von *depthwise convolution* ermöglichen.

Während eine Standardfaltung (3.1.2) Eingänge in einem Schritt filtert und kombiniert, teilt die *depthwise convolution* diese in zwei Schichten auf, eine separate Schicht zur Filterung und eine separate Schicht zur Kombination.

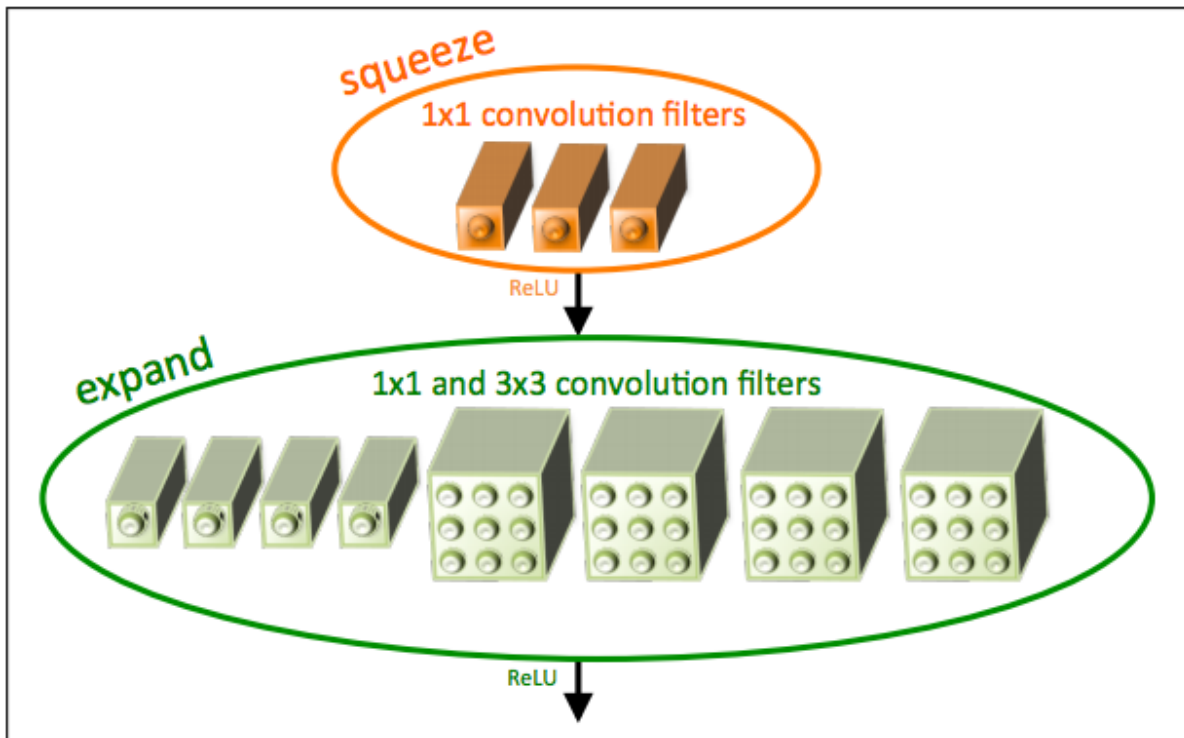


Abbildung 28: fire_module [20]

5.4.4 MobileNet

5.5 Algorithmen zur Optimierung des Gradientenabstiegsverfahren: Optimizer

Die Wahl des Optimierungsalgorithmus für ein CNN kann den Unterschied zwischen guten Ergebnissen in Minuten, Stunden und Tagen ausmachen. Zum besserer Anwendung der Gradientenabstiegsverfahren wurden mehrere Optimierte Lernverfahren entwickelt. Im folgenden wird ein kurzer Einblick über die bekanntesten Lernverfahren (*Optimizer*) gegeben.

Alle heutige Optimizer haben SGD als Vorfahren und der Hauptnachteil von SGD ist , dass es die gleiche Lernrate für die Anpassung aller Netzwerkparameter verwendet und diese Lernrate wird auch während des Trainings nie geändert.

5.5.1 Adaptive Gradient Algorithm (AdaGrad)

AdaGrad bietet während des Netztrainings nicht nur die Möglichkeit, die Lernrate zu verändern, sondern auch für jeden Parameter eine geeignete Lernrate zu finden. Die

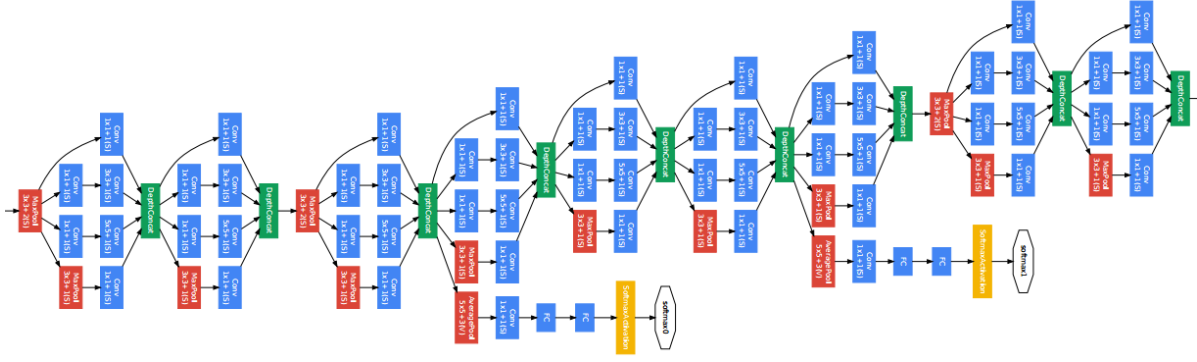


Abbildung 29: Xception Architektur

AdaGrad-Aktualisierungsregel ergibt sich aus der folgenden Formel:

$$\alpha_t = \sum_{i=1}^t (g_{i-1})^2 \quad \theta_{t+1} = \theta_t - \eta_t g_t \quad (5.6)$$

$$\eta_t = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

Voreingestellte Parameter(*KERAS*) : $\alpha_0 = 0.0$ $\eta = 0.001$ $\epsilon = 10^{-7}$

Dabei wird am Trainingsanfang eine Lernrate für jeden Parameter definiert und im Trainingsverlauf separat angepasst. Dieses Verfahren eignet sich gut für spärliche Daten, denn es gibt häufig auftretende Merkmale sehr niedrige Lernraten und seltene Merkmale hohe Lernraten, wobei die Intuition ist, dass jedes Mal, wenn eine seltene Eigenschaft gesehen wird, sollte der Lernende mehr aufpassen. Somit erleichtert die Anpassung das Auffinden und Identifizieren sehr voraussehbarer, aber vergleichsweise seltener Merkmale.[14]. Wie in der Gleichung (5.6) festzustellen, nach einer bestimmten Anzahl von Iterationen haben wir keine Verbesserung der Netzleistung, denn je größer t wird, desto kleiner η_t wird und irgendwann wird η_t so klein, dass $\eta_t g_t$ fast gleich null ist.

5.5.2 Root Mean Square Propagation(RMSProp)

RMSProp wie AdaGrad findet für jeden Parameter eine geeignete Lernrate und zur Anpassung der Netzparameter basiert der RMSProp Optimizer auf den Durchschnitt der aktuellen Größen der Gradienten statt auf der Summe der ersten Moment wie in AdaGrad. Da $E[g^2]_t$ nicht schneller als α_t (5.6) ansteigt, wird die radikal sinkenden Lernraten von Adagrad deutlich verlangsamt. Die Parameteranpassungen richten sich nach der folgenden Gleichung:

$$E[g^2]_t = \alpha E[g^2]_{t-1} + (1 - \alpha) g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t, \quad \epsilon \approx 0 \quad (5.7)$$

Der RMSProp funktioniert besser bei Online- und nicht-stationären Problemen.

Table 1. MobileNet Body Architecture

Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5×	Conv dw / s1	$3 \times 3 \times 512$ dw
	Conv / s1	$1 \times 1 \times 512 \times 512$
		$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool 7×7	$7 \times 7 \times 1024$
FC / s1	1024×1000	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$

Abbildung 30: MobileNet Architektur

5.5.3 Adaptive Moment Estimation(Adam)

Der Adam[15] Optimizer ist auch ein adaptiver Algorithmus, der die ersten und zweiten Momente der Gradienten schätzt, um individuelle adaptive Lernraten für verschiedene Parameter zu berechnen. Adam weist die Hauptvorteile von AdaGrad, das mit spärlichen Gradienten gut funktioniert, und RMSProp, das einige Probleme von AdaGrad löst und das für nicht-konvexe Optimierung geeignet ist, auf. Wie die Parameteranpassung von Adam Optimizer genau funktioniert, ergibt sich aus der folgenden Gleichung:

$$\begin{aligned}
 m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t, & \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\
 v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2, & \hat{v}_t &= \frac{v_t}{1 - \beta_2^t} \\
 \theta_{t+1} &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t
 \end{aligned} \tag{5.8}$$

Voreingestellte Parameter(KERAS) : $\beta_1: 0.9$ $\beta_2: 0.999$
 $\eta: 0.001$ $\epsilon: 10^{-7}$

Zu weiteren Vorteile der Nutzung von Adam gehört auch seine Einfachheit zur Implementierung, effizienter Nutzung der Speicherplatz und seine Invarianz zur diagonalen Neuskalierung der Gradienten.

Kleines Experiment

5.6 Problem beim Training von Convolutional neuronale Netzwerke

5.6.1 Overfitting

Wenn ein von der Maschine gelerntes Modell zu gut auf die Trainingsdaten abgestimmt ist und sehr schlechte Vorhersagen über Daten macht, die es bisher nicht gesehen hat, wird gesagt, dass das Modell an Überanpassung(Overfitting) leidet, anders gesagt, das Modell war nicht in der Lage, die relevanten Merkmale aus den Trainingsdaten zu verallgemeinern, sondern die ganzen Trainingsdaten auswendig zu lernen. Die irrelevanten Informationen aus den Trainingsdaten können z.B die Position des Tellers in einem Bild sein, wenn man Essen klassifizieren sollte. Im folgenden werden einige Mittels vorgestellt, um mit Overfitting umzugehen.

Strategie gegen Overfitting

Data Augmentation Ein großer Datensatz ist entscheidend für die Leistung tiefer neuronaler Netze. Dass ein Datensatz groß oder ausreichend für das Training eines NNs ist, hängt nur von der Größe des NNs ab und da die NNs, die die besten Leistungen aufweisen, Millionen von Parametern haben, ist fast unmöglich für jedes Problem von Maschine Lernen ausreichende Daten zu finden. Anstatt immer neue Daten zur Verbesserung der Netzleistung zu sammeln, können wir die Leistung des Modells verbessern, indem wir neue Daten von den bereits vorhandenen Daten aus erzeugen.

Die populären Techniken oder Transformationen zur Vermehrung des Datensatzes sind die horizontalen oder vertikalen Spiegelungen, Drehungen, Skalierungen, Zuschneiden, Parallelverschiebungen und die Gauß'sches Rauschen. Für diese Arbeit habe ich zwei Ansätze im Gebrauch gehabt, um den Datensatz zu erhöhen: Der erste Ansatz besteht darin vor dem Training neue Daten zu erzeugen. Dabei werden die oben erwähnten Techniken vor dem Training angewendet, um zum Trainingszeitpunkt und zur Testzeit einen großen Datensatz zu haben und die originalen Daten werden zur Validierung verwendet. Der zweite Ansatz besteht darin, zum Trainingszeitpunkt und zur Testzeit die neuen Daten zu erzeugen. Hier haben wir keinen Datensatz größer als den originalen, aber die Daten, die ins Netzwerk eingespeist werden, ändern sich ständig. Angenommen wir die Möglichkeit haben, alle diese Transformationen durchzuführen, dann kann es vorkommen, dass eine Photo während der ersten Epoche horizontal gespiegelt wird und während der zweiten um zwanzig Grad gedreht wird, umso weiter. Zur Implementierung des zweiten Ansatzes bietet *KERAS* Framework die Funktion *ImageDataGeneerator*. Das interessanteste an *ImageDataGeneerator* ist, dass es mehrere Transformationen gleichzeitig anwenden(siehe Abbildung 31).

Abbildung 31: Anwendung von *ImageDataAugmentation*

Je mehr Daten verfügbar sind, desto effektiver können die CNNs sein. Es ist also mehr als wichtig über eine große Datenmenge zu verfügen. Leider können die gesammelten Datensätze nicht alle mögliche Szenarios des reellen Lebens abdecken, deshalb ist es auch bedeutend, CNN mit zusätzlichen synthetisch modifizierten Daten zu trainieren. Die CNNs funktionieren glücklicherweise besser oder immer gut, solange nützliche Daten durch das Modell aus dem ursprünglichen Datensatz extrahiert werden können, selbst wenn die erzeugten Daten von geringerer Qualität sind.

Dropout Künstliche Neurone sind von biologischen Neuronen inspiriert, aber die Beiden unterscheidet sich sehr voneinander und einer der wichtigen Unterschiede ist, dass biologische Neuronen unvollkommene Maschinen sind, die sehr oft nicht richtig funktioniert und das ist a priori nie den Fall bei künstlichen Neuronen. Wir könnten also glauben, dass KNN die biologische übertreffen könnten. Es sei denn, dass diese Funktionsstörung von biologischen Neuronen nicht eine Schwäche ist, sondern eher eine Stärke ist. Eine der verblüffenden Entdeckungen in Künstliche Intelligenz (KI) Bereich ist, dass es wünschenswert ist, künstliche Neuronen von Zeit zu Zeit zu Fehlfunktionen zu bringen[1]. [Jetzt können wir uns fragen, wie Dysfunktion von Neuronen die Performances CNNs verbessern kann.](#) Die zufällige Hinzufügen von Dysfunktionen in einer Schicht der CNN wird *Dropout* benannt und wurde von [3, Geoffrey E. et al] eingeführt.

Funktionsweise von Dropout

Genauer gesagt, Dropout bezeichnet die zeitliche zufällige Ausschaltung von Neuronen (versteckt

und sichtbar) in einem NN [4]. Wie die Abbildung 32 zeigt, wenn ein Neuron zufällig aus dem NN entfernt wird, werden auch all seine ein- und ausgehenden Verbindungen entfernt. In einer Dropout-Schicht wird ein Neuron N unabhängig von anderen Neuronen mit einer Wahrscheinlichkeit p zurückgehalten, d.h. N wird mit einer Wahrscheinlichkeit von p nicht am Ergebnis der Schicht teilnehmen. Während der Testphase werden alle Verbindungen zurückgesetzt, die während des Trainings gelöscht wurden und die ausgehenden Verbindungen gelöschter Neurone mit p multipliziert.

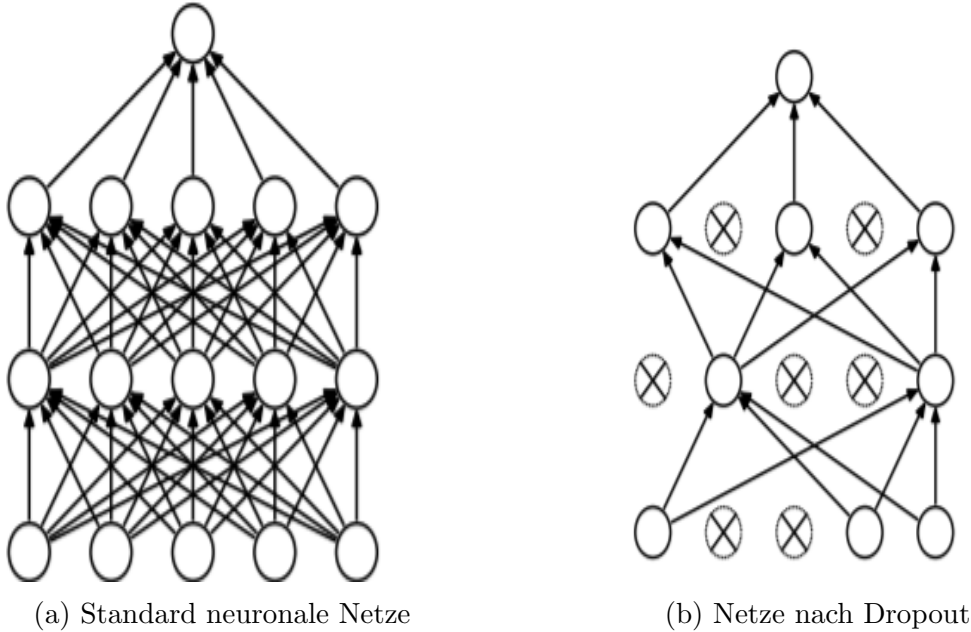


Abbildung 32: Neuronales Netz mit Dropout ausgestattet [4].

Verhinderung der Koadaptationen zwischen Neuronen

Während des Trainings können mehrere Neurone zur Minimierung der Fehlerfunktion so gut zusammenarbeiten, dass die Erkennung bestimmter Merkmale ohne diese Zusammenarbeit nicht mehr möglich ist, aber solche komplexe Koadaptationen können zu einer Überanpassung führen, denn diese komplexe Koadaptationen existieren nicht immer in Testdaten. Da die am Training teilnehmenden Neuronen nach dem Zufallsprinzip nach jeder Epoche ausgewählt werden, haben wir für jedes Training ein neues Modell, was die Neuronen zur Zusammenarbeit zwingt, ohne jedoch voneinander abhängig zu sein, anders gesagt, wird jedes Neuron unabhängig von anderen Neuronen die Muster korrekt lernen können.

Automatische Erhöhung von Training Daten und Regelung

Noch dazu führt die Ausschaltung von Neuronen, wie es in Abbildung 33 angezeigt ist, zu einer automatische Erzeugung neuer Trainingsdaten. Die verwendeten Daten in ausgedünnten Modellen sind also nur eine Abstraktion von echten Daten bzw. Rauschdaten

und da wir für ein Netz mit n versteckten Einheiten, von denen jede fallen gelassen werden kann, 2^n mögliche Modelle haben, haben wir 2^n mögliche Abstraktion von unseren Daten und das sollte einer der Gründe sein, warum Dropout effektiver als andere rechnerisch kostengünstige Regler ist [4] und warum die Trainingszeit von NNs mit Dropout mindestens verdoppelt wird.

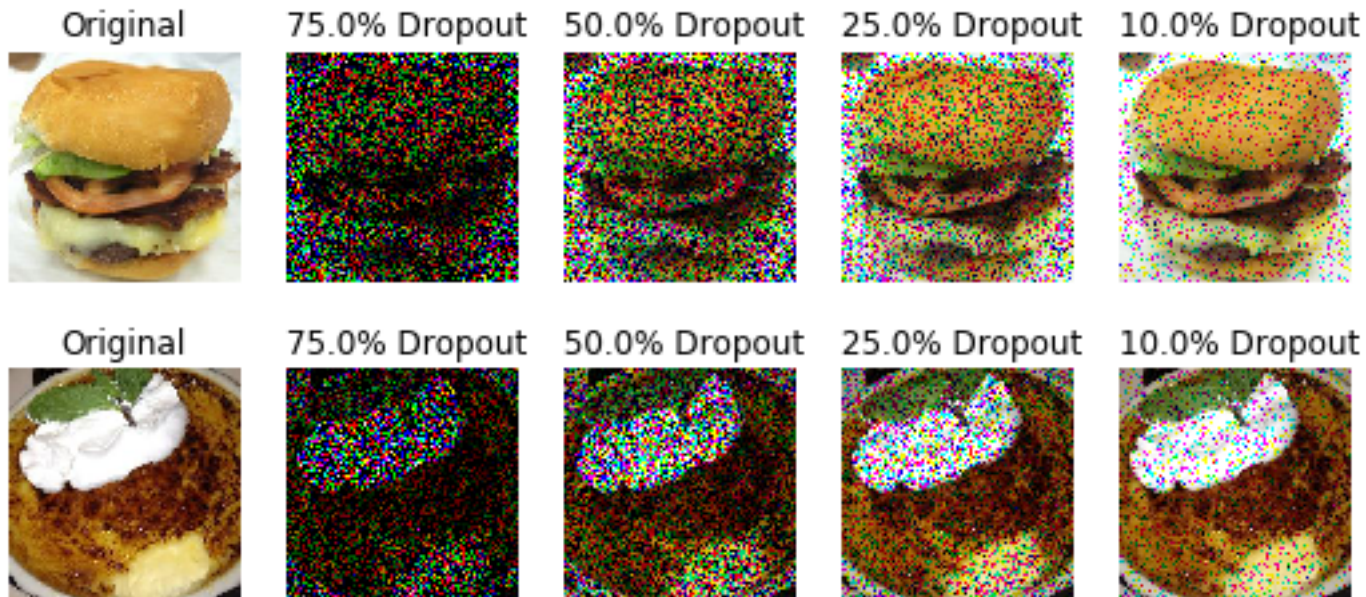


Abbildung 33: Erhöhung des Trainingsdaten durch Dropout

Da heutige CNNs Million von Neurons haben, wäre es unmöglich alle mögliche ausgedünnte Netzwerke zu trainieren, deshalb ist das Modell, das am Ende des Trainings erhalten wird, nur eine durchschnittliche Approximation aller mögliche Modelle, was schon gut, denn es gibt schlechte und gute Modelle.

Batch-Normalisierung Das Training tiefer neuronaler Netze ist sehr kompliziert und ein Grund dafür ist zum Beispiel die Tatsache, dass die Parameter einer Schicht während des Trainings tiefer neuronaler Netze immer unter der Annahme, dass sich die Parameter anderer Schichten nicht ändern, aktualisiert werden und da alle Schichten während des Updates geändert werden, verfolgt das Optimierungsverfahren ein Minimum, das sich ständig bewegt. Ein anderer Grund dafür ist die ständigen Veränderungen im Laufe des Trainings in die Verteilung des Netzeinputs, diese Veränderung wird von [12] als interne kovariante Verschiebung (*Internal Covariate Shift*) genannt. Zur Lösung dieser Probleme schlagen *LeCunn et al*[?] vor dem Training das Netzeinput zu normalisieren. Aber dieser Ansatz bringt nicht so viel, wenn das NN wirklich tief ist, denn nur der Netzeinput profitiert von der Normalisierung und die kleinen Veränderungen in versteckte Schichten werden sich immer mehr verstärken, je tiefer man das Netz durchläuft. Mit der Ausbreitung tiefer NNs dehnt Batch-Normalisierung(BN)[12] diese Idee der Datennormalisierung auf versteckte Schichten tiefer NNs aus. Bei der BN werden die Eingaben in einem

5 Experiment

Netzwerk standardisiert, die entweder auf die Aktivierungen einer vorherigen Schicht oder auf direkte Eingaben angewendet wird, so standardisiert, dass der Mittelwert in der Nähe von null liegt und die Standardabweichung in der Nähe von eins liegt. Die BN wird über Mini-Batches und nicht über den gesamten Trainingssatz durchgeführt, daher enthalten wir nur Näherungen an tatsächliche Werte der Standardabweichung und des Mittelwerts über das Trainingssatzes, aber wir gewinnen an Geschwindigkeit und an Speicherplatzverbrauch. Die Gleichung (5.9) gibt die formale Beschreibung des BN Algorithmus an.

Batch-Normalisierungstransformation, angewendet auf Aktivierung x über einen Mini-Batch

Input: Werte von x über einer Mini-Batch: $B = \{x_{1...m}\}$

Lernbare Parameter β, γ

Output: $\{y_i = BN_{\beta, \gamma}(x_i)\}$

$$\text{Mini-Batch Mittelwert : } \mu_\beta = \frac{1}{m} \sum_{i=1}^m x_i \quad (5.9a)$$

$$\text{Mini-Batch Standardabweichung : } \sigma_\beta^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_\beta)^2 \quad (5.9b)$$

$$\text{Normalisierung: } \hat{x}_i = \frac{x_i - \mu_\beta}{\sqrt{\sigma_\beta^2 + \epsilon}} \quad (5.9c)$$

$$\text{Skalierung und Verschiebung : } y_i = \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta} \quad (5.9d)$$

Wenn $\gamma = \sqrt{\sigma_\beta^2 + \epsilon}$ und $\beta = \mu_\beta$, bekommen wir die gleiche Verteilung wie vor der Batch-Normalisierung, d.h die Eingabe war also schon normalisiert. Interessanterweise kann das Netz während des Trainings eine bessere Verteilung als die erwünschte finden, denn γ und β sind lernbare Parameter.

Durch die BN kann zum einen eine hohe Lernrate verwendet, was in tiefer NNs ohne BN dazu führen kann, dass die Gradienten explodieren oder verschwinden und in schlechten lokalen Minima stecken bleiben. Die Verwendung einer höheren Lernrate ermöglicht eine schnellere Konvergenz. Zum anderen wird die interne kovariante Verschiebung geringer, was das Training beschleunigt, in einigen Fällen durch Halbierung der Epochen oder besser. Noch dazu wird das Netz durch die BN in gewissem Maße reguliert, daher wird die Verwendung von Dropout bzw. Regulierungstechnik reduziert oder sogar überflüssig und somit eine Verbesserung der Verallgemeinerungsgenauigkeit.

5.6.2 Dataset

6 Literatur

Literatur

- [1] P. Kerlirzin, and F. Vallet: Robustness in Multilayer Perceptrons
- [2] Md. Zahangir Alom, Tarek M. Taha, Christopher Yakopcic, Stefan Westberg, Mahmudul Hasan, Brian C. Van Esesn, Abdul A. S. Awwal und Vijayan K. Asari The History Began from AlexNet: A Comprehensive Survey on Deep Learning Approaches
- [3] Geoffrey E. Hinton and Nitish Srivastava and Alex Krizhevsky and Ilya Sutskever and Ruslan R. Salakhutdinov: Improving neural networks by preventing co-adaptation of feature detectors
- [4] Srivastava, Hinton, Krizhevsky, Sutskever and Salakhutdinov: Dropout: A Simple Way to Prevent Neural Networks from Overfitting
- [5] Ian Goodfellow, Yoshua Bengio, Aaron Courville: Adaptive Computation and Machine Learning series Page 342
- [6] Song Han, Huizi Mao, William J. Dally Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding
- [7] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, Hans Peter Graf Pruning Filters for Efficient ConvNets
- [8] Franco Manessi, Alessandro Rozza, Simone Bianco, Paolo Napoletano, Raimondo Schettini Automated Pruning for Deep Neural Network Compression
- [9] Emily Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun and Rob Fergus Exploiting Linear Structure Within Convolutional Networks for Efficient Evaluation
- [10] Pavel Golik, Patrick Doetsch, Hermann Ney Cross-Entropy vs. Squared Error Training: a Theoretical and Experimental Comparison
- [11] Neuronale Netze: Eine Einführung
- [12] Sergey Ioffe, Christian Szegedy Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift
- [13] Wikipedia: Learning Rate
- [14] John Duchi, Elad Hazan, Yoram Singer: Adaptive Subgradient Methods for Online Learning and Stochastic Optimization

- [15] Diederik P. Kingma, Jimmy Ba: Adam: A Method for Stochastic Optimization
- [16] Compressing Models:Quantization
- [17] Yoni Choukroun, Eli Kravchik, Fan Yang, Pavel Kisilev: Low-bit Quantization of Neural Networks for Efficient Inference
- [18] wikipedia: Artificial neuron
- [19] Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton ImageNet Classification with Deep Convolutional Neural Networks
- [20] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, Kurt Keutzer: SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun: Deep Residual Learning for Image Recognition
- [22] François Chollet Xception: Deep Learning with Depthwise Separable Convolutions
- [23] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, Zbigniew Wojna Rethinking the Inception Architecture for Computer Vision
- [24] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, Hartwig Adam MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications

