

Deep Learning

6. Juli 2015

Inhaltsverzeichnis

1	Einleitung	3
2	Grundlagen	4
2.1	Hopfield Netze	4
2.2	Boltzmann Maschinen	4
2.3	Eingeschränkte Boltzmann Maschinen	7
2.4	Kontrastive Divergenz	9
3	Deep-Belief Netze	10
3.1	Greedy Algorithmus zum Trainieren von Deep-Belief Netzen	11
3.2	Deep-Belief Netze zur Klassifikation	13
3.3	Implementation	14
4	Versuch	17
4.1	Achsensymmetrie und Rotationssymmetrie	17
4.2	MNIST Datenbank	17
4.3	Versuchsdurchführung	17
4.3.1	Achsensymmetrie und Rotationssymmetrie	18
4.3.2	MNIST Datenbank	19
4.3.3	MNIST Datenbank	19
5	Mögliche Verbesserung von Deep-Belief Netzen	19
5.1	Probleme	19
5.2	Einfache Verbesserung	19

1 Einleitung

Blümchenwiese wie toll und wichtig das alles ist.

2 Grundlagen

Als Grundlage für Deep-Belief Netze dienen eingeschränkte Boltzmann Maschinen. Diese sind eng mit Hopfield Netzen verwandt und der Weg von Hopfield Netzen zu eingeschränkten Boltzmann Maschinen wird in diesem Kapitel vorgestellt. Hopfield Netze und allgemeine Boltzmann Maschinen werden hierbei jedoch nur kurz erläutert und am Ende des Kapitels werden eingeschränkte Boltzmann Maschinen ausführlich erklärt.

2.1 Hopfield Netze

In diesen Netzen ist jedes Neuron mit allen anderen Neuronen außer sich selber verbunden, wie in Abbildung 1 dargestellt ist. Die Neuronen in einem Hopfield Netz sind binär, das bedeutet, die Ausgabe ist 1, wenn ein Schwellwert θ überschritten wird, und ansonsten 0.

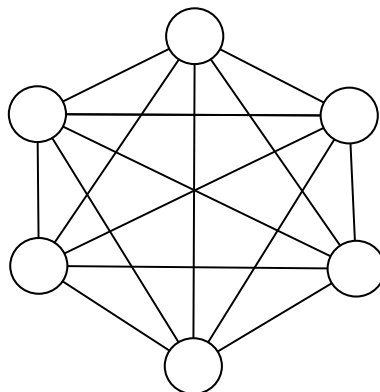


Abbildung 1: Modell eines Hopfield Netzes. Alle Neuronen besitzen eine Ein- und Ausgabe und sind über Gewichte jeweils mit allen anderen Neuronen verbunden

Die Gewichte in einem Hopfield Netz sind symmetrisch, also gilt $w_{ij} = w_{ji}$. Hopfield Netzen wird eine Energie zugewiesen, die beim Lernen des Netzes minimiert werden soll [7]:

$$E = -\frac{1}{2} \sum_{i \neq j} w_{ij} s_i s_j + \sum_i \theta_i s_i \quad (1)$$

s_i bzw. s_j sind hierbei die Zustände der Neuronen. Sobald die Energie unverändert bleibt, wurde ein stabiler Zustand erreicht und man befindet sich in einem lokalen Minimum der Energielandschaft.

2.2 Boltzmann Maschinen

Boltzmann Maschinen sind ähnlich wie Hopfield Netze. Ein wesentlicher Unterschied besteht darin, dass die Neuronen stochastisch binär sind, also nicht mehr feuern, wenn ein Schwellwert überschritten ist, sondern mit einer bestimmten Wahrscheinlichkeit eine 1 als Ausgabe haben. Dies wird in der Praxis so realisiert, dass ein Wurf zwischen 0 und 1 größer

als die Wahrscheinlichkeit sein muss, um das Neuron zu aktivieren. Boltzmann Maschinen wird auch eine Energie zugewiesen, die ähnlich der von Hopfield Netzen ist [1]:

$$E = - \left(\sum_{i < j} w_{ij} s_i s_j + \sum_i \theta_i s_i \right) \quad (2)$$

Desweiteren werden bei Boltzmann Maschinen die Neuronen in sichtbare und in versteckte Neuronen unterteilt. Die sichtbaren Neuronen haben die Möglichkeit eine Eingabe von außen zu bekommen, die versteckten Neuronen nicht. Der Aufbau ist in Abb. 2 dargestellt.

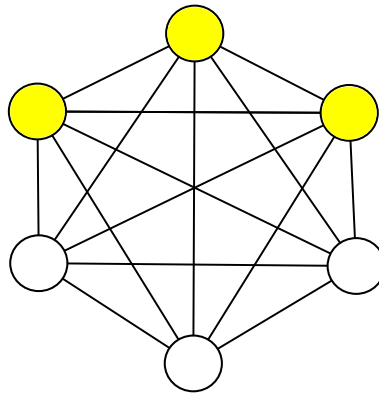


Abbildung 2: Modell einer Boltzmann Maschine. Gelbe Kreise stellen versteckte Neuronen dar, weiße Kreise sind sichtbare Neuronen.

Die Differenz in der globalen Energie abhängig vom Aktivierungszustand eines Neuron lässt sich als ΔE_i beschreiben:

$$\Delta E_i = \sum_j w_{ij} s_j + \theta_i \quad (3)$$

Dies lässt sich als Differenz der Energie der beiden möglichen Zustände eines Neurons schreiben:

$$\Delta E_i = E_{i=\text{off}} - E_{i=\text{on}} \quad (4)$$

Nun kann man mit Hilfe des Boltzmann Faktors, die Zustände mit ihren relativen Wahrscheinlichkeiten substituieren. Der Boltzmann Faktor ist eine Eigenschaft der Boltzmann Verteilung dass die Energie eines Zustands proportional zur negativen logarithmischen Wahrscheinlichkeit des Zustandes ist:

$$\Delta E_i = -k_B T \ln(p_{i=\text{off}}) - (-k_B T \ln(p_{i=\text{on}})) \quad (5)$$

Hierbei ist k_B die Boltzmann Konstante und wird in die Temperatur T absorbiert. Unter der Annahme, dass die Wahrscheinlichkeiten für an und aus summiert 1 ergeben müssen kann man die Formel umstellen:

$$\begin{aligned}
\frac{\Delta E_i}{T} &= \ln(p_{i=\text{on}}) - \ln(p_{i=\text{off}}) \\
\frac{\Delta E_i}{T} &= \ln(p_{i=\text{on}}) - \ln(1 - p_{i=\text{on}}) \\
\frac{\Delta E_i}{T} &= \ln\left(\frac{p_{i=\text{on}}}{1 - p_{i=\text{on}}}\right) \\
\frac{\Delta E_i}{T} &= \ln\left(\frac{1 - p_{i=\text{on}}}{p_{i=\text{on}}}\right) \\
-\frac{\Delta E_i}{T} &= \ln\left(\frac{1}{p_{i=\text{on}}} - 1\right) \\
\exp\left(-\frac{\Delta E_i}{T}\right) &= \frac{1}{p_{i=\text{on}}} - 1
\end{aligned} \tag{6}$$

Nun kann man die Wahrscheinlichkeit $p_{i=\text{on}}$ dass Neuron i aktiviert ist ausrechnen:

$$p_{i=\text{on}} = \frac{1}{1 + \exp(-\frac{\Delta E_i}{T})} \tag{7}$$

Hierbei ist T die Temperatur des Systems. Diese herleitung bietet die Grundlage für die Verwendung der logistischen Funktion, die zum Beispiel bei eingeschränkten Boltzmann Maschinen verwendet wird ¹.

Das Training von Boltzmann Maschinen läuft in zwei Phasen ab: Die erste Phase wird auch als positive Phase bezeichnet und besteht darin, dass man die binären Zustände der sichtbaren Neuronen vorgibt und die Maschine laufen lässt bis sich ein Equilibriumszustand einstellt, sich also die Aktivierungswahrscheinlichkeiten der einzelnen Neuronen nicht mehr ändern.

Danach beginnt man mit der negativen Phase. In dieser Phase lässt man die Maschine ohne eine Eingabe von außen laufen bis diese wieder in einem Equilibriumszustand ist. In den Phasen ergeben sich die Verteilungen $P^+(V)$ und $P^-(V)$ über die sichtbaren Neuronen. Das Ziel ist es, die Differenz zwischen beiden Verteilungen zu minimieren. Um die Differenz der beiden Verteilungen zu berechnen, verwendet man die Kullback-Leibler Divergenz [8]:

$$G = \sum_v P^+(v) \ln\left(\frac{P^+(v)}{P^-(v)}\right) \tag{8}$$

Diese ist immer positiv und nur dann 0, wenn beide Verteilungen genau gleich sind. Daraus ergibt sich eine einfache Lernregel:

$$\frac{\partial G}{\partial w_{ij}} = -\frac{1}{R} [p_{ij}^+ - p_{ij}^-] \tag{9}$$

¹https://en.wikipedia.org/w/index.php?title=Boltzmann_machine&oldid=665335516

Hierbei ist R die Lernrate, und p_{ij} die Wahrscheinlichkeit das Neuron i und j in der positiven, bzw. negativen Phase jeweils aktiviert sind. Theoretisch wären Boltzmann Maschinen sehr gut in der Lage Modelle von Daten zu trainieren und zu erkennen, allerdings sind diese in der Praxis recht untauglich, da sie bei großen Maschinen zu lange brauchen um trainiert zu werden. Dieses Problem wird mit eingeschränkten Boltzmann Maschinen gelöst.

2.3 Eingeschränkte Boltzmann Maschinen

Bei eingeschränkten Boltzmann Maschinen werden die sichtbaren und versteckten Neuronen als Ebenen betrachtet. Die Ebenen sind symmetrisch miteinander verbunden, jedoch existieren zwischen den Neuronen auf einer Ebene keine Verbindungen wie in Abbildung 3 dargestellt.

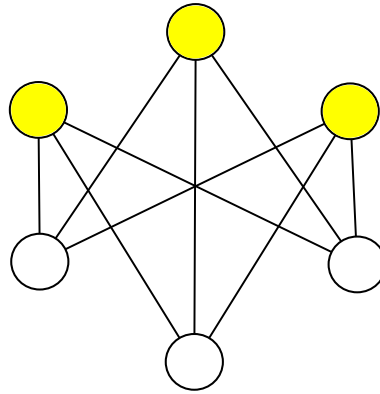


Abbildung 3: Modell einer eingeschränkten Boltzmann Maschine. Oben sind die versteckten Knoten und unten die nach außen sichtbaren Knoten.

Auch eingeschränkte Boltzmann Maschinen benutzen stochastisch binäre Neuronen wie normale Boltzmann Maschinen.

Als Eingabewerte in so ein Netzwerk sind reelle Werte möglich. Diese würden dann als Aktivierungswahrscheinlichkeiten der sichtbaren Neuronen interpretiert und man müsste dann noch die binären Zustände der sichtbaren Neuronen bestimmen. Im Folgenden wird davon ausgegangen, dass binäre Eingabevektoren vorliegen. Eine Konfiguration (\vec{v}, \vec{h}) aus sichtbaren und unsichtbaren Neuronen wird auch bei eingeschränkten Boltzmann Maschinen eine Energie zugewiesen:

$$E(\vec{v}, \vec{h}) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i,j} v_i h_j w_{ij} \quad (10)$$

Hierbei sind v_i und h_j die binären Zustände des sichtbaren Knotens i und des versteckten Knotens j . a_i und b_j sind die Bias der entsprechenden Knoten und w_{ij} ist das Gewicht zwischen den beiden Knoten. Mit Hilfe dieser Energiefunktion weist das Netz jedem möglichen Paar von sichtbaren und versteckten Knoten eine Wahrscheinlichkeit zu:

$$p(\vec{v}, \vec{h}) = \frac{1}{Z} e^{-E(\vec{v}, \vec{h})} \quad (11)$$

Die Zustandssumme (partition function) Z ergibt sich aus der Summe über alle möglichen Werte aus sichtbaren und versteckten Knoten und stellt einen Normierungsfaktor dar, damit die Gesamtwahrscheinlichkeit 1 ergibt:

$$Z = \sum_{\vec{v}, \vec{h}} e^{-E(\vec{v}, \vec{h})} \quad (12)$$

Die Wahrscheinlichkeit, dass sich das Netzwerk einem Eingabebild anpasst, ergibt sich aus der Summe über alle versteckten Vektoren [2].

$$p(\vec{v}) = \frac{1}{Z} \sum_{\vec{h}} e^{-E(\vec{v}, \vec{h})} \quad (13)$$

Die Wahrscheinlichkeit, dass sich das Netz einem Trainingsbild anpasst, kann erhöht werden, indem man die Energie der anderen Bilder erhöht. Vor allem für Bilder mit einer niedrigen Energie ist dies wichtig, da diese einen hohen Beitrag zur Zustandssumme haben. Die Ableitung der Gewichte der logarithmischen Wahrscheinlichkeiten von Trainingsdaten ergibt sich wie folgt:

$$\frac{\partial \log p(\vec{v})}{\partial w_{ij}} = \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model} \quad (14)$$

Hierbei ist $\langle \dots \rangle$ der Erwartungswert der Verteilung über die Eingabedaten und wird als *data* bezeichnet. Die vom Netz rekonstruierte Verteilung wird als *model* bezeichnet. Dies entspricht vom Prinzip her der positiven und der negativen Phase von normalen Boltzmann Maschinen. Auch bei eingeschränkten Boltzmann Maschinen ist das Ziel, dass die Differenz der beiden Verteilungen möglichst gering ist. Sollte diese 0 betragen ist das Modell der Eingabedaten perfekt gelernt, dies geschieht in der Praxis jedoch äußerst selten. Der Erwartungswert ist das Mittel der Wahrscheinlichkeitsverteilung und kann durch das Ziehen von ausreichend vielen Knotenpaaren approximiert werden. Dies führt zu einer einfachen Lernregel für den stochastisch stärksten Aufstieg in der logarithmischen Wahrscheinlichkeit der Trainingsdaten:

$$\Delta w_{ij} = \epsilon (\langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model}) \quad (15)$$

ϵ ist die Lernrate und kann zur Steuerung, wie stark die Gewichte verändert werden, frei gewählt werden. In der Praxis hat eine zu hohe Lernrate zur Folge, dass die Gewichte zu groß werden und nichts gelernt wird. Eine zu kleine Lernrate führt dazu, dass das Training sehr langsam ist, jedoch empfiehlt es sich für ein besseres Ergebnis gegen Ende des Trainings die Lernrate zu verkleinern [?]. Da die versteckten Knoten innerhalb einer RBM nicht miteinander verbunden sind, sind diese statistisch voneinander unabhängig. Dies führt dazu, dass man sehr einfach eine Stichprobe ohne Bias für $\langle v_i h_j \rangle_{data}$ erhalten

kann. Wenn ein zufällig ausgewähltes Trainingsbild \vec{v} gegeben ist, wird der binäre Zustand h_j jedes versteckten Knotens j mit folgender Wahrscheinlichkeit auf 1 gesetzt:

$$p(h_j = 1|\vec{v}) = \sigma(b_j + \sum_i v_i w_{ij}) \quad (16)$$

$\sigma(x)$ ist die logistische sigmoide Funktion $1/(1 + e^{-x})$. $v_i h_j$ ist dann eine Stichprobe der Verteilung ohne Bias a_i , da die einzelnen versteckten Knoten nicht miteinander verbunden und nur abhängig von den Eingabedaten sind.

Da auch die sichtbaren Knoten keine Verbindungen untereinander haben und damit voneinander unabhängig sind, ist es genau so einfach eine Stichprobe für den Zustand eines sichtbaren Knotens zu erhalten, wenn ein versteckter Vektor \vec{h} gegeben ist:

$$p(v_i = 1|\vec{h}) = \sigma(a_i + \sum_j h_j w_{ij}) \quad (17)$$

Eine Stichprobe ohne Bias b_j für $\langle v_i h_j \rangle_{model}$ zu erhalten, ist aufwändiger und kann über abwechselndes “Gibbs Sampling” über einen langen Zeitraum berechnet werden. Beim Gibbs Sampling beginnt man mit einem zufälligen Trainingsvektor und hört auf, wenn ein Equilibriumszustand erreicht wird. Eine Iteration des Gibbs sampling besteht daraus, dass man zuerst parallel die Zustände aller versteckten Knoten mit Gleichung 16 berechnet und im Anschluss mit Gleichung 17 die sichtbaren Knoten neu berechnet, siehe dazu Abbildung 4

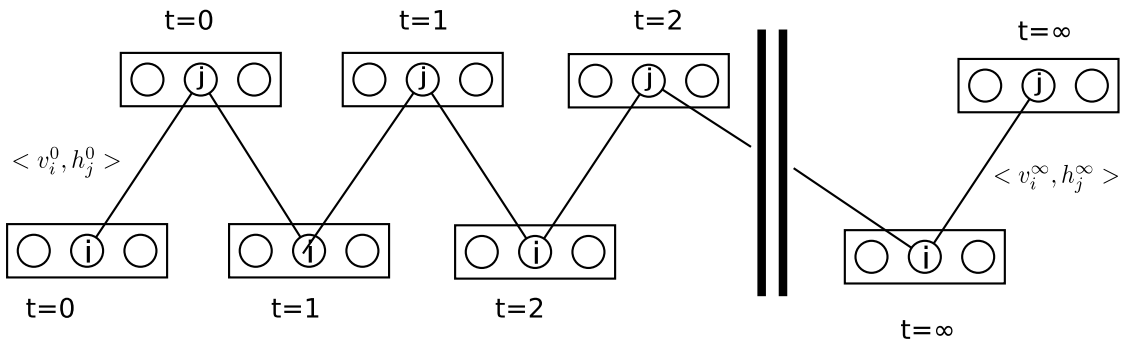


Abbildung 4: Markov-Kette mit Gibbs Sampling. Diese wird mit einem Trainingsvektor initialisiert. Die sichtbaren und unsichtbaren Neuronen werden abwechselnd neu berechnet bis ein Equilibriumszustand erreicht wird.

2.4 Kontrastive Divergenz

Ein schnelleres Verfahren als das Gibbs sampling besteht daraus, die Zustände der sichtbaren Knoten mit einem Trainingsvektor zu belegen. Anschließend werden die versteckten Knoten mit Gleichung 16 berechnet. Nachdem die binären Zustände der versteckten Knoten gefunden sind, wird eine “Rekonstruktion” angefertigt indem man jeden sichtbaren

Knoten mit einer Wahrscheinlichkeit aus Gleichung 17 auf 1 setzt. Die Änderung in einem einzelnen Gewicht wird dadurch zu:

$$\Delta w_{ij} = \epsilon (\langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{recon}) \quad (18)$$

Diese Lernregel wird Kontrastive Divergenz (Contrastive Divergence(CD)) genannt und liefert bessere Ergebnisse als die vorherige Regel. Zunächst wurde angenommen, dass die Lernregel grob den Gradienten der logarithmischen Wahrscheinlichkeit approximiert. Danach nahm man an, dass die Differenz zweier Kullback-Leiber Divergenzen approximiert wurden. Was so eine Divergenz ist, wird in [8] erklärt und würde hier zu weit führen. Da hierbei ein wichtiger Term jedoch vernachlässigt wird, folgt es auch diesem Gradienten nicht genau und Untersuchungen ergaben, dass die Lernregel gar keinem Gradienten einer Funktion folgt und trotzdem gut funktioniert [4].

Die Idee hinter Kontrastiven Divergenz besteht darin, dass beim normalen Gibbs Sampling die Markov-Kette sich immer weiter vom gewünschten Modell entfernt, da diese immer unabhängiger von den Eingabedaten wird. Um dem entgegen zu wirken, wird für die Lernregel nicht eine Probe aus dem Gleichgewichtszustand der Boltzmann Maschine verwendet, der entsteht, wenn man das Gibbs Sampling lange genug betreibt, sondern eine Probe aus der rekonstruierten Wahrscheinlichkeitsverteilung nach einem Schritt des Gibbs Sampling. Dies sorgt dafür dass die Markov-Kette näher an der tatsächlichen Wahrscheinlichkeitsverteilung der Eingabe bleibt. Da die rekonstruierte Eingabe näher am Gleichgewichtszustand der Boltzmann Maschine ist als die Eingabe im Gibbs Schritt davor, ist garantiert, dass die kontrastive Divergenz niemals negativ wird und nur Null erreicht, wenn das Modell perfekt ist. [3]

Ein weiterer Schritt, um das Lernen eines Modells auf RBMs zu verbessern, besteht darin, n Schritte des Gibbs samplings zu machen, bevor man die Statistik für $\langle v_i h_j \rangle_{recon}$ ermittelt. Dies wird mit CD_n notiert, wobei n die Anzahl der Schritte des Gibbs samplings angibt.

3 Deep-Belief Netze

Deep-Belief Netze sind eine Form von künstlichen neuronalen Netzen. Diese werden meist unüberwacht trainiert, da angenommen wird, dass Labels zu wenig Informationen enthalten über die vorhandenen Daten oder bei großen Datenmengen sehr oft keine Labels für alle Daten vorhanden sind. Zum einen ermöglichen Deep-Belief Netze das Rekonstruieren von einmal trainierten Modellen, z.B. kann ein Deep-Belief Netz, wenn es das Modell von Bildern lernt teile von Bildern oder auch ganze Bilder rekonstruieren. Desweiteren können Deep-Belief Netze zur Klassifikation benutzt werden oder zum Vortrainieren von Feedforward Netzen, da diese nicht so leicht in lokalen Minima stecken bleiben und somit das Netz näher am globalen Minimum konvergiert. Feedforward Netze sind neuronale Netze die ihre Daten nur in eine Richtung propagieren. Zum Trainieren von Deep-Belief Netzen hat sich ein einfacher Algorithmus basierend auf eingeschränkten Boltzmann Maschinen etabliert.

3.1 Greedy Algorithmus zum Trainieren von Deep-Belief Netzen

Ein effizienter Weg ein kompliziertes Modell zu lernen, besteht darin, eine Menge von einfacheren Modellen nacheinander zu lernen. Die Idee des hier vorgestellten greedy Algorithmus besteht darin, dass jedes Modell eine andere Repräsentation der Daten darstellt. Dazu berechnet jedes Modell eine nichtlineare Transformation auf seine Eingabe und die Ausgabe eines Modells wird als Eingabe des nächsten Modells verwendet.

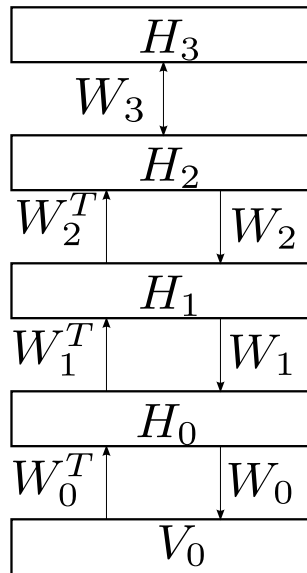


Abbildung 5: Hybrid Netzwerk. Die Ebenen H_3 und H_2 sind mit ungerichteten Kanten verbunden und bilden einen Assoziativspeicher. Die anderen Ebenen sind mit gerichteten Kanten verbunden.

In Abbildung 5 ist ein generatives Modell mit mehreren Ebenen abgebildet. Als generatives Modell bezeichnet man Modelle, die mit Hilfe von versteckten Parametern zufällige und beobachtbare Daten generieren, wie zum Beispiel bei einer Markov-Kette. Die oberen beiden Ebenen kommunizieren über ungerichtete Kanten und simulieren damit unendlich viele weitere gerichtet verbundene Ebenen mit gleichen Gewichten durch die Daten propagiert werden. Das Prinzip wird in Abbildung 4 dargestellt, nur dass diesmal jeder neue Schritt als neue Ebene interpretiert wird. Es gibt keine Verbindungen zwischen den Knoten einer Ebene, sodass diese unabhängig voneinander sind. Man kann gute Parameter für W_0 finden, indem man annimmt, dass die Gewichte zwischen den höheren Ebenen eine komplementäre Verteilung (complementary prior) abbilden, um den “Explained Away“-Effekt für W_0 auszulöschen. Dies gilt jedoch nur wenn die Ebenen die gleiche Knotenzahl besitzen, wenn dies nicht der Fall ist, werden die Gewichte der höheren Ebene mit Zufallszahlen einer Gauss-Verteilung initialisiert [5].

Der “Explained Away“-Effekt beschreibt hierbei dass zwei Knoten stark antikorrelieren,

also je höher die Wahrscheinlichkeit in einem Knoten ist, desto geringer ist sie in einem anderen. Bildhafter in Abbildung 6 beschrieben ist die Ursache für ein wackelndes Haus durch einen hineinrasenden LKW sehr gering wenn bereits ein Erdbeben das Haus wackeln lässt. Hierbei wäre das wackelnde Haus ein Ausgangsknoten und Erdbeben und der LKW jeweils ein Eingangsknoten. Wenn einer der beiden Eingangsknoten also aktiviert ist, reicht ein Ereignis aus zur Erklärung, da das Eintreten beider Ereignisse sehr unwahrscheinlich ist. Dieses Verhalten macht es schwierig Schlussfolgerungen des Netzes nachzuvollziehen.

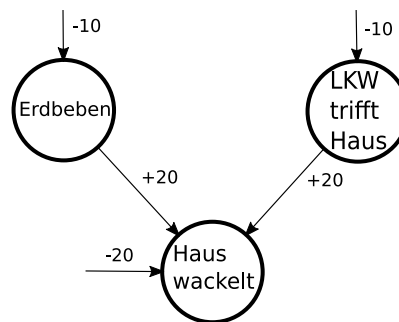


Abbildung 6: Beispiel für den Explained Away Effekt. Der Bias von -10 am Erdbebenknoten bedeutet dass dieser ohne beobachtet zu werden e^{10} mal wahrscheinlicher aus als aktiviert ist. Beide Ereignisse zusammen haben eine Wahrscheinlichkeit von e^{-20} , da diese unabhängig voneinander sind. Dies ist also sehr unwahrscheinlich. Deshalb erklärt die Aktivierung eines Knotens den anderen weg.

Die Annahme, dass die höheren Ebenen eine komplementäre Verteilung bilden, führt dazu, dass das Lernen von W_0 lediglich das Trainieren einer RBM darstellt und mithilfe der Contrastive Divergence eine gute Approximation gefunden wird. Sobald W_0 gelernt ist, kann man W_0^T , also die transponierte Gewichtsmatrix, als Eingabe für die erste versteckte Ebene benutzen, indem man die Eingabedaten mit dieser multipliziert.

Im allgemeinen findet die eingeschränkte Boltzmann Maschine kein perfektes Modell für die Eingabedaten. Das Modell kann durch einen einfachen greedy Algorithmus verbessert werden:

1. Lerne W_0 .
2. Benutze W_0^T um die Verteilung der einzelnen Variablen in der ersten versteckten Ebene zu approximieren.
3. Lerne eine weitere Restricted Boltzmann Maschine mithilfe der "Daten", die durch W_0^T generiert wurden, für die nächste Ebene.

Jede durch den Algorithmus trainierte Ebene verbessert das Modell. Gleichzeitig bedeutet dies auch, dass jedes Modell durch Hinzufügen von neuen Ebenen bis zu einem bestimmten Punkt verbessert werden kann, wenn man jede Ebene sorgfältig genug trainiert. Wenn alle Ebenen die gleiche Anzahl von Knoten haben, besteht eine Verbesserung darin, dass jede

neue Ebene mit den bereits gelernten Gewichten initialisiert wird [5]. Bei unterschiedlicher Knotenzahl werden die Gewichte mit Zufallszahlen initialisiert.

3.2 Deep-Belief Netze zur Klassifikation

Nachdem man das Netz mit dem im vorigen Kapitel beschriebenen Algorithmus trainiert hat, muss man immer noch die Labeldaten ins Netz speisen wenn man dieses für eine Klassifikation nutzen möchte. Dazu gibt es im wesentlichen zwei verschiedene Methoden: Die erste Methode besteht darin, beim Training auf der vorletzten Ebene eine Gruppe von Softmax-Neuronen hinzuzufügen. Diese haben eine andere Aktivierungsfunktion als die normalen Neuronen:

$$p_j = \frac{e^{x_j}}{\sum_{i=1}^K e^{x_i}} \quad (19)$$

Diese Aktivierungsfunktion sorgt dafür, dass die Aktivierung eines Neurons von der Aktivierung der anderen Neuronen abhängig ist und somit immer nur ein Neuron in der Gruppe aktiviert ist. Indem man abhängig vom Label ein Neuron fest auf 1 setzt beim Trainieren, lernt das Netz welche Aktivierungszustände zu einer Klasse gehören. Dieses Training ist unüberwacht. Der Netzaufbau für diese Art der Klassifikation wird in Abbildung 7 gezeigt. Wenn man nach dem Training einen Datensatz durchs Netz propagiert und anschließend ein paar Gibbs Samplings zwischen den letzten beiden Ebenen macht, kann man in der Softmax Gruppe das erkannte Label ablesen.

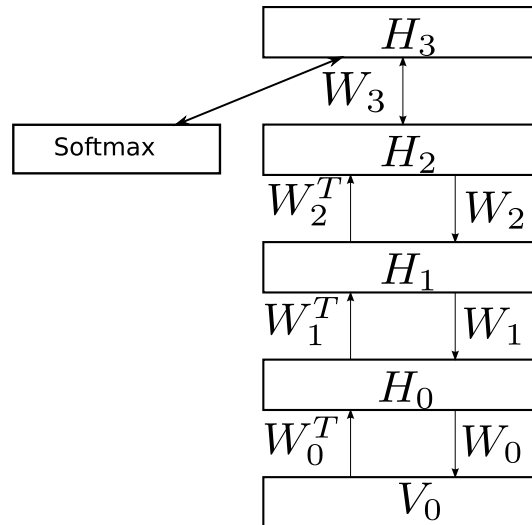


Abbildung 7: Netzwerk mit Softmax-Gruppe zum Klassifizieren. Aus Sicht von H_3 werden die Softmax-Neuronen als normale Eingabe betrachtet. An der Softmax-Gruppe kann später die gefundene Klasse abgelesen werden.

Die zweite Methode besteht darin, dass Deep-Belief Netz normal zu trainieren und nach dem Training eine Ebene mit Ausgabenneuronen an die letzte Ebene hinzuzufügen. Diese

Neuronen können entweder die sigmoide Funktion als Aktivierung nutzen oder auch eine Softmax Gruppe bilden. Nun trainiert man die Gewichte zwischen der neuen Ebene und der letzten Ebene mit Backpropagation um die Labelinformationen mit den Aktivierungszuständen des Netzes zu verbinden. Der Aufbau des so entstandenen Netzes wird in Abbildung 8 gezeigt. Um das Ergebnis im Anschluß noch weiter zu verbessern, kann man die Gewichte des ganzen Netzes nochmals mit Backpropagation feiner einstellen, allerdings sollte das Netz auch ohne diesen Schritt bereits positive Ergebnisse liefern[6].

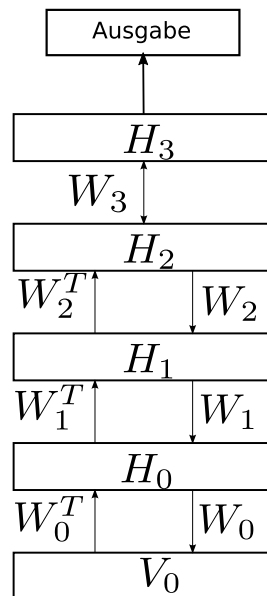


Abbildung 8: Netzwerk mit Ausgabeschicht zur Klassifizierung. Die Gewichte zwischen H_3 und der Ausgabe werden mit Backpropagation initialisiert.

3.3 Implementation

Als Grundgerüst wurde ein bereits implementiertes Feedforward-Netz, welches bereits Backpropagation kann, genutzt. Abbildung 9 zeigt die Klassenhierarchie des verwendeten Grundgerüsts. In rot sind die für den Aufbau und die Verwendung von Deep-Belief Netzen benötigten Klassen markiert.

Die Klasse `ContrastiveDivergenceConfig` liest die Parameter für die Kontrastive Divergenz aus einer Datei ein. Die Klasse kann einfach erweitert werden und setzt momentan die Anzahl der Gibbschritte, die Anzahl der Epochen, die Lernrate und die Größe der zu trainierenden Batches.

`DBNLayer` erweitert die Layerklasse soweit, dass statt normalen Neuronen Neuronen für das Deep-Belief-Netz benutzt werden und neben den rückwärtsgerichteten Gewichten der normalen Neuronen auch vorwärtsgerichtete Gewichte besitzt. Diese Gewichte stellen jedoch lediglich Zeiger auf die rückwärtsgerichteten Gewichte dar. Desweiteren wurde eine Möglichkeit geschaffen die Gewichte für das Deep-Belief Netz zu initialisieren. Dazu wird

für jedes Gewicht eine zufällige Zahl einer gausschen Normalverteilung mit einem Erwartungswert von 0 und einer Standardabweichung von 0,1 gezogen.

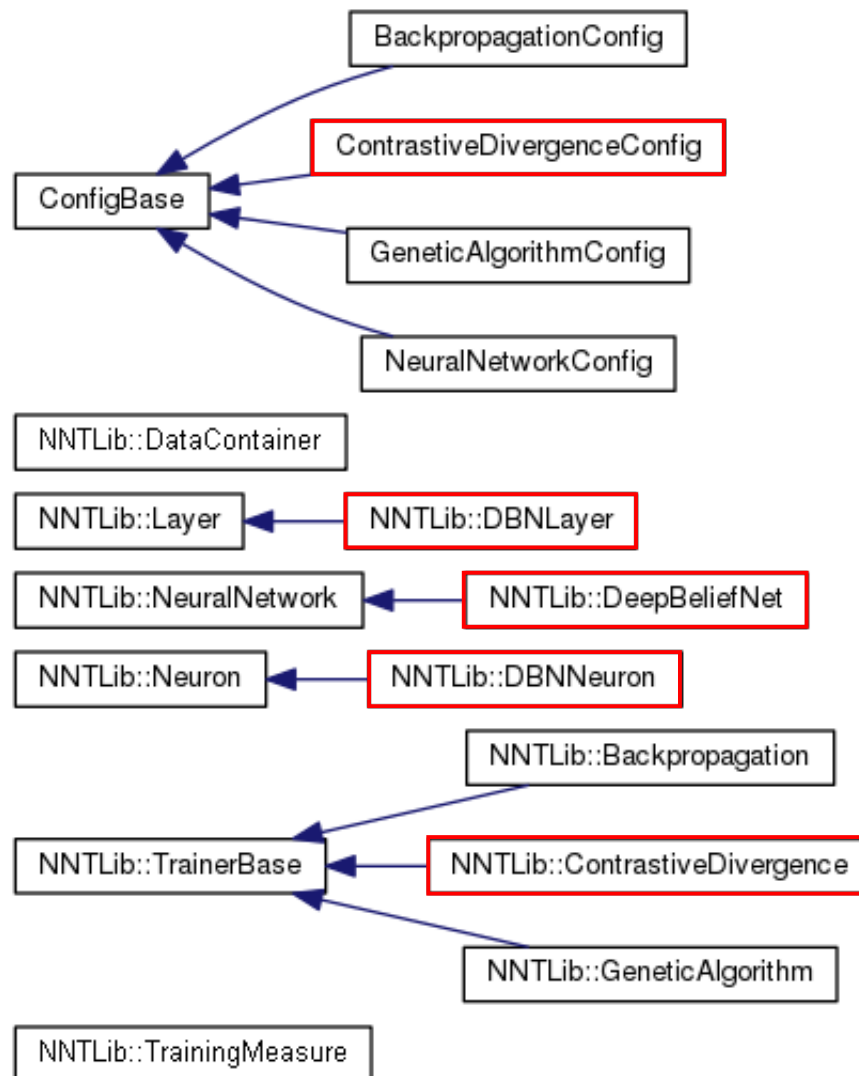


Abbildung 9: Klassenhierarchie des verwendeten Frameworks, rot umrandete Klassen wurden in dieser Arbeit hinzugefügt.

DeepBeliefNet stellt die Netzstruktur für Deep-Belief Netze zur Verfügung. Im Gegensatz zu den normalen Netzwerken des Frameworks, besitzen Deep-Belief Netze auch auf der Eingabeebene eigene Neuronen. Dies ist wichtig, da man für die Kontrastive Divergenz die Aktivierungswahrscheinlichkeiten der einzelnen Neuronen speichern muss. Desweiteren besitzt das Netz die Funktion seine Gewichte so abzuspeichern, dass diese von einem normalen neuronalen Netz wieder eingelesen werden können. Dabei gehen die Biasgewichte in Rückwärtsrichtung verloren, da diese von dem normalen neuronalen Netz mit Backpropagation nicht benötigt werden.

DBNNeuron ist eine Erweiterung der Neuron Klasse. DBNNeuronen können ihre Aktivierungswahrscheinlichkeit speichern und besitzen die Möglichkeit auf die Gewichte zum nächsthöheren Layer zuzugreifen. Außerdem kam eine Funktion zum Initialisieren von Biasgewichten hinzu. Diese bekommt entweder den Datencontainer für die Eingabedaten des Layers oder NULL übergeben. Wenn kein Container übergeben wird, wird der Bias auf 0 gesetzt. Dies ist der Fall, wenn der Biasknoten für die versteckte Ebene zuständig ist. Wenn ein Container übergeben wird, wird aus den Eingabedaten der Initialwert des Bias für jeden sichtbaren Knoten i berechnet:

$$\log[p_i/(1 - p_i)] \quad (20)$$

p_i ist hierbei die durchschnittliche Aktivierungswahrscheinlichkeit des Knoten i über alle Eingabedaten.

Als letzte wichtige Erweiterung wurde die Klasse ContrastiveDivergence implementiert, diese enthält den Lernalgorithmus für das Deep-Belief Netz.

```

Data: Eingabedaten, Epochenanzahl, Anzahl der Gibsschritte, Lernrate
Result: Trainiertes Deep Belief Netz
Initialisiere leere Datencontainer für Eingabedaten für jede Ebene
for jede Ebene do
    Initialisiere Bias für versteckte Ebene mit 0
    Initialisiere Bias für sichtbare Ebene mit  $\log[p_i/(1 - p_i)]$ 
    for Jeden Datensatz do
        Setze Eingabedaten als Aktivierungswahrscheinlichkeit für sichtbare Ebene
        Bestimme Aktivierungen der sichtbaren Knoten
        Berechne Wahrscheinlichkeiten für versteckte Knoten:
         $p(h_j = 1|\vec{v}) = \sigma(b_j + \sum_i v_i w_{ij})$ 
        Sammle Statistik für  $\langle v_i h_j \rangle_{data}$ 
        for Anzahl Gibbsschritte do
            Berechne Wahrscheinlichkeiten für sichtbare Knoten:
             $p(v_i = 1|\vec{h}) = \sigma(a_i + \sum_j h_j w_{ij})$ 
            Berechne Wahrscheinlichkeiten für versteckte Knoten
        end
        Sammle Statistik für  $\langle v_i h_j \rangle_{model}$ 
        Passe Gewichte an:  $w_{ij} = w_{ij} + \epsilon (\langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{recon})$ 
        Passe Bias für sichtbare Knoten an:  $a_i = a_i + \epsilon (\langle v_i \rangle_{data} - \langle v_i \rangle_{recon})$ 
        Passe Bias für versteckte Knoten an:  $b_i = b_i + \epsilon (\langle h_i \rangle_{data} - \langle h_i \rangle_{recon})$ 
    end
    for Jeden Datensatz do
        Propagiere Eingabedaten durch trainierte Ebene und speichere Ergebnis in
        Datencontainer für die nächste Ebene
    end
end

```

Algorithm 1: Implementation der Kontrastiven Divergenz

Im Anschluss an die Kontrastive Divergenz werden die Gewichte zwischen der Ausgabe- und der vorletzten Ebene mit Hilfe von Backpropagation trainiert.

Das Framework reagiert auf verschiedene Eingabeparameter zum Setzen von Konfigurationsdaten für den Aufbau des Netzes und welcher Lernalgorithmus verwendet werden soll. Wenn neben einer Konfiguration für die Kontrastive Divergenz auch eine Konfiguration für die Backpropagation angegeben ist, wird das Netz im Anschluss nochmal feiner trainiert.

4 Versuch

Es sollen die verschiedenen Methoden zum Klassifizieren von Dateneingaben mit Deep-Belief Netzen untersucht werden. Dazu werden verschiedene Netzgrößen mit ansonsten den gleichen Parametern getestet. Zur Verfügung stehen zwei verschiedene Datensätze.

4.1 Achsensymmetrie und Rotationssymmetrie

Die erste Datenbank besteht aus 1346 Datensätzen und beinhaltet Merkmale zu einzelnen Punkten. Der Datensatz wurde von Prof. Dr. Christoph Dalitz zur Verfügung gestellt. Jeder Datensatz enthält vier Merkmale zu einem einzelnen Punkt. Dieser Punkt ist entweder der Kategorie Achsensymmetrisch oder Rotationssymmetrisch zugeordnet. Die einzelnen Klassen wurden mit 1 und 0 kodiert. Die vier Merkmale sind auf das Intervall $[0, 1]$ skaliert. Für jedes Merkmal gibt es ein Eingabeneuron und jede Klasse hat ein Ausgabeneuron. Aufgrund fehlender Testdaten wird eine einfache Kreuzvalidierung mit $K = 5$ gewählt. Daraus ergibt sich für jeden Trainingsdatensatz eine Größe von 1076 und jeder Testdatensatz besteht aus 270 Datensätzen.

4.2 MNIST Datenbank

Die Mixed National Institute of Standards and Technology² beinhaltet 70000 Datensätze. Diese Datenbank ist ein Standarddatensatz zum Vergleichen von verschiedenen Verfahren der Mustererkennung und bietet daher eine große Vergleichsmöglichkeit. Der Datensatz besteht aus 60000 Trainingsdaten und 10000 Testdaten. Die Daten bestehen aus $28 \cdot 28$ Grauwertpixeln und ergeben damit eine Eingabe von 784 Neuronen. Da die Bilder zur Erkennung von Ziffern sind ergeben sich 10 Ausgabeneuronen, für jede Klasse ein Neuron. Das Ziel ist es, dass das Neuron zu dem die Klasse des Bildes gehört aktiviert wird und die anderen Neuronen nicht. Zur Eingabe ins Netz wurden die Grauwerte auf das Intervall $[0, 1]$ skaliert.

4.3 Versuchsdurchführung

Die beiden Datenbanken werden mit verschiedenen Topologien getestet. Als Grundlage werden die Ergebnisse des Backpropagation Algorithmus verwendet und diese dann mit

²<http://yann.lecun.com/exdb/mnist/>

dem Vortraining durch Kontrastive Divergenz verglichen und mit einem Deep-Belief Netz mit einer Softmax Gruppe. Da es zu Aufwändig ist, alle möglichen Netztopologien zu testen, werden nur ein paar Stichproben ausprobiert.

4.3.1 Achsensymmetrie und Rotationssymmetrie

Bei diesem Datensatz werden Topologien mit 1-2 versteckten Ebenen und verschiedenen Anzahlen von Neuronen auf den versteckten Ebenen untersucht. Für die Backpropagation wird eine Epochenanzahl von 1000 gewählt mit einer Lernrate von 0,5. Die Kontrastive Divergenz hat eine Epochenanzahl von 100 und eine Lernrate von 0,1. Es wird ein Gibbs Schritt verwendet. Für das Trainieren der letzten Ebene nach der Kontrastiven Divergenz wird eine Epochenanzahl von 1000 und eine Lernrate von 0,1 verwendet. In Tabelle 1 sieht man....

NR	Topologie	Backpropagation			BP mit Vortraining		
		MSE	Fehler	Laufzeit	MSE	Fehler	Laufzeit
1	4-2-1						
2	4-4-1						
3	4-6-1						
4	4-8-1						
5	4-2-2-1						
6	4-4-4-1						
7	4-4-8-1						
8	4-8-8-1						

Tabelle 1: Vergleich verschiedener Netztopologien mit einem Ausgabeneuron

In den Arbeiten zu Deep-Belief Netzwerken wurden für die Backpropagation Softmaxneuronen als Ausgabeneuronen um ein besseres Ergebnis zu liefern. Daher wurde ein weiterer Versuch durchgeführt. Dieser Vergleicht auch ein Deep-Belief Netz mit einer Softmax Gruppe in der vorletzten Ebene die auch als Ausgabe verwendet wird.

NR	Topologie	Backpropagation		BP mit Vortraining		DBN mit Softmax Gruppe	
		MSE	Fehler	MSE	Fehler	MSE	Fehler
1	4-2-2						
2	4-4-2						
3	4-6-2						
4	4-8-2						
5	4-2-2-2						
6	4-4-4-2						
7	4-4-8-2						
8	4-8-8-2						

Tabelle 2: Vergleich verschiedener Netztopologien mit Softmax-Neuronen als Ausgabe

4.3.2 MNIST Datenbank

Bei der MNIST Datenbank wird zunächst die in [5] verwendete Topologie getestet und mit dem Ergebnis der Arbeit verglichen. Da diese nur den relativen Fehler enthält, wird hier auch nur dieser gemessen. Diese Topologie besteht aus 5 Ebenen. Die erste Ebene hat 784 Eingabeknoten für die Bilddaten. Danach kommen drei versteckte Ebenen mit 500, 500 und 2000 Knoten und eine Ausgabebene mit 10 Knoten zur Klassifizierung. Dabei wurde die Anzahl der Gibbs Schritte variiert und auf 1, 10 und 30 gesetzt. Eine Frage die sich stellt ist, ob die dritte versteckte Ebene mit 2000 Knoten so groß sein muss, oder ob das Netz auch mit weniger Knoten auskommt. Daher wurden noch zwei Topologien mit 1000 bzw. 500 versteckten Knoten getestet.

NR	Topologie	Gibbs Schritte	BP mit Vortraining	DBN mit Softmax Gruppe
1	784-500-500-2000-10	1		
2	784-500-500-2000-10	10		
3	784-500-500-2000-10	30		

Tabelle 3: Vergleich verschiedener Netztopologien mit Softmax-Neuronen als Ausgabe

4.3.3 MNIST Datenbank

5 Mögliche Verbesserung von Deep-Belief Netzen

5.1 Probleme

5.2 Einfache Verbesserung

Literatur

- [1] G. E. Hinton, T. J. Sejnowski, *Learning and relearning in Boltzmann machines*, Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Volume 1: Foundations, MIT Press, Seite 282-317, 1986
- [2] Geoffrey Hinton, *A Practical Guide to Training Restricted Boltzmann Machines*, Department of Computer Science, University of Toronto, 2010
- [3] Geoffrey Hinton, *Training Products of Experts by Minimizing Contrastive Divergence*, Neural Computation 14 Seite 1771-1800, 2002
- [4] Ilya Sutske, Tijmen Tieleman, *On the convergence properties of contrastive divergence*, Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS), 2010
- [5] Geoffrey Hinton, Yee-Whye Teh, *A Fast Learning Algorithm for Deep Belief Nets*, Neural Computation 18 Seiten 1527-1554, 2006
- [6] G. E. Hinton, R. R. Salakhutdinov, *Reducing the dimensionality of data with neural networks*, Science Vol. 313. no. 5786n Seite 504 - 507, 2006
- [7] J. J. Hopfield, *Neural networks and physical systems with emergent collective computational properties*, Proceedings of the National Academy of Sciences (USA) 79, Seite 2554-2558, 1982
- [8] S. Kullback, R.A Leibler, *On information and sufficiency*, Annals of Mathematical Statistics 22 No. 1, Seite 79–86, 1951

Abbildungsverzeichnis

1	Modell eines Hopfield Netzes. Alle Neuronen besitzen eine Ein- und Ausgabe und sind über Gewichte jeweils mit allen anderen Neuronen verbunden . . .	4
2	Modell einer Boltzmann Maschine. Gelbe Kreise stellen versteckte Neuronen dar, weiße Kreise sind sichtbare Neuronen.	5
3	Modell einer eingeschränkten Boltzmann Maschine. Oben sind die versteckten Knoten und unten die nach außen sichtbaren Knoten.	7
4	Markov-Kette mit Gibbs Sampling. Diese wird mit einem Trainingsvektor initialisiert. Die sichtbaren und unsichtbaren Neuronen werden abwechselnd neu berechnet bis ein Equilibriumszustand erreicht wird.	9
5	Hybrid Netzwerk. Die Ebenen H_3 und H_2 sind mit ungerichteten Kanten verbunden und bilden einen Assoziativspeicher. Die anderen Ebenen sind mit gerichteten Kanten verbunden.	11
6	Beispiel für den Explained Away Effekt. Der Bias von -10 am Erdbebenknoten bedeutet dass dieser ohne beobachtet zu werden e^{10} mal wahrscheinlicher aus als aktiviert ist. Beide Ereignisse zusammen haben eine Wahrscheinlichkeit von e^{-20} , da diese unabhängig voneinander sind. Dies ist also sehr unwahrscheinlich. Deshalb erklärt die Aktivierung eines Knotens den anderen weg.	12
7	Netzwerk mit Softmax-Gruppe zum Klassifizieren. Aus Sicht von H_3 werden die Softmax-Neuronen als normale Eingabe betrachtet. An der Softmax-Gruppe kann später die gefundene Klasse abgelesen werden.	13
8	Netzwerk mit Ausgabeschicht zur Klassifizierung. Die Gewichte zwischen H_3 und der Ausgabe werden mit Backpropagation initialisiert.	14
9	Klassenhierarchie des verwendeten Frameworks, rot umrandete Klassen wurden in dieser Arbeit hinzugefügt.	15