

Maschinelles Lernen

Als letzten Themenblock dieser Vorlesung wollen wir uns einer Klasse von Algorithmen zuwenden, deren Ziel es ist, künstliche Formen von Wissen zu generieren. Derartiges Wissen soll dabei aufgrund von Beispielen in einer "Lernphase" angeeignet werden und anschließend auf neue Situationen / Datensätze angewandt werden. Diese Idee des "maschinellem Lernens" gehört damit zu algorithmischen Ansätzen, große Datensätze ("big data") dahingehend zu analysieren, daß charakteristische Merkmale extrahiert werden und damit eine Dimensionsreduktion durchgeführt wird.

Eine wichtige Rolle spielt die Form der Wissensrepräsentation. Hier unterscheiden wir zwei Typen – symbolische vs. subsymbolische Systeme. In symbolischen Systemen wird das Wissen, d.h. sowohl die Beispiele als auch die induzierten Regeln, explizit repräsentiert. In subsymbolischen Systemen wie sogenannten artifiziellen neuronalen Netzen wird Wissen implizit repräsentiert, d.h. es wird ein berechenbares Verhalten "antrainiert", ein direkter Einblick in die Lösungswege bleibt aber oft verwahrt.

Wir wollen uns hier mit letzterem Zugang beschäftigen und künstliche neuronale Netze besprechen. Betrachten wir dazu das folgende Beispiel:

Ziffernerkennung + künstliche neuronale Netze

Wir wollen dem Computer beibringen, Ziffern zu erkennen und damit Postleitzahlen zu lesen

Im visuellen System des Menschen wird diese Aufgabe im primären visuellen Kortex bewältigt – einer Ansammlung von 140 Millionen Neuronen mit über 10 Milliarden Verbindungen.

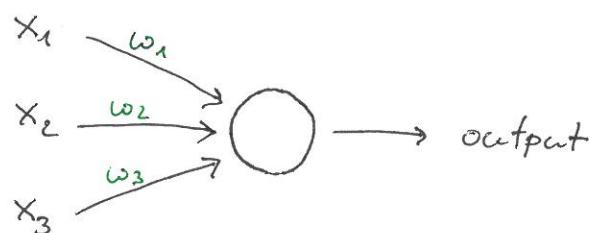
Eine derartige neuronale Struktur wird in künstlichen neuronalen Netzen modelliert und was die Zahl der Neuronen und Verbindungen betrifft erheblich vereinfachten Form implementiert. Tatsächlich reichen schon weniger als 100 Zeilen Code, um ein künstliches neuronales Netz zu implementieren, das Ziffern mit einer Genauigkeit von 96% erkennt. Kommerzielle Codes erreichen eine noch weit höhere Genauigkeit und werden etwa in der Postleitzahlenerkennung seit gut einem Jahrzehnt eingesetzt.

In folgenden wollen wir drei Hauptaspekte von künstlichen neuronalen Netzen besprechen:

- künstliche Neuronen : Perzeptronen vs. Sigmoidalne Neuronen
- Netzarchitekturen
- Optimierung neuronaler Netze : überwachtes Lernen

Künstliche Neuronen - Perzeptronen (Frank Rosenblatt 1950-1960's)

Sogenannte Perzeptronen stellen die historisch erste und einfachste Form künstlicher Neuronen dar. Ein Perzepron akzeptiert multiple binäre Eingaben x_1, x_2, \dots, x_n und produziert einen einzigen binären Output, was wir schematisch darstellen können als



Der binäre Output eines solchen Perzeptions berechnen wir wie folgt

$$\text{output} = \begin{cases} 0 & \sum_i w_i \cdot x_i \leq \text{threshold} \\ 1 & \sum_i w_i \cdot x_i > \text{threshold} \end{cases}$$

Gewichte

Beispiel: Sie wollen entscheiden, ob Sie am Wochenende zu einem Jazz-Konzert gehen wollen und betrachten dabei folgende Kriterien

1. Ist das Wetter gut ? x_1
2. Möchte Ihre Freundin / Ihr Freund mitkommen ? x_2
3. Gibt es eine gute ÖPNV - Verbindung ? x_3

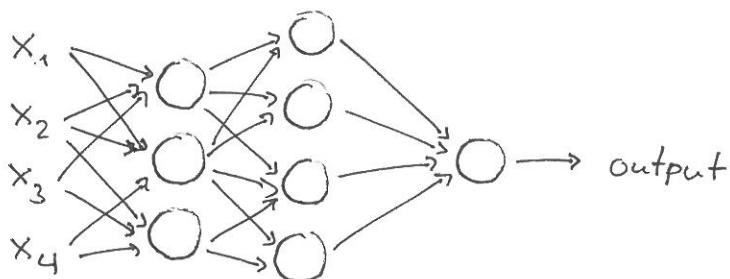
Option 1: $w_1 = 6, w_2 = 2, w_3 = 2$ threshold = 5

Sie mögen Jazz wirklich sehr und gehen, wenn das Wetter gut ist.

Option 2: $w_1 = 6, w_2 = 2, w_3 = 2$ threshold = 3

Sie gehen, wenn gutes Wetter ist oder wenn Ihr Freund mitkommt und ÖPNV gut ist.

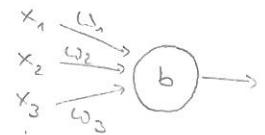
An dem Beispiel erkennen wir, wie ein einziges Perzepton genutzt werden kann, um eine (relativ einfache) Entscheidung zu treffen. Entsprechend können wir uns jetzt schon vorstellen, daß ein Netzwerk von Perzeptoren komplexere Entscheidungen herbeiführen kann.



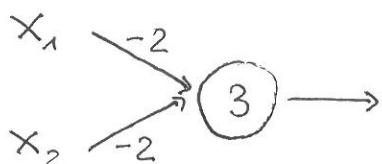
Zur weiteren Formalisierung wollen wir schließlich noch obige Notation für die Berechnung des Outputs vereinfachen

$$\text{output} = \begin{cases} 0 & \vec{w} \cdot \vec{x} + b \leq 0 \\ 1 & \vec{w} \cdot \vec{x} + b > 0 \end{cases}$$

Bias



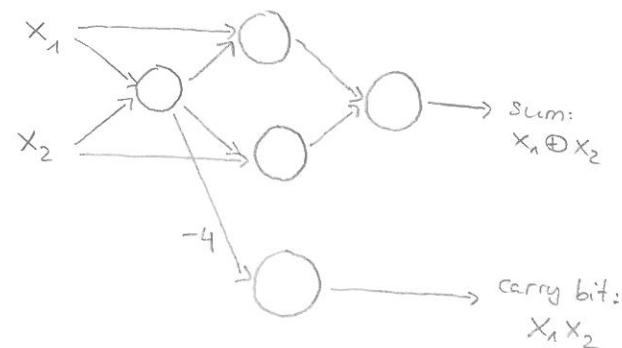
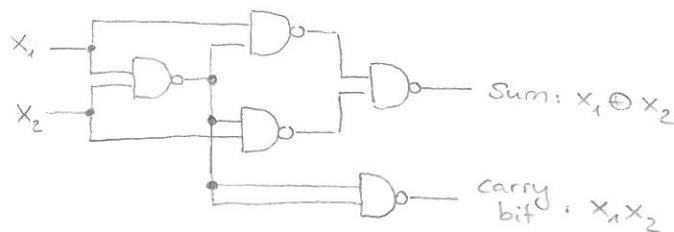
Um die algorithmische Mächtigkeit künstlicher neuronaler Netze abschätzen zu können, ist es wichtig festzustellen, daß Perzeptoren elementare logische Boolefunktionen berechnen können. Dies sei am sogenannten NAND Gate veranschaulicht.



x_1	x_2	output
0	0	1
0	1	1
1	0	1
1	1	0

Das ist eine wichtige Einsicht, denn das NAND-Gate ist ein sogenanntes universelles Gate, d.h. eine beliebig komplexe logische Funktion von n binären Inputs läßt sich auf einen Schaltheorie solcher NAND-Gates reduzieren. Ähnliches können wir also auch auf Netzwerke von Perzeptoren schließen – sie erlauben universelles Rechnen und sind somit genau so mächtig wie alle anderen (universellen) Rechenmodelle.

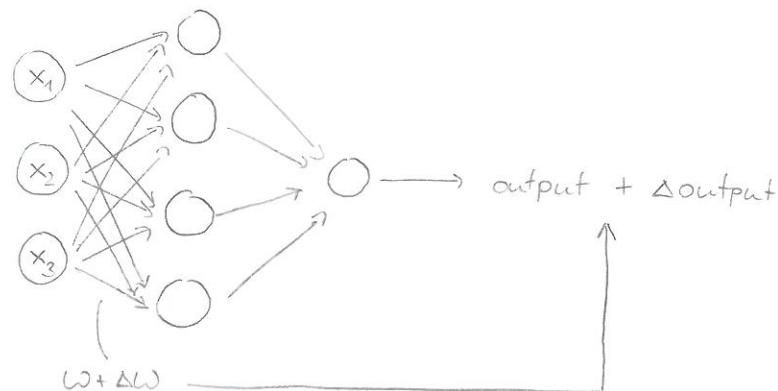
Beispiel: Binäre Addition



Jenseits ihrer prinzipiellen Universalität liegt die algorithmische Mächtigkeit von neuronalen Netzen allerdings darin, daß ihre Gewichte so optimiert werden können, daß das Netz bestimmte Funktionalität bekommt. Dieser Optimierungsprozess kann automatisiert werden, etwa indem eine große Menge an Datensätzen von Input/Output Parametern betrachtet werden und das Netz, bzw. dessen Gewichte, so optimiert werden, daß es die vorgegebenen Input/Output-Paare reproduziert kann. Man spricht hierbei von sogenanntem "überwachten Lernen".

Sigmoid Neuronen

Ein solches Lernprozess erfolgt durch eine inkrementelle Variation der Gewichte. Hierfür ist es essentiell, dass kleine Änderungen eines beliebigen Gewichts auch in einer kleinen Änderung des Outputs resultieren.

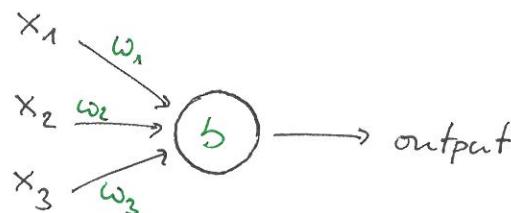


kleine Änderung Δw resultiert in kleine Änderung Δoutput

Wir werden sehen, dass dies eine notwendige Bedingung sein wird, um automatisiertes Lernen zu ermöglichen.

Bedauerlicherweise ist dies aber für ein Netzwerk von Perzeptoren nicht der Fall! Eine kleine Änderung in einem Gewicht oder Bias kann dann führen, dass der binäre Output von 0 nach 1 oder umgekehrt flippt – eine erhebliche Änderung! Dies macht es unmöglich, die Gewichte graduell so zu verändern, dass der gewünschte Output für eine große Input-Menge möglich ist.

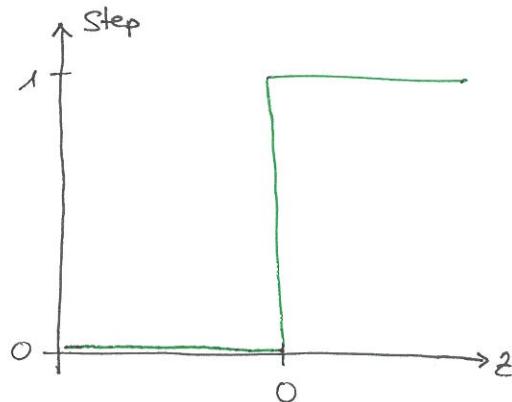
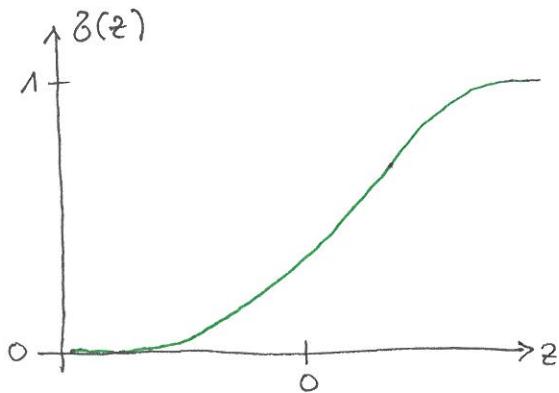
Zum dieses Problem zu beheben, verändern wir die Charakteristik unseres elementaren Neurons und betrachten sogenannte Sigmoid Neuronen:



Wiederum betrachten wir n Input Variablen $\vec{x} = (x_1, x_2, \dots, x_n)$, erlauben aber nunmehr dass die Komponenten $x_i \in \mathbb{R}$ in $[0,1]$ sind. Ebenfalls gewichten wir die Eingaben mit Gewichten $\vec{w} = (w_1, w_2, \dots, w_n)$ und schreiben dem Neuron einen Bias-Wert b zu.

Der Output allerdings sei nicht mehr binär, sondern ebenfalls -6- eine reelle Zahl in $[0, 1]$. Diesen Output berechnen wir über die Sigmoidische Funktion $\delta(z) = \frac{1}{1+e^{-z}}$ als

$$\delta(\vec{x}, \vec{\omega}, b) = \frac{1}{1 + \exp(-\sum_i \omega_i x_i - b)}$$

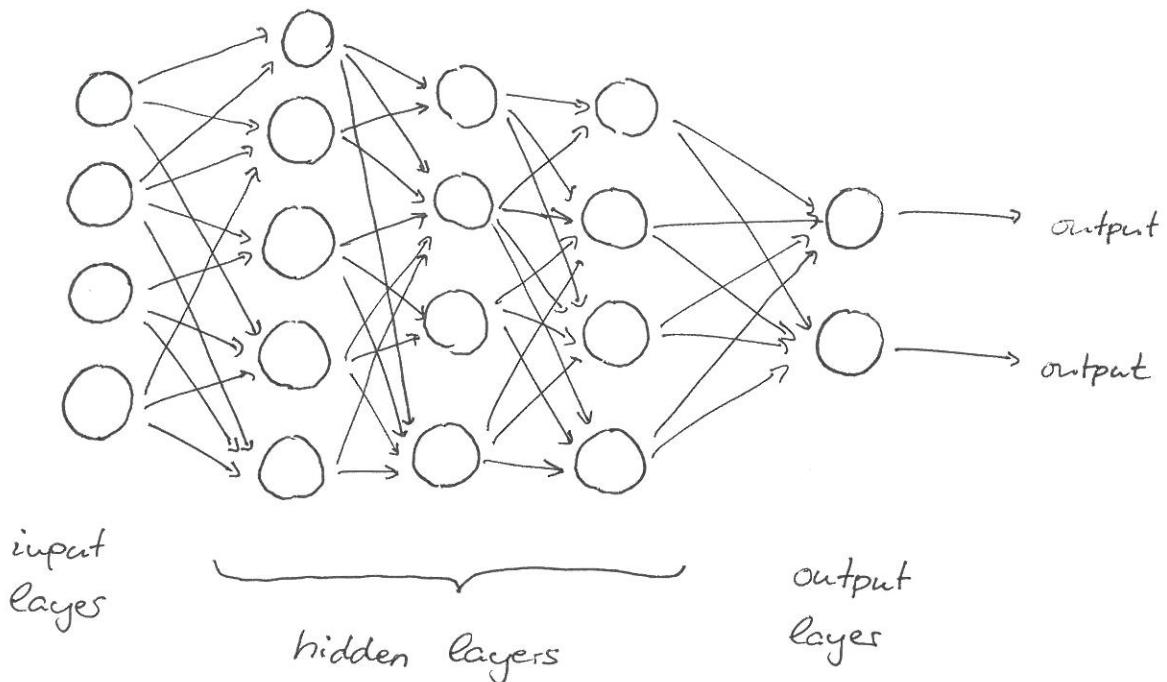


Die sigmoidische Funktion $\delta(z)$ ist eine verbreiterte Stufenfunktion. Im Falle einer solchen Stufenfunktion reduziert sich das sigmoidische Neuron auf das Perzepton. Die genaue Form von $\delta(z)$ ist tatsächlich nicht wichtig, relevant ist vor allem ihre Stetigkeit, die nunmehr garantiert, daß

$$\Delta \text{output} \approx \sum_i \frac{\partial \text{output}}{\partial \omega_i} \Delta \omega_i + \frac{\partial \text{output}}{\partial b} \Delta b$$

Das bedeutet, daß Δoutput eine lineare Funktion der (kleinen) Änderungen $\Delta \omega_i$ und Δb ist, was eben eine notwendige Voraussetzung für das Lernen, die Optimierung des neuronalen Netzwerks ist.

Die Architektur von künstlichen neuronalen Netzen unterteilt sich typischerweise in input layer, sogenannte hidden layers und den output layer:



Die hidden layers bezeichnen dabei all jene Layer zwischen input und output Layer. Die einzelnen Layer sind hierarchisch angeordnet – der Output eines Layers ist der Input des nächsten Layers. Solche Netzwerke, Netzwerke mit nur einem hidden layer nennt man feedforward Netzwerke. Netzwerke mit nur einem hidden layer nennt man "flach" und bezeichnet das maschinelle Lernen solcher Netzwerke als "shallow learning". Netzwerke mit mehreren hidden layers hingegen nennt man "tief" und es hat sich hierfür die Bezeichnung "deep learning" eingebürgert. Die genaue Struktur des Netzwerks, etwa im Bezug auf die Zahl der Neuronen pro Layer oder die Art der Verbindungen, wird im allgemeinen heuristisch bestimmt.

Optimierung des Netzwerks - überwachtes Lernen

Abschließend wollen wir uns dem eigentlichen Lernprozess zuwenden. Dazu verstehen wir eine Optimierung der Gewichte auf den Kanten des Netzwerks und den Bias-Werten einzelner Neuronen für eine gegebene Netzwerkarchitektur.

Beim sogenannten überwachten Lernen geschieht dies anhand eines sogenannten Trainingssets von Daten. Für jeden Datensatz x aus diesem Set sei der gewünschte Output $y(x)$ des Netzwerks bekannt. Die Optimierung der Gewichte w und Bias-Werte b kann dann anhand einer Kostenfunktion der Form

$$C(w, b) = \frac{1}{2n} \sum_x \| y(x) - a \|^2$$

↑ gewünschter Output
↑ tatsächlicher Output

erfolgen, d.h. wir suchen das globale Minimum dieser quadratischen Kostenfunktion $C(w_{\text{opt}}, b_{\text{opt}}) \approx 0$. Diese Optimierung erfolgt gewöhnlich über einen sogenannten gradient descent Algorithmus, d.h. wir folgen jeweils der größten Ableitung der Kostenfunktion. Dazu ist es unabdingbar effizient die Gradienten der Kostenfunktion für kleine Abweichungen Δw und Δb zu berechnen. Dies ist für die komplexe Netzwerkarchitektur kein offensichtlich triviales Unterfangen, wurde aber durch den sogenannten back propagation Algorithmus in idealer Weise gelöst, auf den wir hier leider nicht mehr eingehen können.

Mit den obigen Erklärungen haben wir aber nun alle Elemente in der Hand, um erste eigene Schritte zum maschinellen Lernen selbst zu gehen → siehe Programmier-Tutorial.