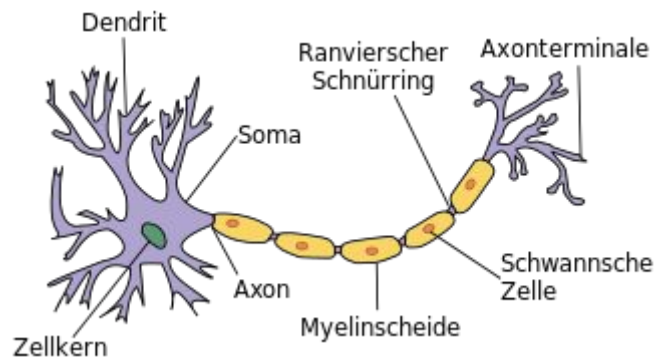


Biologische Grundlagen:

Ein **Neuron** empfängt mehrere Signale über seine Dendriten.

Nach der Verarbeitung der Signaleingänge sendet das Neuron ein Signal über das (nur einfach vorhandene) **Axon**.

Der Informationsübergang findet an der **Synapse** statt: dieser Übergang kann das Signal erregen oder hemmen.



Modellierung im Computer

Beispiel:

Wir haben die drei Neuronen A, B und X, wobei A und B jeweils ihr Signal an X weiterleiten.

Die Synapsen haben die **Gewichte** 0.6 und -0.4.

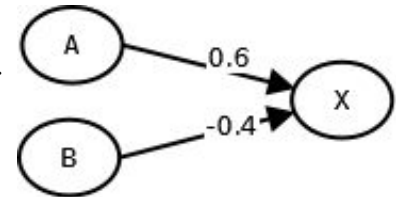
Wenn zum Zeitpunkt $t=0$ die Zelle A nun einen

Aktivierungszustand von 0.5 besitzt und Zelle B einen Aktivierungszustand von 0.2 besitzt, so erhält die Zelle X folgenden Input:

$$0.6 \cdot 0.5 + (-0.4) \cdot 0.2 = 0.22$$

Dies erinnert an ein Skalarprodukt (oder allgemeiner: an ein Matrix*Vektor-Produkt), und tatsächlich könnte man die Rechnung auch folgendermaßen schreiben:

$$\begin{pmatrix} 0,6 & -0,4 \end{pmatrix} \cdot \begin{pmatrix} 0,5 \\ 0,2 \end{pmatrix} = 0.22$$



(Dabei ist der erste Faktor als Matrix mit einer Zeile aber zwei Spalten zu verstehen)

Diese Rechenart wird sich später als nützlich erweisen.

Die Entscheidung, dass alle Ausgänge (multipliziert mit den Gewichten) einfach in der Zelle X addiert werden ist eine Möglichkeit unter vielen (warum nicht multiplizieren?). Die Funktion, die die n Eingabewerte (im Beispiel: $n=2$) auf eine reelle Zahl abbildet (also eine Funktion, welche die Zahlenmengen $\mathbb{R}^n \rightarrow \mathbb{R}$ abbildet) heißt **Propagierungsfunktion**.

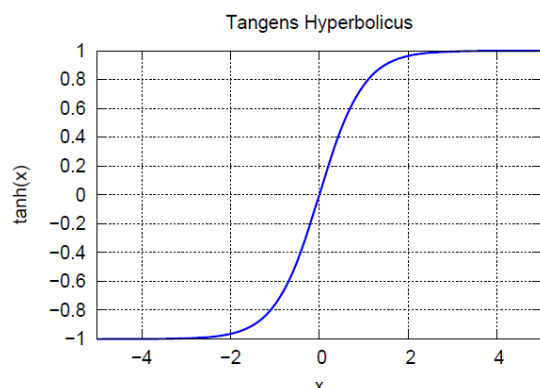
Auf das Ergebnis der Propagierungsfunktion wird nun die **Aktivierungsfunktion** angewendet um den neuen Aktivierungszustand zu berechnen. Da das Ergebnis der Propagierungsfunktion ja beliebig groß werden kann (z.B. wenn man nicht nur zwei Neuronen A und B hat, die auf X einwirken, sondern mehrere tausend) bewirkt die Aktivierungsfunktion oft eine "Normierung" des Signals auf einen bestimmten Wertebereich. Die einfachste Aktivierungsfunktion ist die sogenannte **Schwellenwertfunktion**. Man wählt einen Schwellwert, z.B. 0.15 und legt fest:

$$f(x) = \begin{cases} 0 & \text{für } x < 0,15 \\ 1 & \text{für } x \geq 0,15 \end{cases}$$

In unserem Beispiel würde für das Ergebnis 0.22 aus unserer Propagierungsfunktion der Wert 1 erzeugt, da $0.22 > 0,15$ ist.

Damit erhält man nur noch die Aktivierungszustände 0 und 1.

In der Praxis benötigt man für bestimmte Lernverfahren, mit dem nan das künstliche neuronale Netz trainieren will, aber **differenzierbare** Aktivierungsfunktionen. Beliebte ist z.B. der Tangens Hyperbolicus.



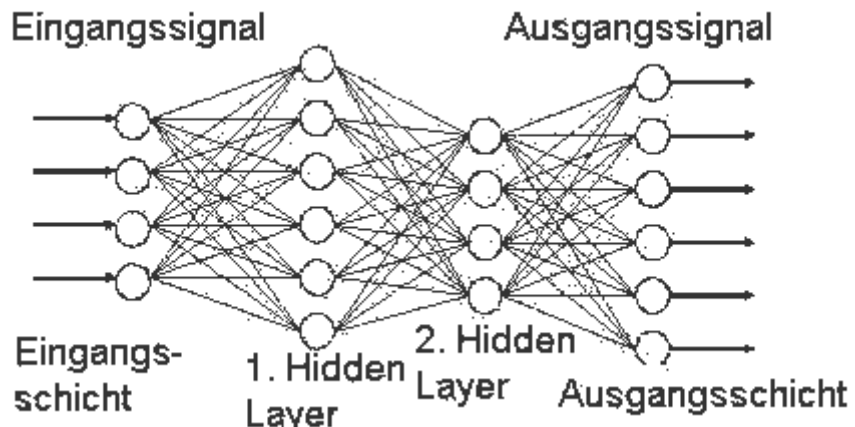
Man kann schließlich auch noch den Aktivierungszustand in eine Ausgabefunktion stecken, die entscheidet, welcher Wert an die Synapsen geschickt wird. In der Regel nimmt man hier aber keine Veränderungen vor, d.h. man wählt hier die Identitätsfunktion, also die Funktion $f(x) = x$, also eine Funktion die keine Veränderung vornimmt.

Grundlegendes Funktionsprinzip im Allgemeinen

Lebewesen haben Sinneszellen, die ihr Signal an das Gehirn weiterleiten können. Ein künstliches neuronales Netz hat (oft) ebenfalls eine **Eingabeschicht**, die man sich als "Sinneszellen" vorstellen kann, welche z.B. die Pixelwerte eines Bildes aufnimmt. Im Beispielmodell vom Anfang wären dies die Neuronen A und B.

Indem die Aktivierungen der Zellen sich über das Netz ausbreiten, gelangen die Daten irgendetwas an die **Ausgabeschicht** (im Beispiel das Ausgabeneuron X).

Oft baut man dazwischen sogenannte **Hidden Layer**, also weitere Verarbeitungsschichten. Das Netz in der Abbildung hat vier

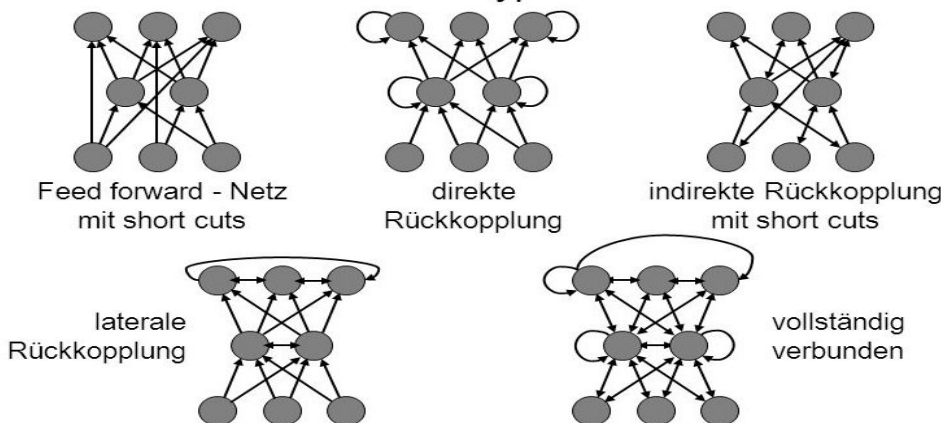


Neuronenschichten, und damit drei **Gewichtsschichten**. Die Abbildung lässt schon erahnen, dass die Anzahl der Synapsen stark ansteigt, je mehr Neuronen dazukommen. Man kann zwar nicht so einfach pauschal sagen, dass ein neuronales Netz umso leistungsfähiger ist, je mehr Neuronen und Synapsen es hat, jedoch hat man in der Praxis meistens wesentlich mehr Elemente als nur ein paar Dutzend wie in den Abbildungen angedeutet.

Mathematisch ausgedrückt ist ein neuronales Netz nichts weiter als eine (komplizierte) Funktion, welche Eingabewerte auf Ausgabewerte abbildet, also eine Abbildung der Art $\mathbb{R}^n \rightarrow \mathbb{R}^m$, wobei n natürlich eine sehr große Zahl sein kann (wobei m möglicherweise sehr klein ist: Möglicherweise nur eine einzige Zahl, beispielsweise ein Gewissheitsindikator dafür, dass ein großes Eingangssignal etwas bestimmtes darstellt (z.B. ein Bild ein Foto von Donald Trump ist).

Topologie von NN

weitere Typen



Unser Beispielnetz vom Anfang mit drei Neuronen hatte also keinen Hidden Layer. Selbstverständlich kann man sich noch weitere **Topologien** (also räumliche Strukturen) ausdenken (und je nach Anwendung wird das auch getan) – siehe Abb.

Falls ein Input immer weiter „nur nach rechts durchgeschubst“ wird

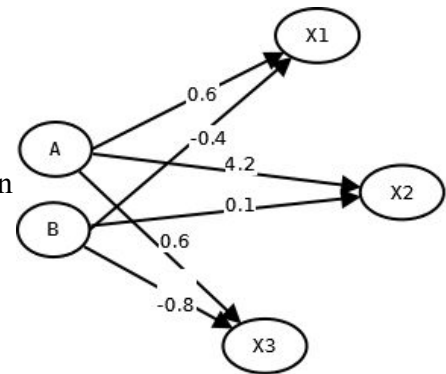
(wenn also keine Rückkopplungen eingebaut sind), so erhält man ein **Feed-Forward-Netz**. Diese sind besonders einfach zu programmieren und gut theoretisch untersuchbar.

Lernen

Die Fähigkeit eines neuronalen Netzes ist in den Gewichten gespeichert. Durch geschickte Änderung der Gewichte trainiert man das Netz so, dass es zu einem bestimmten Input (Aktivierung der Eingangsschicht) einen gewünschten Output (Aktivierung der Ausgabeschicht) liefert. Eine Möglichkeit ist das **überwachte Lernen**: Man präsentiert dem Netz mehrere Eingaben zu denen man die gewünschte Ausgabe kennt. Liefert das Netz die falsche Ausgabe, "dreht" man etwas an den Gewichten. Wenn man genügend gutes Trainingsmaterial hatte (und die Netztopologie, die Aktivierungsfunktion etc. günstig gewählt hat), kann das Netz dies nun generalisieren und auch bei neuen, unbekanntem Eingaben richtige Ergebnisse erzeugen.

Konkrete Betrachtungen und Beispiele

Beispiel 1: Eine etwas ausführlichere Beispielrechnung. Wir nehmen an, dass das abgebildete Netz in der Eingangsschicht jeweils die Aktivierungszustände 1 besitzt, die Aktivierung der Eingangsschicht also durch den Vektor $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ beschrieben werden kann. Darüberhinaus nehmen wir an, dass die Propagierungsfunktion bei jedem Neuron einfach (wie auf Seite 1 beschrieben) die gewichtete Summe berechnet. Die Aktivierungsfunktion soll (bei jedem Neuron) tanh sein, die Ausgabefunktion (bei jedem Neuron) die Identität.



$$\begin{pmatrix} 0,6 & -0,4 \\ 4,2 & 0,1 \\ 0,6 & -0,8 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0,2 \\ 4,3 \\ -0,2 \end{pmatrix}$$

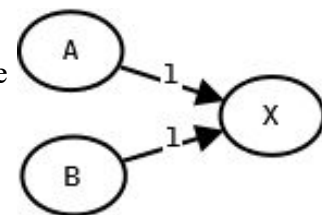
Die Aktivierung von X1 liegt also nun bei $\tanh(0,2) \approx 0,197$
Die Aktivierung von X2 liegt also nun bei $\tanh(4,3) \approx 0,9996$
Die Aktivierung von X3 liegt also nun bei $\tanh(-0,2) \approx -0,197$

Wir stellen fest: Die Aktivierung von A hat also auf das Neuron X2 einen erheblichen Einfluss.

Beispiel 2: Berechnung der logischen "ODER" Funktion.

Wir nehmen für das Neuron X folgendes an:

- Propagierungsfunktion ist die gewichtete Summe (wobei die Gewichte wie im Diagramm jeweils gleich 1 sind)
- Aktivierungsfunktion ist die Schwellenwertfunktion mit Schwellenwert 0.5.
- Ausgabefunktion ist die Identität.



Eingabe (0 0) liefert bei X die Aktivierung 0
Eingabe (1 0) liefert bei X die Aktivierung 1
Eingabe (0 1) liefert bei X die Aktivierung 1
Eingabe (1 1) liefert bei X die Aktivierung 1

Damit hat man eine "ODER"-Operation modelliert. "X ist 1, wenn A oder B 1 ist".

Übung 1: Modelliere ein Netz, welches die "UND"-Operation modelliert.

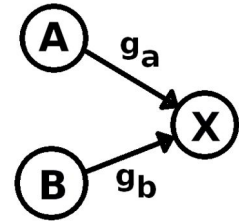
Übung 2: Modelliere ein Netz, welches die "NOT"-Operation modelliert (Das "!" in Java)
(Hinweis: Auch negative Zahlen sind als Gewichte und Schwellwerte erlaubt)

Beispiel 3: Das XOR-Problem

Die XOR-Operation ("entweder oder – aber nicht beides!") verknüpft zwei Bits und liefert als Ergebnis eine 1, wenn genau eines der beiden Eingabebits 1 war – sonst ist das Ergebnis 0. Ein Versuch, diese Operation mit drei Neuronen, (mit der gewichteten Summe als Propagierungsfunktion und einer Schwellenwertfunktion als Aktivierungsfunktion) abzubilden, scheitert. Dazu machen wir uns zunächst folgendes klar:

Kurzer Exkurs: Darstellung einer Aktivierung im Koordinatensystem

Wir nehmen an, dass es zwei Neuronen A und B, gibt, die jeweils eine Aktivierung im Wertebereich zwischen 0 und 1 haben. Mit den entsprechenden Gewichten g_a und g_b bestimmen Sie die Aktivierung von X. Diese Kombination eines „A-Wertes“ und eines „B-Wertes“ lässt sich in einem zweidimensionalen Koordinatensystem (mit A-Achse und B-Achse) als Punkt darstellen. Im Koordinatensystem weiter unten sehen wir z.B. den Punkt (1|0,5), welcher eine Aktivierung von A mit 1 und von B mit 0,5 veranschaulicht.



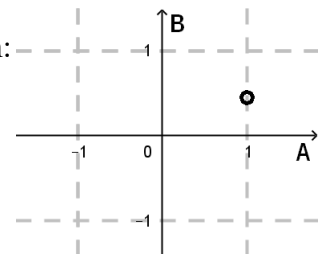
Nehmen wir nun weiter an, dass die beiden Gewichte folgende Werte haben:

$g_a = 2$ und $g_b = 1$. Außerdem soll der Schwellenwert bei 1,5 liegen.

Dann können wir die Aktivierung von X ausrechnen:

Propagierungsfunktion anwenden: $A \cdot g_a + B \cdot g_b = 1 \cdot 2 + 0,5 \cdot 1 = 2,5$

Aktivierungsfunktion anwenden: Ergebnis 2,5 ist größer als 1,5 also ist die Aktivierung von X gleich dem Wert 1.



Weiteres Zahlenbeispiel:

Es seien die Aktivierungen $A=0,4$ und $B=0,3$ gegeben. Mit den gleichen Gewichten $g_a = 2$ und $g_b = 1$ ergibt sich nun: Der Wert der Propagierungsfunktion ist 1,1 und ist damit kleiner als der Schwellenwert 1,5. Damit hat die Aktivierung von Neuron X den Wert 0.

Übung 3:

Veranschauliche für $g_a = 2$ und $g_b = 1$ im Koordinatensystem all jene Werte für A und B, bei denen X den Aktivierungswert 1 erhält. (Hinweis für diejenigen, die in Mathe gut aufgepasst haben: Das Konzept „Koordinatenform“ hilft hier sehr weiter!)

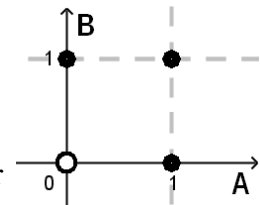
Übung 4:

Zeichne ein entsprechendes Diagramm für die in Beispiel 2 beschriebene „Oder“-Funktion

Weiter mit Beispiel 3:

Wir nehmen noch einmal die Darstellung aus dem Exkurs von vorhin. Es dürfte mittlerweile klar geworden sein, dass die Gewichte zusammen mit dem Schwellenwert eine Grenze in Form einer Geraden im Koordinatensystem beschreiben, welche die 1-Aktivierungen von den 0-Aktivierungen trennt.

Bei den bekannten booleschen Funktionen AND, OR und XOR beschränken wir den Aktivierungswert von A und B üblicherweise auf 0 oder 1 (wir lassen keine Zwischenwerte zu). Im Bild rechts sind drei der vier möglichen Wertekombinationen von A und B schwarz eingefärbt (dies soll den Wert „1“ des Ausgabeneurons X symbolisieren). Das Bild passt damit zur „OR“-Funktion aus Beispiel 2.



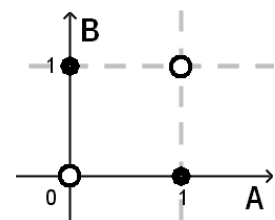
Das entsprechende Bild zur XOR-Funktion ist in folgendem Bild dargestellt:

Wenn $A=0$ und $B=0$, dann ist auch $X=0$

Wenn $A=1$ und $B=1$, dann ist auch $X=0$

Wenn $A=0$ und $B=1$, dann ist auch $X=1$

Wenn $A=1$ und $B=0$, dann ist auch $X=1$



Das Problem beim XOR ist nun, dass man keine Gerade (d.h. keine Gewichte und Schwellwert) finden kann, die die vier Punkte so teilt, dass nur die beiden Punkte auf den Koordinatenachsen (was ja dem Muster "nur genau ein Neuron ist aktiviert" entspricht) den Schwellwert überschreiten. Der Profi sagt: Das Problem ist nicht **linear separierbar**.

Es zeigt sich, dass man dieses Konzept erweitern kann und man kommt zum Schluss: Ein neuronales Feed-Forward-Netz ohne Hidden Layer (also nur eine Eingabe- und eine Ausgabeschicht) kann nur linear separierbare Probleme lösen.

Diese Erkenntnis hat im Jahr 1969 dazu geführt, dass man Forschungsgelder auf dem Gebiet der KI gestrichen hat, weil das Konzept als Ganzes nicht mehr vielversprechend erschien.

Allerdings kann man das Problem ganz einfach durch Einführung von Hidden Layers lösen (siehe Abbildung 5.9)

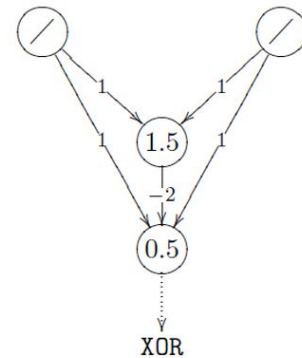
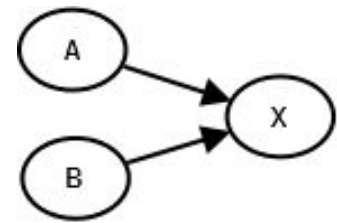


Abbildung 5.9: Neuronales Netz, welches die XOR-Funktion realisiert. Schwellenwerte stehen (soweit vorhanden) innerhalb der Neurone.

Training eines neuronalen Netzes: Die Delta-Regel

Wenn wir uns für eine Struktur eines Netzes entschieden haben, so möchten wir es auch trainieren. Die Frage ist, wie wir nun die Gewichte des Netzes geschickt modifizieren, um ihm etwas "beizubringen". Der Übersichtlichkeit halber wählen wir wieder ein ganz einfaches Netz mit drei Neuronen. Propagierungsfunktion ist die gewichtete Summe, Ausgangsfunktion ist die Identität, Aktivierungsfunktion ist ebenfalls die Identität.



Angenommen, wir möchten dem Netz beibringen, in der Ausgabe X dann eine hohe Aktivität zu zeigen, wenn die Aktivität von B größer ist als die von A.

Als Trainingsset wählen wir:

(0,1) soll (1) liefern,

(1,0) soll (-1) liefern.

Da wir zu Beginn noch gar nicht wissen, wie unsere Gewichte aussehen sollen, wählen wir eine zufällige Belegung: Gewicht $g_{AX} = 0,2$ und Gewicht $g_{BX} = 0,4$

Wir testen nun die beiden Trainingsdatensätze:

Datensatz 1: $0,2 \cdot 0 + 0,4 \cdot 1 = 0,4$

Datensatz 2: $0,2 \cdot 1 + 0,4 \cdot 0 = 0,2$

Fehler:
aber eigentlich sollte 1 rauskommen! $(0,4-1) = -0,6$

aber eigentlich sollte -1 rauskommen! $(0,2 - (-1)) = 1,2$

Wir berechnen nun den Gesamtfehler (als quadratischer Fehler, vgl. z.B. Standardabweichung!):

Gesamtfehler = $(-0,6)^2 + (1,2)^2 = 1,8$

Übung 5: Berechne für obigen Trainingsdatensatz den Gesamtfehler für die Gewichte $g_{AX} = -0,1$ und Gewicht $g_{BX} = 0,6$

Hätten wir einen Gesamtfehler von Null, so wären wir fertig. Weil dem nicht so ist, müssen wir unsere Gewichte anpassen. Den Gesamtfehler kann man auch als Funktion auffassen, welche in Abhängigkeit von g_{AX} und g_{BX} eine reelle Zahl liefert:

Gesamtfehler(g_{AX} , g_{BX}) ist eine Abbildung $\mathbb{R}^2 \rightarrow \mathbb{R}$

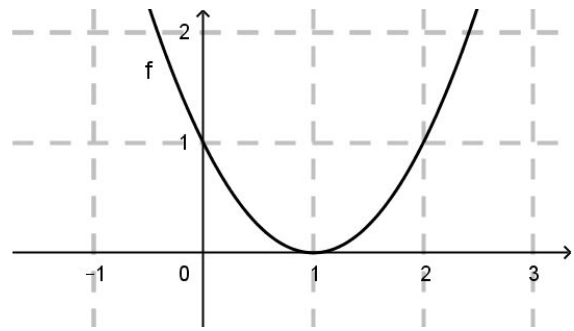
Von dieser Funktion sucht man nun ein Minimum!

Aus der Analysis wissen wir, dass die erste Ableitung für dieses Problem sehr hilfreich ist. Dies ist der Grund dafür, dass man in der Regel differenzierbare (ableitbare) Aktivierungsfunktionen wählt! Die erste Ableitung einer Abbildung $\mathbb{R}^2 \rightarrow \mathbb{R}$ ist ein Vektor der keine zwei Zahlen, sondern zwei Funktionen enthält; er heißt **Gradient**.

Das Standardverfahren "Extremwertsuche" verlangt eigentlich zunächst die Nullstellen der ersten Ableitung. Leider ist es in der Praxis viel zu kompliziert, die Nullstellen des Gradienten zu bestimmen. Daher bedient man sich des sogenannten **Gradientenabstiegsverfahrens**. Man bestimmt die Ableitung an einem Punkt (in unserem Fall an dem Punkt ($g_{AX} = 0,2$, $g_{BX} = 0,4$)) erhält damit einen Vektor, der die (größte) "Steigung" der Funktion widerspiegelt. Ändern wir unsere Gewichte ein wenig in die "Gegenrichtung", so wird der Gesamtfehler kleiner. Dies machen wir nun so oft, bis der Gesamtfehler unter eine bestimmte Schwelle sinkt.

Die Deltaregel im eindimensionalen Fall

Da wir wenig Erfahrung mit Funktionen der Art $\mathbb{R}^2 \rightarrow \mathbb{R}$ haben, machen wir uns das Verfahren an der Funktion $f(x) = (x-1)^2 = x^2 - 2x + 1$ klar:



Der x-Wert ist hier nun die Entsprechung für ein einziges Gewicht und der y-Wert ist die Entsprechung für den Gesamtfehler.

Angenommen, unsere Synapse hat momentan das Gewicht $x = 0,5$, und man erhält den Gesamtfehler $f(0,5) = 0,25$.

Wir können $f'(x) = 2x - 2$ bestimmen und sehen: $f'(0,5) = -1$

Wenn man ausgehend von $x=0,5$ nun also (leicht) nach rechts wandert, so muss $f(x)$ kleiner werden.

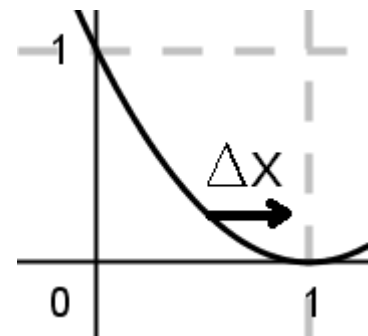
Wir berechnen zunächst die Intensität, um die wir den x-Wert verändern:

$$\Delta x = -c \cdot f'(x)$$

Die Konstante c legt fest, wie stark im allgemeinen die Änderung ausfallen soll. Wir wählen $c=0,4$.
Übrigens: " Δx " (sprich: Delta-x) ist der Namensgeber für diese Lernregel.

Der nächste x-Wert ergibt sich so: $x_{\text{neu}} = x_{\text{alt}} + \Delta x$

Wir erhalten mit unserem Zahlenbeispiel: $\Delta x = -0,4 \cdot f'(0,5) = 0,4$
 $x_{\text{neu}} = 0,9$



Falls der Gesamtfehler (jetzt: $f(0,9) = 0,01$) immer noch zu groß ist, wiederholen wir den Vorgang.

Eine erneute Anwendung der Regel liefert

$$\Delta x = -0,4 \cdot f'(0,9) = -0,4 \cdot -0,2 = 0,08$$
$$x_{\text{neu}} = 0,98$$

Die Regel sorgt also "automatisch" dafür, dass – wenn wir uns dem Minimum annähern, das Δx kleiner wird und wir nicht darüber "hinwegsauen".

Nun beträgt der Fehler $f(0,98) = 0,0004$

Wenn wir c ungeschickt wählen, kann es aber natürlich dennoch passieren, dass wir das Minimum "überschreiten"

Übung: Wende die Deltaregel zwei mal auf obige Funktion $f(x)$ an, wobei der Startwert $x = 2,8$ und $c = 2$ sein soll.

Fazit:

Mit der Delta-Regel kann man anhand eines Trainingsdatensatzes ein neuronales Netz trainieren. Dies klappt nur bei Feed-Forward-Netzen, welche nur eine Gewichtsschicht haben.

Hat man mehrere Gewichtsschichten, so benutzt man die Delta-Regel "mehrmals" hintereinander. Dies Verfahren nennt man **Backpropagation**.