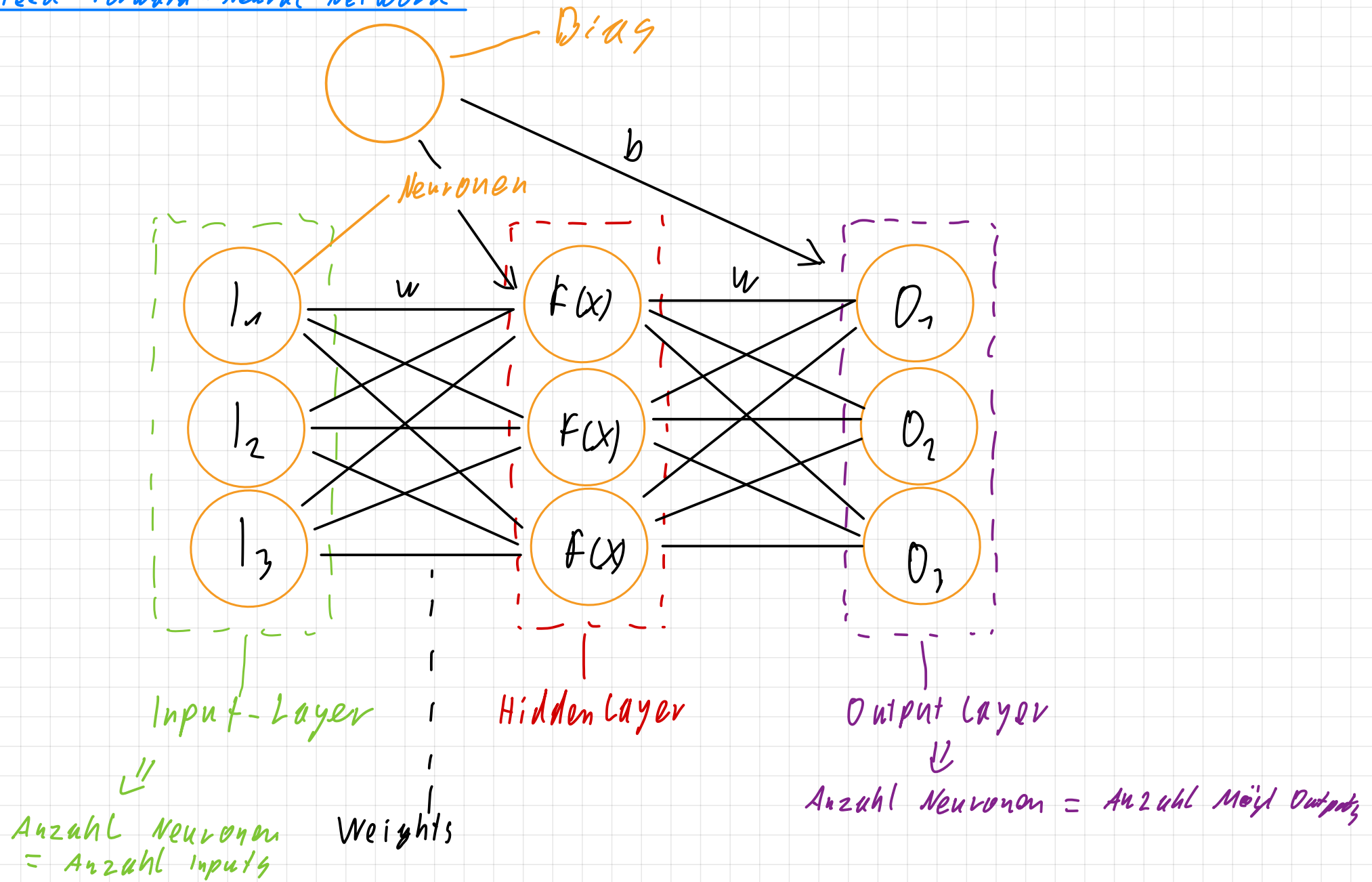


# Feed-Forward-Neural Network



Neuron

$$\begin{matrix} 0 < x_n < 1 \\ \vee \\ f(x) \end{matrix}$$

Beinhaltet meist eine Zahl zwischen 0 und 1.  
Kann auch als mathematische Funktion  $f(x)$  ergänzt werden.



„Aktivierungsfunktion“

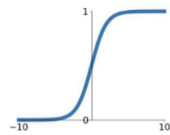
Wie wird das „Learning“ gemacht?

- Weights und Biases der Aktivierungsfunktion können angepasst werden
- Das „Lernen“ ist das Finden der passenden weights und biases

## Activation Functions

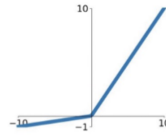
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



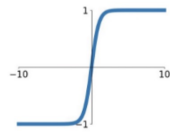
**Leaky ReLU**

$$\max(0.1x, x)$$



**tanh**

$$\tanh(x)$$

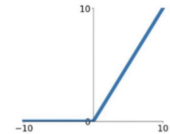


**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

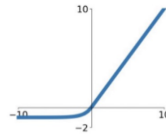
**ReLU**

$$\max(0, x)$$



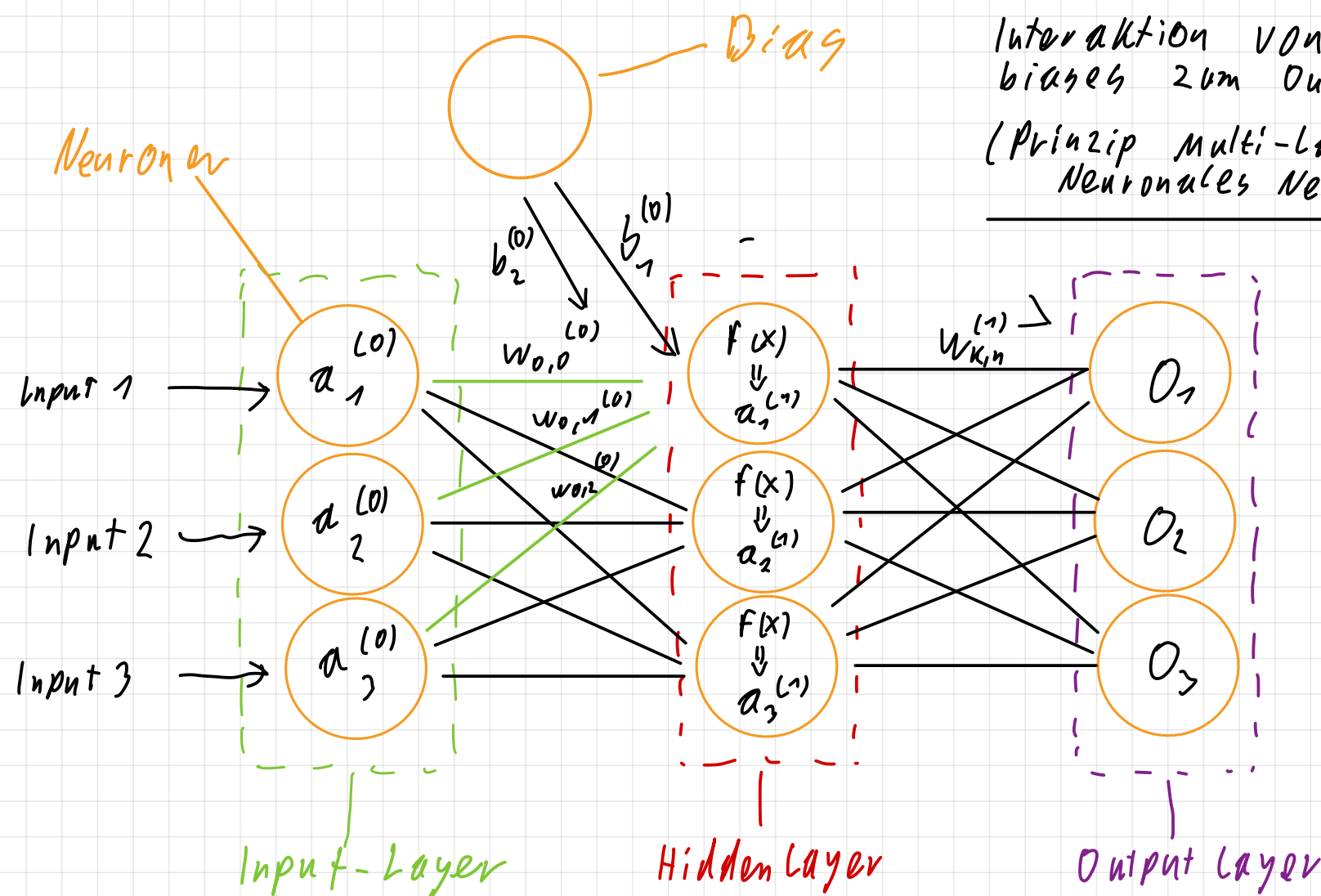
**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Mathe matik

Interaktion von Neuronen, weights, biases zum Output signal!  
(Prinzip Multi-Layer - Perception)  
Neuronales Netzwerk



$$\Rightarrow I = \begin{pmatrix} i_1^{(0)} \\ i_2^{(0)} \\ i_3^{(0)} \\ \vdots \\ i_n^{(0)} \end{pmatrix} ; W = \begin{pmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{pmatrix} ; B = \begin{pmatrix} b_1^{(n)} \\ b_2^{(n)} \\ b_3^{(n)} \\ \vdots \\ b_k \end{pmatrix}$$

Inputs

$$\Rightarrow \begin{cases} a_n^{(0)} = f(W^{(0)} \cdot 1 + B) \\ a_n^{(k)} = \underbrace{f(W^{(k)} \cdot a_n^{(k-1)} + B)} \end{cases}$$

Erste „Berechnung“ mit ersten Weights

Weitere Berechnung,  $k > 0 \in \mathbb{N}$

Aktivierungsfunktion (siehe S. 2)

Die Aktivierungsfunktion erzeugt die Non-Linear Klassifizierung.

Siehe Coding bsp.

## Kostenfunktion

Die Kostenfunktion ist ein Indikator dafür, wie gut / schlecht das neuronale Netzwerk den Input erkennt. Je niedriger die Kosten, desto "besser" ist es.

Sei  $O$  = Outputs,  $G$  = Ziel Outputs,  $C$  = Kosten, dann gilt:

$$C = \sum_{i=0}^{|O|} d(O_i, G_i)^2$$

(Output - Zieloutput)<sup>2</sup>

Da wir daran interessiert sind wie die Kosten über mehrere Daten sind, gilt:

$$\bar{C} = \frac{\sum_{i=0}^{|Daten|} C_i}{|Daten|}$$

Summe aller Kosten: Anzahl der Daten

Siehe Coding bsp.

# Gradient Descent

Da wir nun eine Kostenfunktion haben, interessieren wir uns für das Minimum der Kostenfunktion in Abhängigkeit von weights und biases

$$\Rightarrow C(w, b) \rightarrow \text{Kosten}$$

$$\Rightarrow C'(w, b) = \left( \frac{\partial C}{\partial w}, \frac{\partial C}{\partial b} \right) \quad \begin{array}{l} \text{"Partielle Ableitung"} \\ \text{"}\nabla C\text{"} \leftarrow \text{Gradient} \end{array}$$

Theoretisch würden wir bekannterweise

$$C'(w, b) = (0, 0) \text{ setzen um das Minimum zu bekommen.}$$

In der Praxis wird das nicht gemacht, da dies zu "schwer", unpraktisch und manchmal nicht durchführbar ist besonders für sehr große Dimensionen.

Deshalb kommt das Prinzip vom Gradient Descent "Gradienten Abstieg".

Eine andere Betrachtung des Gradienten, ist die lokale Steigung der Funktion.

$$\nabla C(w, b) = \text{lokale Steigung}$$

Somit können wir  $-\nabla C$ -schrittige Iterationen durchführen um ein lokales Minimum zu erhalten, welches abhängig von den weights und biases ist.

$$C'(w_{k,n}, b_n) = \left( \frac{\partial C}{\partial w_{k,n}} \quad \frac{\partial C}{\partial b_n} \right)$$

Gradient weights der Kosten:

$$\frac{\partial C}{\partial w_{k,n}} = \lim_{h \rightarrow 0} \frac{C(w_{k,n} + h, b_n) - C(w_{k,n}, b_n)}{h}$$

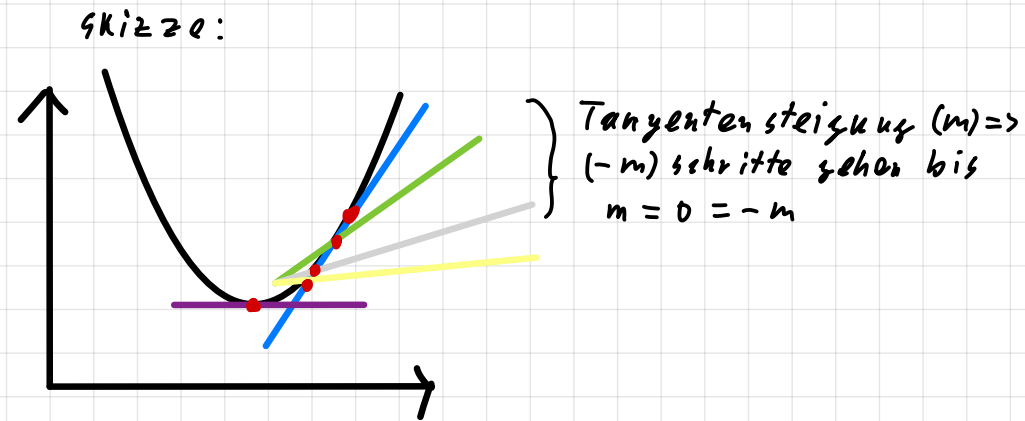
Gradient biases der Kosten:

$$\frac{\partial C}{\partial b_n} = \lim_{h \rightarrow 0} \frac{C(w_{k,n}, b_n + h) - C(w_{k,n}, b_n)}{h}$$

$$f(x) = x^3$$

$$f'(x) = 3x^2$$

↑      ↑  
 $w_i$



Nun müssen alle weights und biases „angepasst“ werden.

Dabei gilt:

$$w_{k,n} = w_{k,n} - \frac{\partial C}{\partial w_{k,n}} \cdot L$$

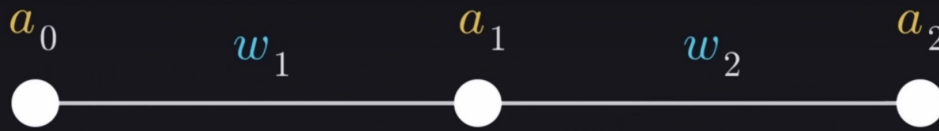
$$b_n = b_n - \frac{\partial C}{\partial b_n} \cdot L$$

, L = Lernrate  $\in (0, 1)$

⇓

Lernrate ist sozusagen die Lerngeschwindigkeit bzw. Schrittweite.

# Backpropagation



(Für Output-Layer)

$$z_1 = a_0 w_1 + b_1$$

$$a_1 = A(z_1)$$

$$z_2 = a_1 w_2 + b_2$$

$$a_2 = A(z_2)$$

$$c = C(a_2, y)$$

$$\frac{\partial c}{\partial w_2} = \frac{\partial z_2}{\partial w_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial c}{\partial a_2}$$

$$\frac{\partial c}{\partial b_2} = \frac{\partial z_2}{\partial b_2} \times \frac{\partial a_2}{\partial z_2} \times \frac{\partial c}{\partial a_2}$$

$$\frac{\partial c}{\partial a_2} = 2 \cdot d(O_i, G_i)$$

Ableitung Kostenfunktion.

$$\frac{\partial a_2}{\partial z_2} = A'(z_2) = \begin{cases} 0, & z_2 < 0 \\ 1, & z_2 > 0 \end{cases}$$

(für Relu)

Ableitung Aktivierungsfunktion.

$$\frac{\partial z_2}{\partial w_2} = a_1$$

$$\frac{\partial z_2}{\partial b_2} = 1$$

Ableitung „Signal“ bevor Aktivierungsfunktion.

$$\Rightarrow \underline{\underline{\frac{\partial c}{\partial w_2} = a_1 \cdot \begin{pmatrix} 0, & z_2 < 0 \\ 1, & z_2 > 0 \end{pmatrix} \cdot 2 \cdot (O_i - G_i)}}$$



$$\Rightarrow \underline{\underline{\frac{\partial C}{\partial b_i} = 1 \cdot \begin{pmatrix} 0, z_2 < 0 \\ 1, z_2 > 0 \end{pmatrix} \cdot 2 \cdot (0_i - b_i)}}$$

(Für HiddenLayer)

daher wie oben.

wird fortgeführt beim ankommen zur Backpropagation.