

Deep Learning

Grundlagen Grundbegriffe Grundübungen

Alfatraining Kurs 2024-08 Dozent: Karsten Flügge

Themen des Tages

Tag 2

Perceptron

Parameter,

Fehler,

Optimierung

Rechenschemen: Backpropagation

Modell

Parameter (freie Variablen)

Feste Werte (bias)

Optimierung

Training: Anpassen an Daten

Einfachstes Modell: ein bit "ja/nein"

Newton Verfahren

Idee: gehe in Richtung der größten Änderung (Ableitung)

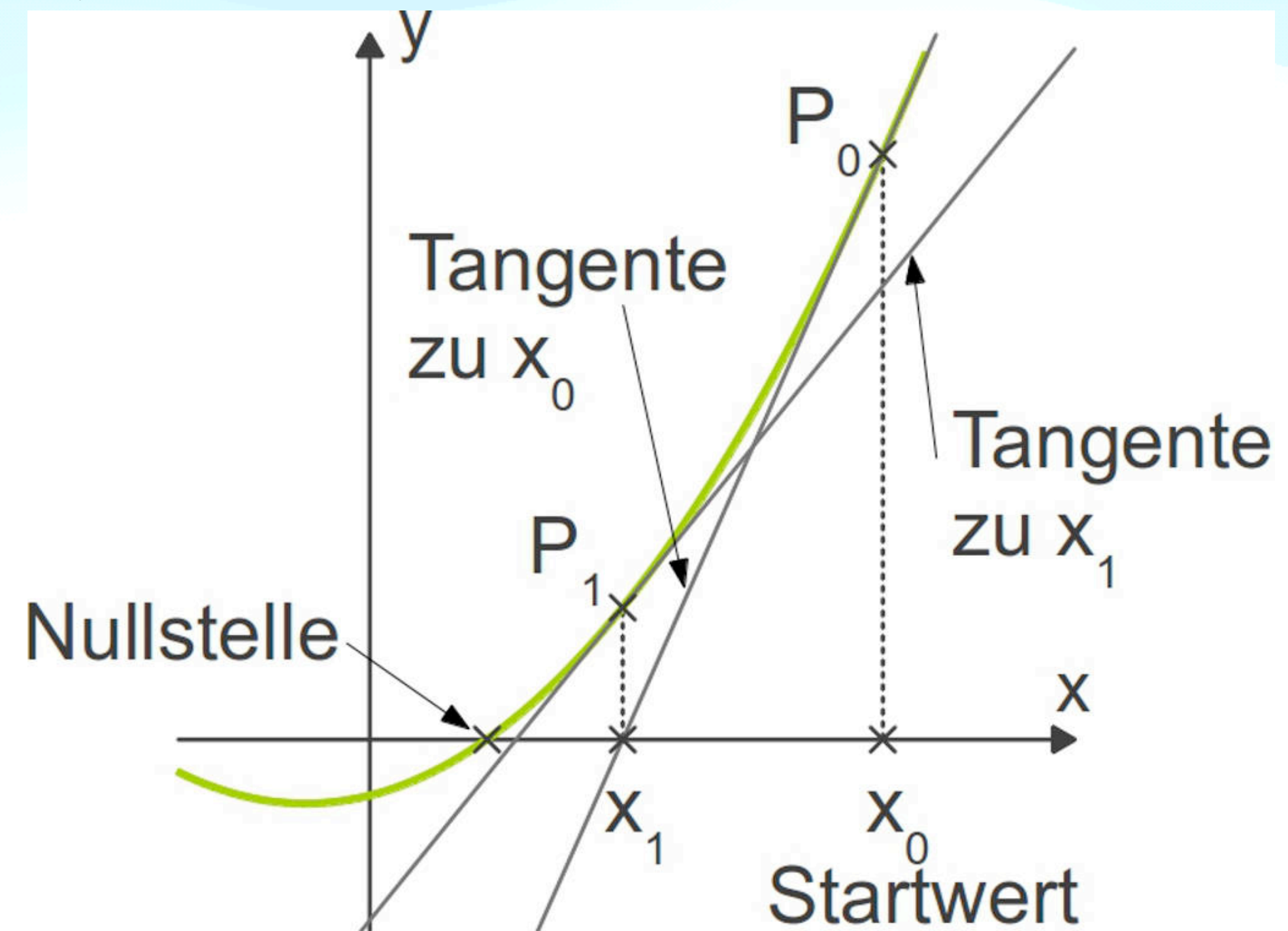
vereinfachtes newton verfahren

```
while fehler(x) > 0.01:
```

```
    ableitung_von_fehler = 2*x # Tangente
```

```
    x = x - ableitung_von_fehler * Schrittweite
```

```
print(fehler(x))
```



Newton Verfahren in \mathbb{R}^n

Idee: gehe in Richtung der größten Änderung (Ableitung)

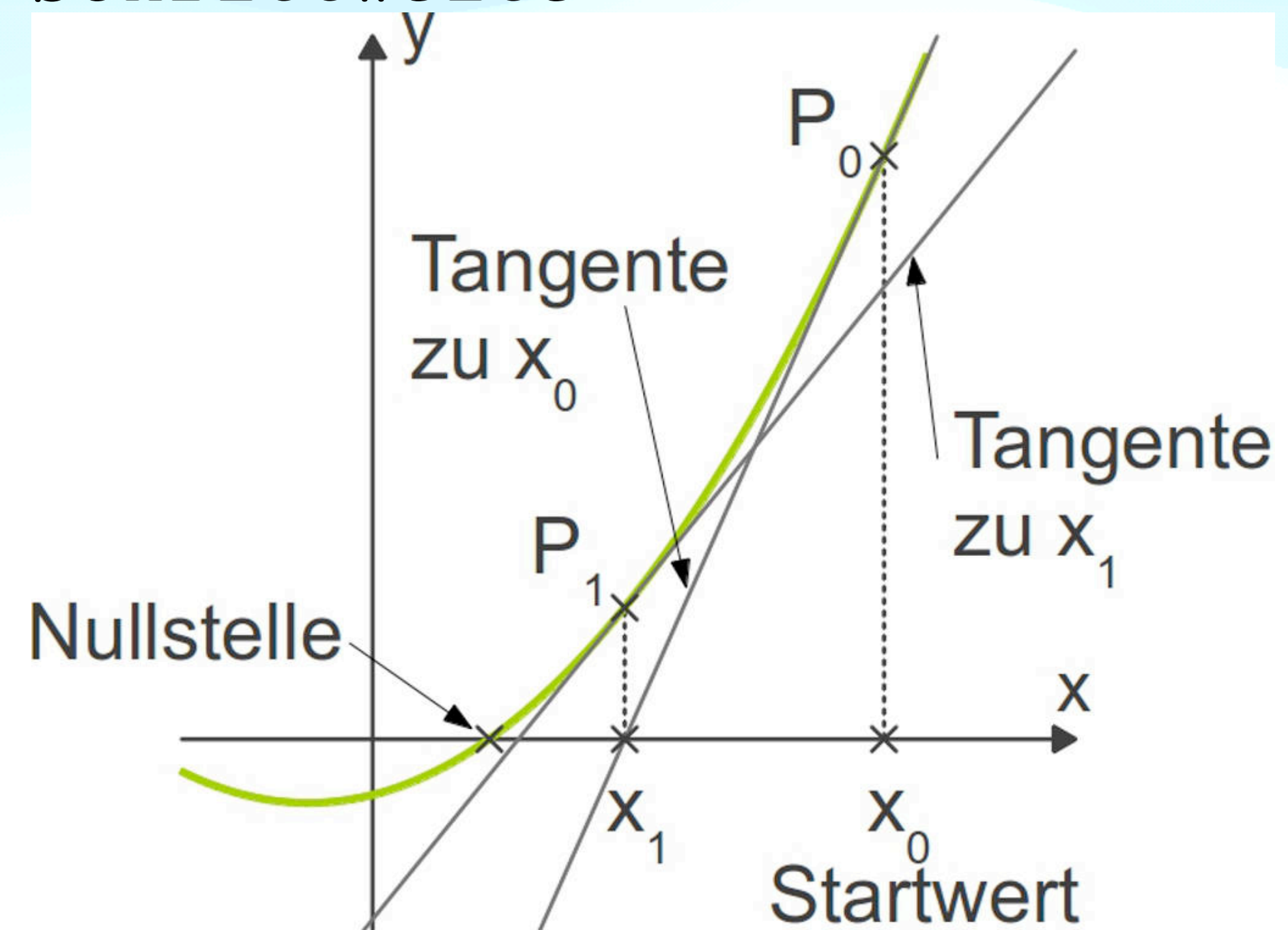
vereinfachtes newton verfahren

```
while fehler(x) > 0.01:
```

```
    gradient_vom_fehler = ...
```

```
    x = x - gradient_vom_fehler * Schrittweite
```

```
print(fehler(x))
```



Aufgaben

Berechnung von Aktivierungen

1) Manuell

Modell Linear : $y = \sum w_i * x_i + w_0$ $f(\text{input}, \text{weights}) = \dots$

a) input = $[\frac{1}{2}, \frac{1}{2}]$ $w_0 = \frac{1}{2}$ weights = $[\frac{1}{2}, \frac{1}{2}]$

b) input = $[1, 2]$ weights = $[1, 2, 3]$ (erstes Element w_0 ist bias)

Modell Linear mit Aktivierung : $y = \max(0, \sum w_i * x_i)$ (x_0 ist 1 per convention, für bias w_0)

c) input = $[-\frac{1}{2}, \frac{1}{2}]$ weights = $[\frac{1}{2}, -\frac{1}{2}, -\frac{1}{2}]$

d) input = $[1, 2]$ weights = $[-1, -2, -3]$

Modell Linear mit Aktivierung : $y = \tanh(\sum w_i * x_i)$

e) input = $[1, 1]$ weights = $[0, 0, 0]$

f) input = $[0, 0]$ weights = $[\frac{1}{2}, \frac{1}{2}, \frac{1}{2}]$

g) input = $[\frac{1}{2}, \frac{1}{2}]$ weights = $[\frac{1}{2}, \frac{1}{2}, \frac{1}{2}]$

h) input = $[1, 2]$ weights = $[1, 2, 3]$

2) Mit Python

so wie e...f nur mit $y = \text{sigmoid}(\sum w_i * x_i)$

Hier sind weitere Aufgaben zur Wiederholung im Gradienten:
1. Gradient einer einfachen Funktion
Bestimmen Sie den Gradienten der Funktion $f(x, y) = x^2 + y^2$.
2. Gradient einer Potenz
Bestimmen Sie den Gradienten der Funktion $f(x, y) = x^3 + y^3$ im Punkt $(2, -1)$.
3. Gradient einer Bruchfunktion
Bestimmen Sie den Gradienten der Funktion $f(x, y) = \frac{1}{x^2 + y^2}$.
4. Gradient einer quadratischen Funktion
Bestimmen Sie den Gradienten der Funktion $f(x, y) = x^2 - 2xy + y^2$.

$$\text{Gradient } \nabla f(x, y) = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

Aufgaben

Berechnung von mehreren Aktivierungen gleichzeitig mit **EINER OPERATION** (mit Python):

```
y=tanH(W*x)
```

```
input = [1,1]
```

```
e) weights = [0,0,0]
```

```
f) weights = [1/2,1/2,1/2]
```

```
g) weights = [1/2,1/2,1/2]
```

```
h) weights = [1,2,3]
```

dann für

```
y=sigmoid(W*x)
```

Hier sind weitere Aufgaben zur Wiederholung im Gradienten:
1. Gradient einer einfachen Funktion
Bestimmen Sie den Gradienten der Funktion $f(x,y) = x^2 + y^2$.
2. Gradient einer Potenz
Bestimmen Sie den Gradienten der Funktion $f(x,y) = \sin(x) + \exp(y)$.
3. Gradient einer Bruchfunktion
Bestimmen Sie den Gradienten der Funktion $f(x,y) = \frac{1}{x^2 + y^2}$.
4. Gradient einer quadratischen Funktion
Bestimmen Sie den Gradienten der Funktion $f(x,y) = x^2 - 2xy + y^2$.

$$\text{Gradient } \nabla f(x, y) = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

Aufgaben

Vorblick: Tensor Operationen

Berechnung von mehreren Aktivierungen gleichzeitig mit EINER OPERATION (mit Python):

$y = \text{tanh}(W * X)$

e) input = [1,1] weights = [0,0,0]

f) input = [0,0] weights = [$\frac{1}{2}, \frac{1}{2}, \frac{1}{2}$]

g) input = [$\frac{1}{2}, \frac{1}{2}$] weights = [$\frac{1}{2}, \frac{1}{2}, \frac{1}{2}$]

h) input = [1,2] weights = [1,2,3]

$X = [[1,1,1], [1,0,0], [1,1/2,1/2], [1,1,2]]$

$W = [[0,0,0], [1/2,1/2,1/2], [1/2,1/2,1/2], [1,2,3]]$

so wie e...f nur mit $y = \text{sigmoid}(\sum w_i * x_i)$

Hier sind weitere Aufgaben zur Wiederholung im Gradienten:
1. Gradient einer einfachen Funktion
Bestimmen Sie den Gradienten der Funktion $f(x,y) = x^2 + y^2$.
2. Gradient einer Potenz
Bestimmen Sie den Gradienten der Funktion $f(x,y) = \sin(x) + \exp(y)$.
3. Gradient einer Bruchfunktion
Bestimmen Sie den Gradienten der Funktion $f(x,y) = \frac{1}{x^2 + y^2}$.
4. Gradient einer quadratischen Funktion
Bestimmen Sie den Gradienten der Funktion $f(x,y) = x^2 - 2xy + y^2$.

$$\text{Gradient } \nabla f(x, y) = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

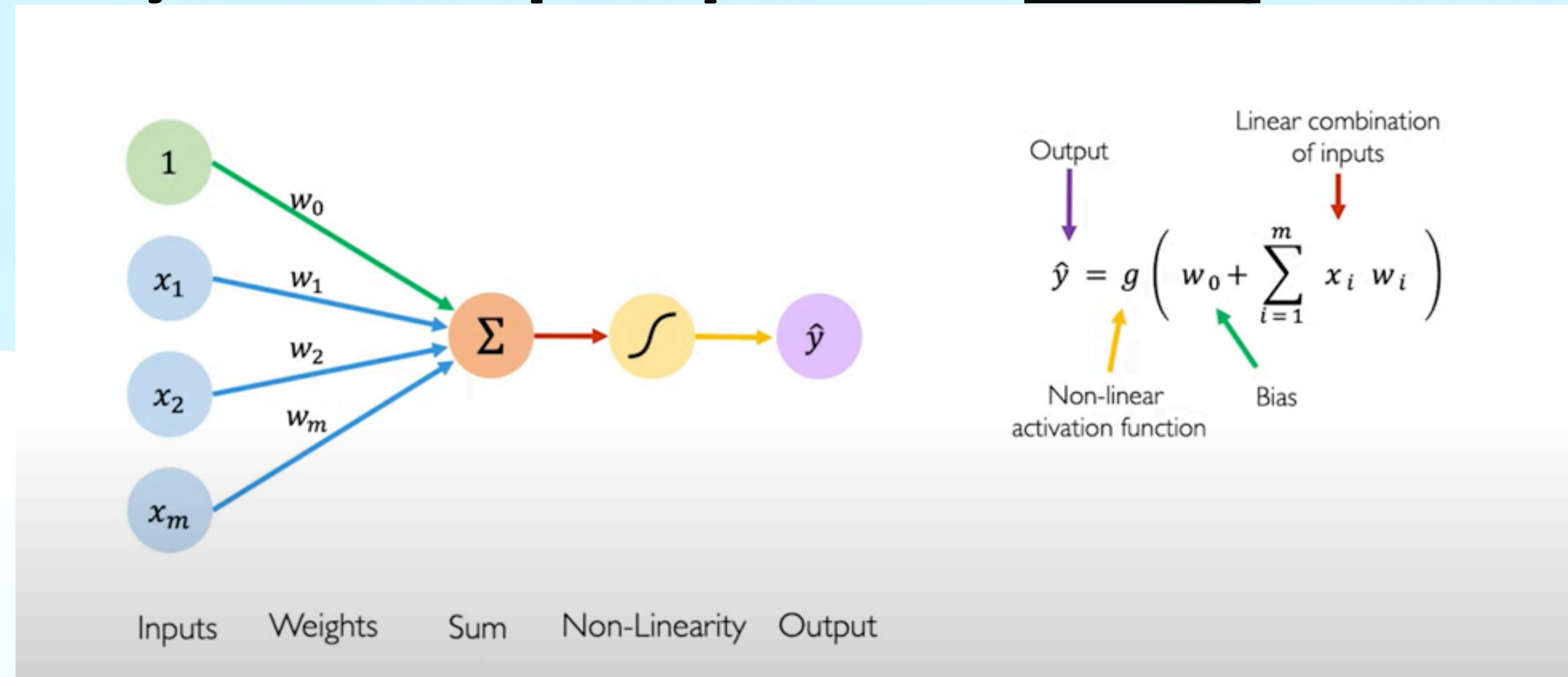
One Layer Perceptron

One Layer Perceptron $\hat{y} = g(W*x + w_0)$ (w_0 constant bias)

Die Werte der Matrix W werden als Gewichte $w_0 \dots w_n$ bezeichnet

Als Aktivierungsfunktion / Nicht-Linearität g wird oft ReLU, TanH, Sigmoid oder 'Step' verwendet.

Das Ergebnis eines Perceptron Layers wird als "Aktivierung" bezeichnet



Um die Aktivierung \hat{y} vom gewünschten output y (training-target) zu unterscheiden schreibt man oft einen Strich / Hut drüber, Elemente-weise \hat{y}_i vs y_i (y true labels vs y predicted)

One Layer Perceptron

Idee: reagiere auf verschiedene Eingaben (input x) mit verschiedener Ausgabe (output y)

Dazu wird jedem Teil der Eingabe (Element im Vektor) ein Gewicht zugeordnet und multipliziert:

$x_1 * w_1$ w_1 gross $\Rightarrow x_1$ ist wichtig
 $x_1 * w_1$ w_1 klein $\Rightarrow x_1$ ist unwichtig

diese gewichteten Elemente werden dann zusammen addiert:

Summe = $x_1 * w_1 + x_2 * w_2 + x_3 * w_3 \dots$

Um etwas mehr Flexibilität zu haben gibt man noch ein Gewicht unabhängig vom input hinzu:

Summe = $x_1 * w_1 + x_2 * w_2 + x_3 * w_3 \dots + \underline{1 * w_0}$

Das ist der sogenannte Bias w_0

Zusammen $\sum x_n w_n$ mit $x_0 = 1$

"Samme die Eingaben und gewichte sie"

One Layer Perceptron

Idee: reagiere auf verschiedene Eingaben (input) mit verschiedener Ausgabe (output)

Die Summe $\sum x_n w_n$ nennt sich "lineare" Schicht.

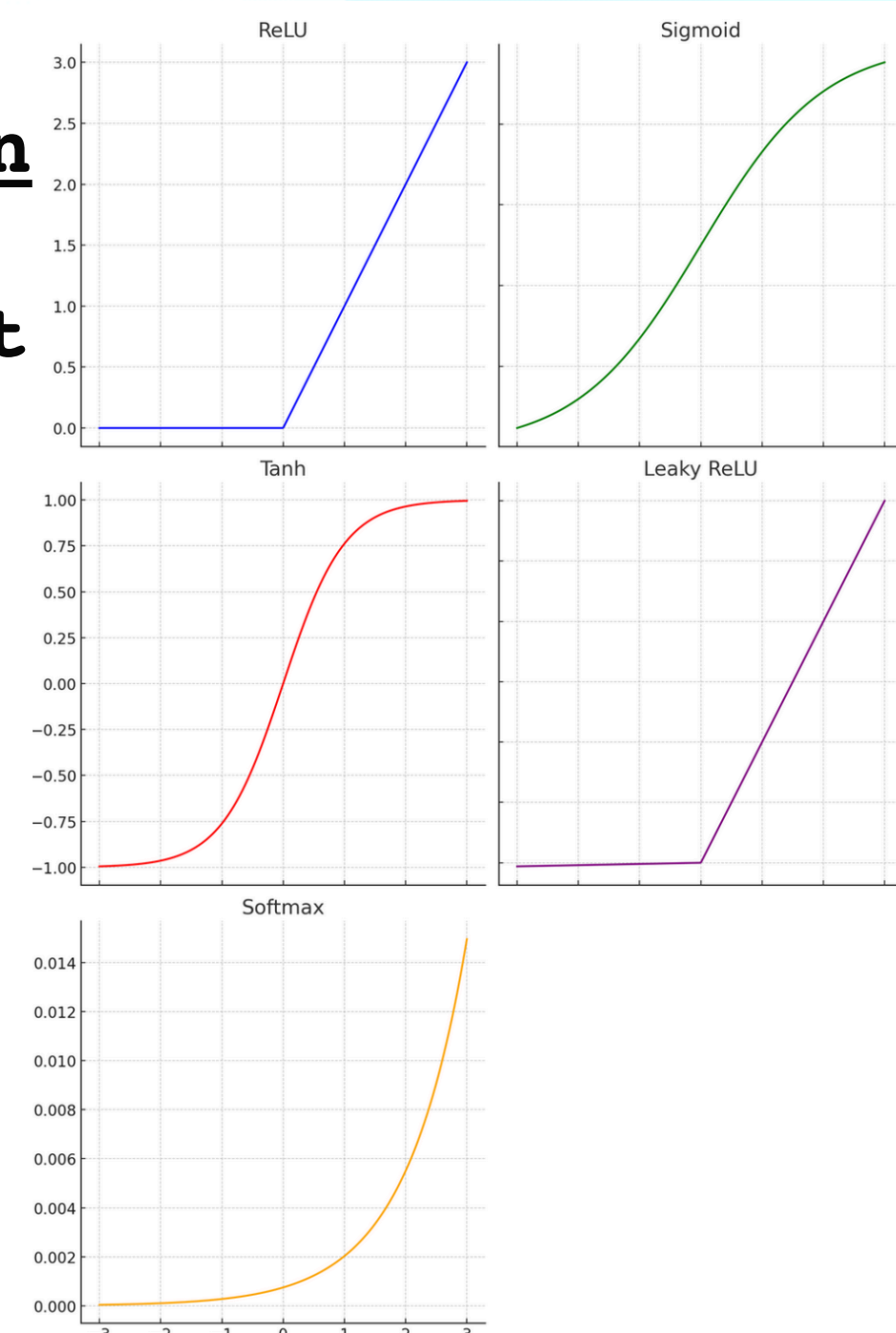
Das kommt daher dass jeder Parameter $x_n w_n$ wie eine lineare Gerade ist.

Typischerweise wird nach der Aufsummierung noch geschaut, ob das Ergebnis über einen gewissen Schwellwert kommt, also ob das Perceptron 'feuern' soll oder nicht, beziehungsweise wie stark.

Dies geschieht mit der sogenannten Nichtlinearität oder Aktivierungsfunktion

Nicht-linear heisst, dass die Aktivierungsfunktion einen Knick hat oder krumm ist

💡 mit einer Geraden allein lässt sich kein Schwellwert modellieren



One Layer Perceptron

Idee: reagiere auf verschiedene Eingaben (input) mit verschiedener Ausgabe (output)

Die Summe $\sum x_n w_n$ nennt sich "lineare" Schicht.

Wenn man den input x und die Gewichte w als Vektor betrachtet dann gibt es eine andere mathematische Schreibweise hierfür:

$x \circ w = \sum x_n w_n$ Skalarprodukt "dot-product"

One Layer Perceptron

$w \circ x = \sum w_n x_n$ Skalarprodukt "dot-product" 1-d output

in python:

```
import numpy as np

# Example vectors
x = np.array([1, 2, 3]) # x[0] ist 1 und gehört zum Bias w[0]
w = np.array([4, 5, 6])

# Scalar product
dot_product = np.dot(x, w)
print(dot_product)
g = np.tanh # or np.sigmoid Activation function
y = g(dot_product)
```

One Layer Perceptron als Matrix multiplikation

$W * X = \sum W_n^i X_n$ Skalarprodukt "matrix-product" n-d output

```
# Example vectors
```

```
W = np.array([[4, 5, 6], [7, 8, 9]])
```

```
x = np.array([1, 2, 3])
```

```
# Matrix product
```

```
#  $W * x = [w1 \cdot x, w2 \cdot x]$  in numpy :  $W @ x$ 
```

```
# Scalar product
```

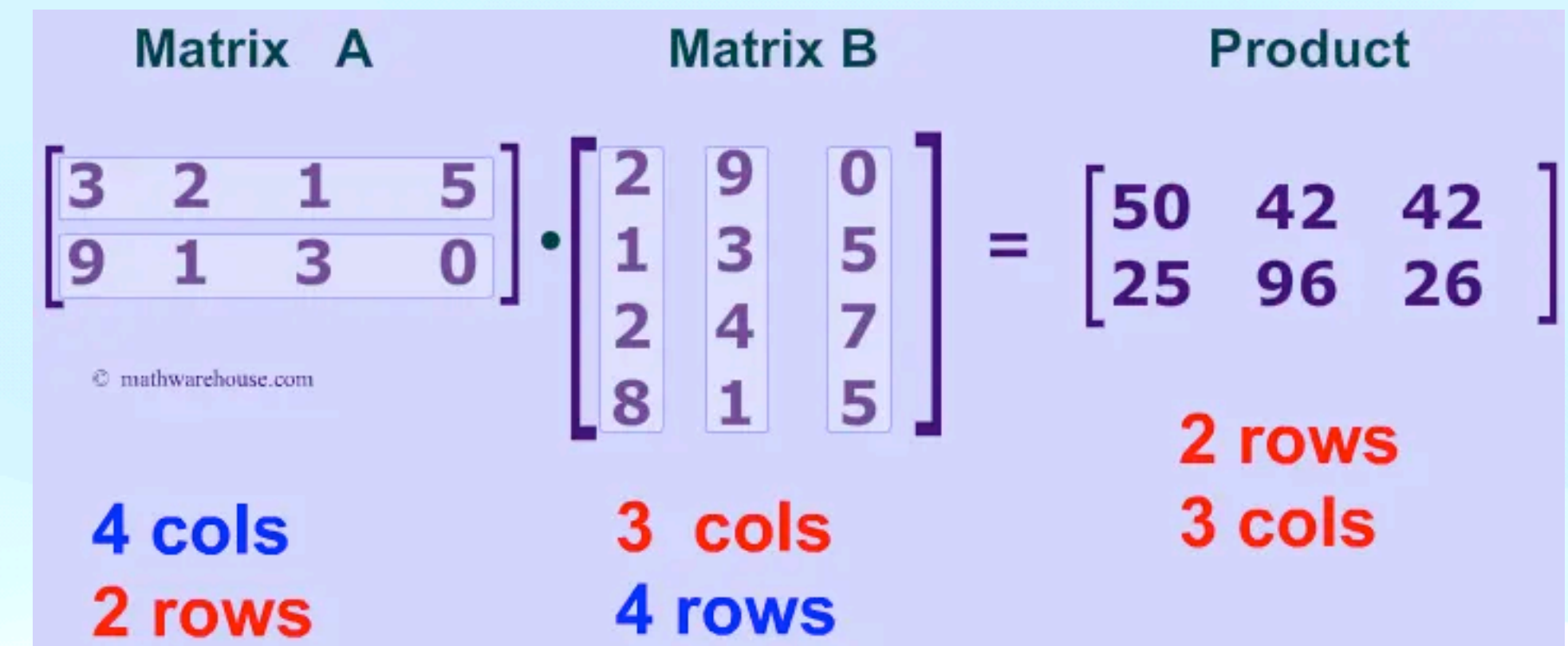
```
matrix_product = np.matmul(W, x) # or [np.dot(w1, x), np.dot(w2, x)]
```

```
print(matrix_product)
```

```
g = np.tanh # or np.sigmoid Activation function
```

```
y = g(matrix_product) # broadcast "erweitert Aktivierungsfunktion auf alle Elemente"
```

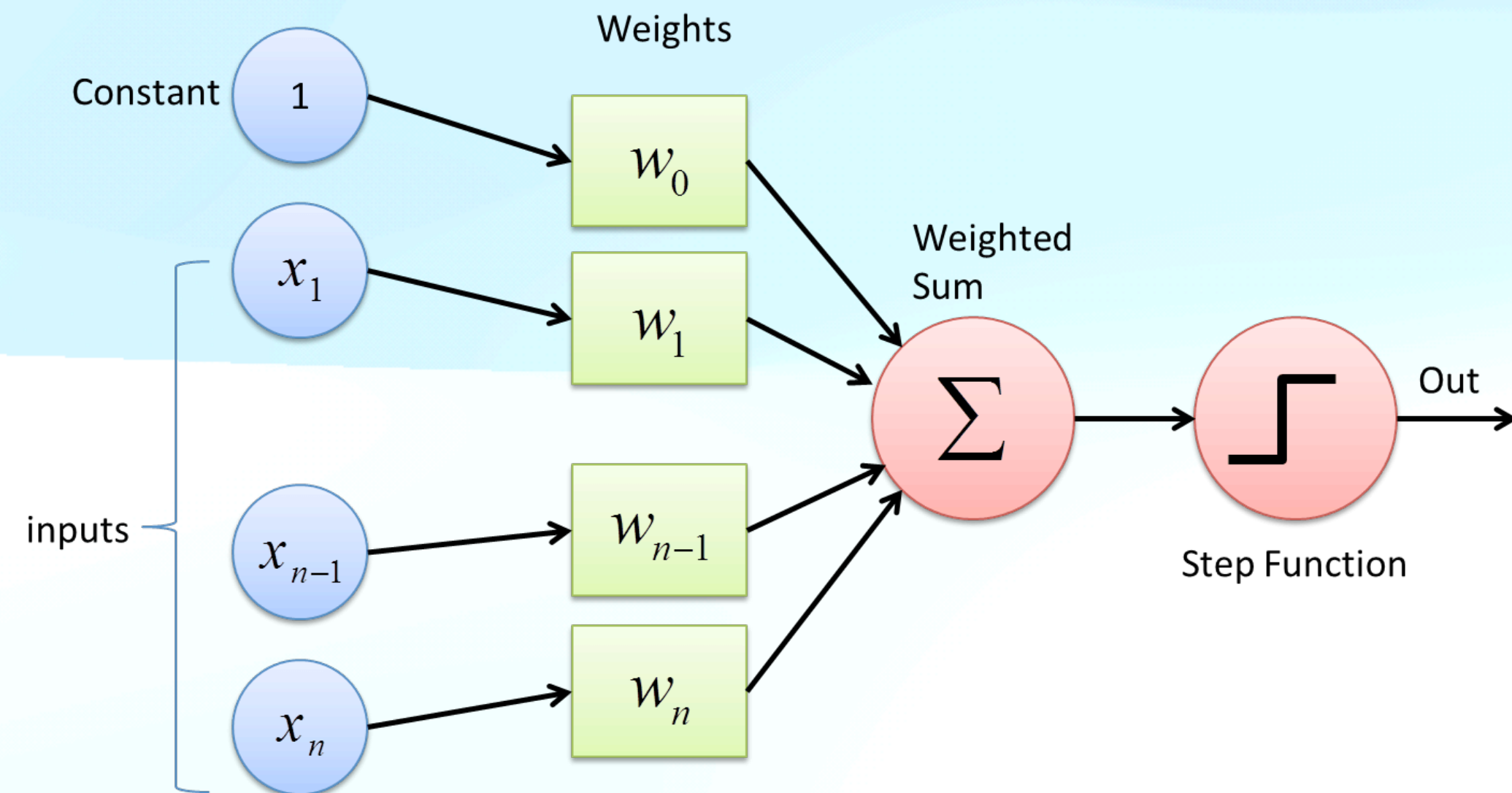
```
print(y)
```



One Layer Perceptron

One Layer Perceptron $y = g(M \cdot x + w_0)$ (w_0 constant bias)

Analogie zu Neuron: Die Gewichte verbinden andere Ausgänge (Axone) und werden ähnlich in Synapsen zusammenaddiert. Die "Aktivierungsfunktion" entspricht ungefähr den binären Spikes von Neuronen. Ein Neuron allein ist natürlich viel komplizierter, aber im Ensemble kann unser Perceptron vielleicht Grundfunktionen vom Neuron abbilden.



Aktivierungen

Sinn von Aktivierungsfunktionen:

- Unser "Neuron" "feuert" wenn ein Schwellwert erreicht ist
- (Optional) "Neuron" kann "gesättigt" sein (sigmoid/tanh)

$\text{Relu}(x) = \max(0, x)$ Rectified Linear Unit

$\text{Step}(x) = x > 0$

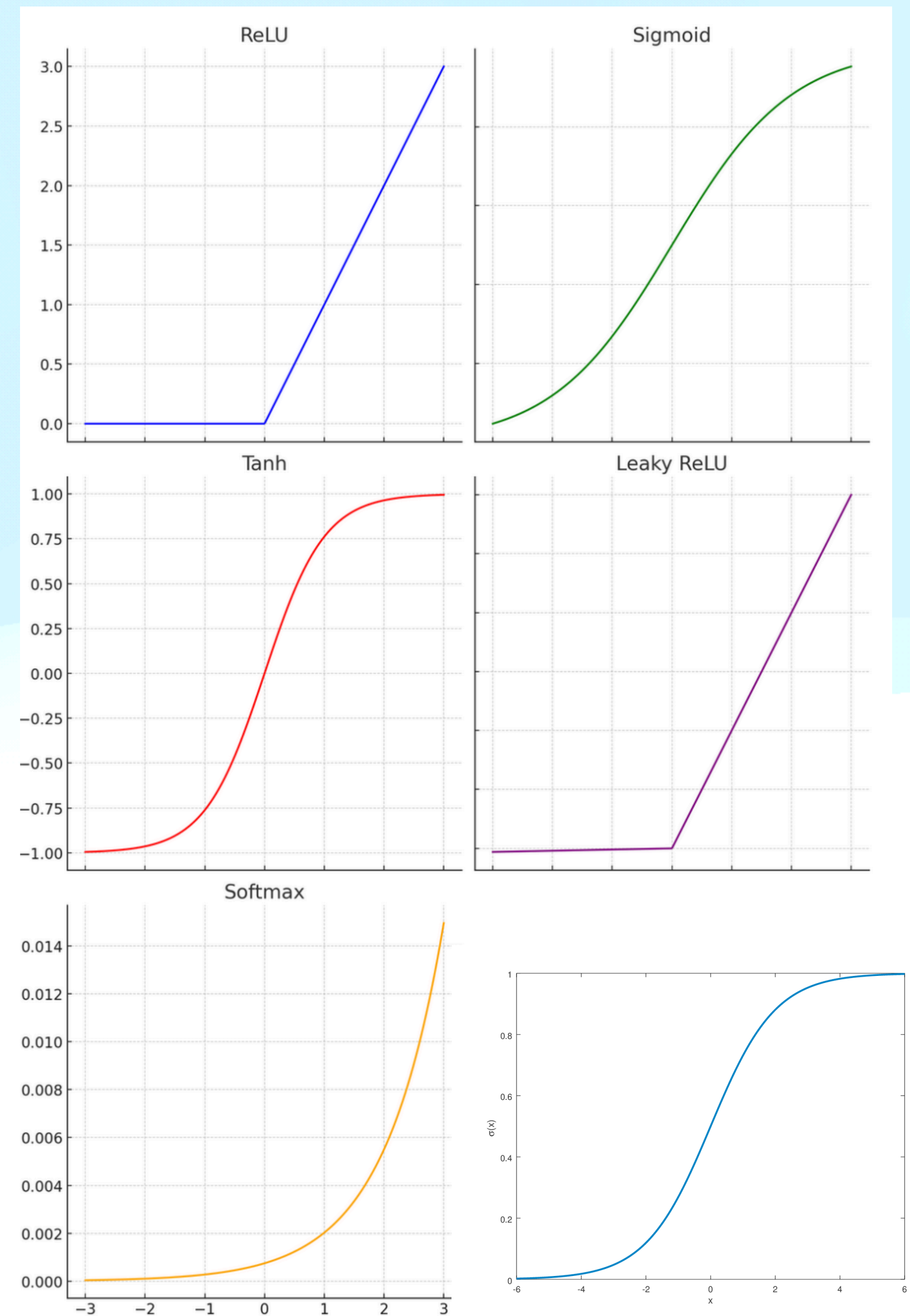
$\text{Sigmoid}(x) = 1 / (1 + e^{-x})$

$\text{TanH}(x)$

$\text{LeakyRelu}(x) = \max(\epsilon x, x)$ $\epsilon \approx 0$

Warum ist Step eine schlechte Aktivierungs Funktion?

Auch Link-Funktion



Fehler

Um zu messen wie gut unser Netzwerk die Trainingsdaten modelliert, messen wir entweder die Genauigkeit (**accuracy**), oder den so genannten **Fehler**.

Es gibt verschiedene Arten den **Fehler** (**error, loss, cost, risk**) eines Netzes zu berechnen.

$\text{abs}(\mathbf{y} - \bar{\mathbf{y}})$

Neben den **standard** Methoden (absolute **Differenz**, $\sum (\mathbf{y}_i - \bar{\mathbf{y}}_i)^2 / N$ mean square error MSE)

Gibt es auch etwas fortgeschrittene Berechnungen:

- Cross Entropy (next slide)
- Custom Errors (next slide)

💡 Während die **Genauigkeit** stets in einem Bereich **zwischen 0 und 100%** liegt, ist der **Fehler eher abstrakt** und kann beliebige Werte annehmen allerdings mit der Eigenschaft dass er monoton fallen sollte. Geringere Fehler entsprechen idealerweise höherer Genauigkeit.

Generell kann man das Vorzeichen beliebig vertauschen also wenn man eine Messgröße hat welche man nach oben optimieren möchte ist es das selbe als wenn man die negative Messgröße nach unten optimiert.

⚠️ Es ist Standard fallende Größen zu verwenden, also Verlust/**loss** zu **reduzieren**.

Lernen

Die Begriffe **lernen** und **trainieren** bedeuten im Zusammenhang von Deep Learning das **Anpassen der Modell Parameter** ("Gewichte") so dass eine Zielfunktion optimiert wird zum Beispiel erhöhen der Genauigkeit oder **Reduzierung des Fehlers**.

Wie Lernen?

Wie können wir unser Netzwerk **lernen** lassen / **trainieren**?

Anpassen der Modell Parameter ("Gewichte") so dass eine Zielfunktion optimiert wird zum Beispiel erhöhen der Genauigkeit oder **Reduzierung des Fehlers**.

Naiv / Zufällig: (blind) an den Gewichten wackeln und schauen ob der Fehler kleiner wird.

Mathematisch: Berechnen, welche Änderung der Gewichte zu kleinerem Fehler führt.

Stochastisch: dem berechneten Wert ein Stück weit vertrauen, aber nur ein wenig an den Gewichten wackeln um nicht 'übers Ziel hinaus zu schießen' (siehe Überanpassung Tag 4)

Aufgaben Perceptron

Berechne für den einfachsten Fall $n=2$ die Aktivierung bei 'zufälligen' initial Gewichten $w_0 = w_1 = w_2 = 1/2$ und input vector $(1,1)$

Das gewünschte Ergebnis der Aktivierung soll 1 sein. Wie könnte man die Gewichte anpassen?

Das gewünschte Ergebnis der Aktivierung für input $(1,1)$ soll 1 sein.

Das gewünschte Ergebnis der Aktivierung für input $(1,0)$ soll 1 sein.

Das gewünschte Ergebnis der Aktivierung für input $(0,1)$ soll 1 sein.

Das gewünschte Ergebnis der Aktivierung für input $(0,0)$ soll 0 sein.

Welche Funktion lernt unser Netz?

Wie könnte man die Gewichte konsistent anpassen so dass alle Bedingungen erfüllt sind?

Nachdem obiges per Hand berechnet wurde, vergleiche mit torch:

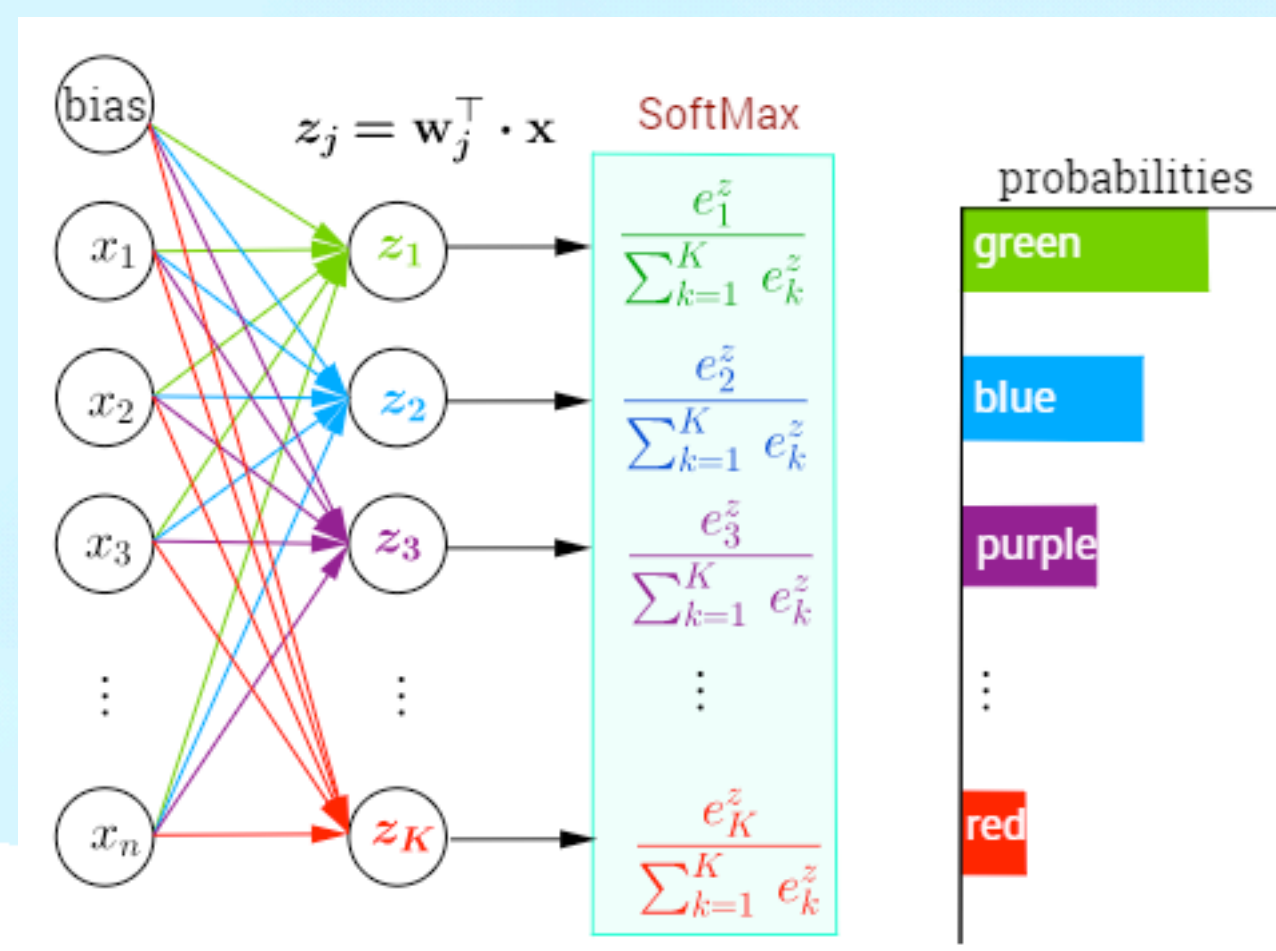
```
model = torch.nn.Linear( 2 , 1, bias=True) # most trivial
```

Wie könnte man Fehler und Optimierer selber schreiben?

Warum kann XOR nicht mit einem einfachen Linear gelernt werden?

Fehler

CrossEntropy für 'Klassen' bewertet zB wie unsere Zahlen 0..9 in verschiedenen **Töpfen** landen. Hierzu werden am Ende alle Töpfe durch die Summe der Aktivierungen geteilt, was man **Softmax** nennt. Das erinnert an Wahrscheinlichkeitsvariablen, deren **Summe stets 1** ergeben.



10 Zielvariablen haben Zustände 0/1, praktisch dazwischen [0,1] (ja ist Ziffer '7', oder nicht, oder vielleicht)

Eine weiche (arg)maximum Funktion schaut, in welchem der Töpfe der größte Wert ist.

Ein gutes Maß zum Schauen, ob die Werte in unseren Töpfen p_i unseren gewünschten Werten y_i entsprechen ist

Normal: $y_i * \bar{y}_i$ solange bei einem Zustand beide 1 sind, kommt 1 raus, was gewünscht/anzustreben ist.

wir formulieren es als Minimierungsproblem um:

CrossEntropy (logistische Verlustfunktion) $-\sum y_i \log(p_i)$

(log ist monoton, man verliert nur Skalierung, lässt sich leichter ableiten!)

Der Cross-Entropy-Verlust ist effektiv, weil er große Strafen für sicher vorhergesagte falsche Antworten verhängt, was die Konvergenzgeschwindigkeit des Lernprozesses beschleunigt.

Custom Errors

Neben der Nähe zum gewünschten Ergebnis kann man auch andere **(meta)Kriterien** in den "**Fehler**" mit einbauen:

besserer_Fehler = Fehler + a * custom_Fehler (z.B. $7 \neq 2$ $1 \neq 7$)

In der Tat kann man mit einer eigenen Fehler Funktion das Verhalten vom Neuronalen Netz fein steuern ("Belohne glatte Kurven" "Bestrafe krakelige Schrift")

Dazu genügt es oft verschiedene **Maße** zu **addieren**, mit **Gewichtungsfaktor** (nachdem sie mit einem Faktor auf ähnliche **Größenordnungen** gebracht wurden)

Der Optimierer reduziert den Fehler dann langfristig so, dass alle Kriterien ähnlich gut erfüllt werden.

Back Propagation

Wenn wir einen Fehler berechnet haben können wir sehen welche **Gewichte** den größten **Einfluss** auf den **Fehler** hatten und die Gewichte dementsprechend anpassen. Präzise funktioniert das über die sogenannte **Backpropagation**, also **Rückwärtsausbreitung** oder Rückverfolgung des Fehlers.

Mathematisch sagt uns die Ableitung wie stark die Sensitivität des Fehlers bezüglich kleiner Änderungen dieser Gewichte ist. Um die Ableitung zu berechnen erinnern wir uns an die einfache

Kettenregel: $g(f(x))' = g'(f(x)) * f'(x)$

Also egal wie tief wir im Netz sind, welches für einen input x einen output y berechnet, wir müssen nur mehrfach $y=g(f(x))$ betrachten.

Wie stark die Gewichte tatsächlich angepasst werden hängt vom Optimierer ab, welcher zum Beispiel die **learning rate** mit der Änderung multipliziert.

Back Propagation

Backpropagation (Rückwärtsausbreitung) ist ein zentrales Konzept in der Funktionsweise tiefer neuronaler Netze und bezieht sich auf den Prozess, durch den Gewichte in einem Netzwerk angepasst werden, um den Fehler zu minimieren.

Es nutzt den **Gradientenabstieg** (Gradient Descent),

um den **Fehler Schritt für Schritt zu reduzieren**,

indem es die **Ableitungen der Verlustfunktion** (Loss Function) bezüglich der Gewichte berechnet.

Backpropagation ist ein Verfahren, das den Fehler (Loss) in einem neuronalen Netz rückwärts durch die Schichten propagiert, um die Gewichte anzupassen. Dies geschieht in zwei Schritten: Forward Propagation und Backward Propagation.

- **Vorwärtsausbreitung** (Forward Propagation):

Die Eingabedaten werden durch das Netzwerk geleitet, um eine Vorhersage zu generieren.

Der **Fehler (Loss)** zwischen der Vorhersage und dem tatsächlichen Zielwert wird berechnet.

- **Rückwärtsausbreitung** (Backward Propagation):

Der Fehler wird rückwärts durch das Netzwerk geleitet, indem die Gradienten der Verlustfunktion in Bezug auf die Gewichte berechnet werden.

Diese Gradienten werden verwendet, um die **Gewichte zu aktualisieren**, sodass der Fehler bei der nächsten Vorhersage reduziert wird.

Back Propagation

Probleme der Kettenregel (gelöst!):

Vanishing / Exploding Gradient

<https://karpathy.medium.com/yes-you-should-understand-backprop-e2f06eab496b>

Mathematisch sagt uns die Ableitung wie stark die Sensitivität des Fehlers bezüglich kleiner Änderungen dieser Gewichte ist. Um die Ableitung zu berechnen erinnern wir uns an die einfache **Kettenregel**:

$$f(g(x)) = (f \cdot g)' = g' * f' \cdot g$$

$$y = f\left(\sum w_i * x_i + w_0\right)$$

$$g(x) := \sum w_i * x_i + w_0$$

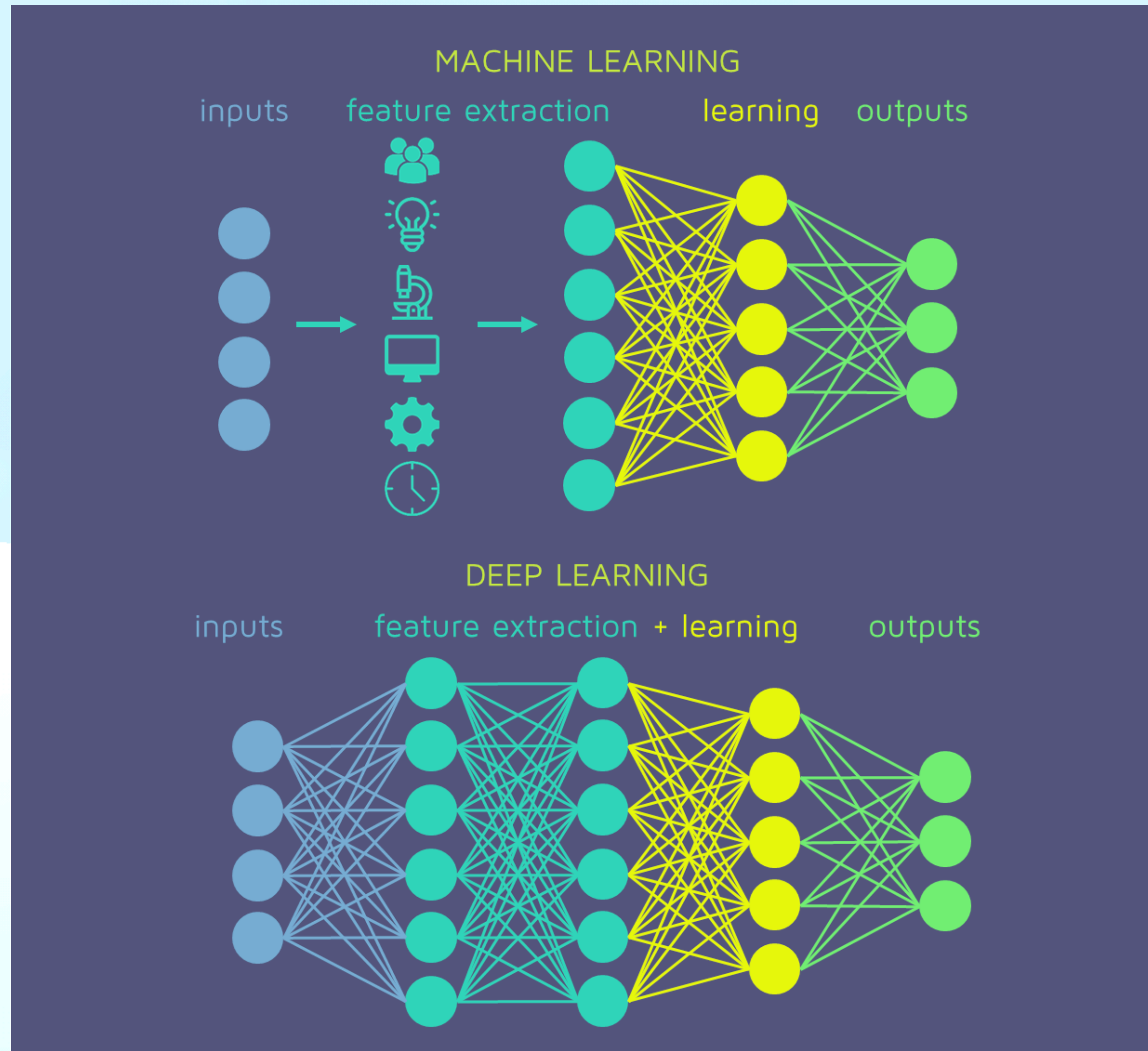
$$y = f(g(x))$$

f unsere Aktivierungsfunktion $\max(0, x)$

$$f'(x) = 0 \text{ für } x < 0 \quad 1 \text{ für } x \geq 0$$

$$g'(x_i) = w_i$$

Deep Learning



Erfolg 2010...

Gelöste Probleme der AI:

Vanishing / Exploding Gradient heisst numerische Instabilitäten

Overfitting (no more early stopping)

Lösungen:

Regularisierung

- **Dropout** (zufällige Deaktivierung von Perzeptronen, Verteilung des Wissens "Multiplikatoren")
- **Residual** Connections (Schicke Signal unverändert zu tieferen Schichten "Überbleibsel der original Informationen des Input")
- Batch Normalization (Mean/Zentrum wird subtrahiert und Varianz dividiert)
- Convolutional Networks
- Architecture does not matter (**mehr compute + mehr trainingsdaten**) Andrew Ng
- Precision does NOT matter
 - Last paper: 1.8 bit averaged information is **better(more efficient) than** f16, f32
 - Quantisierte LLMs **4 bit** "floats" pseudo reals bit \approx Eigenschaft ja/nein
 - 7 GB 30 GB LLMs sind brauchbar für Spezialanwendungen auf eigenem Rechner

20 questions

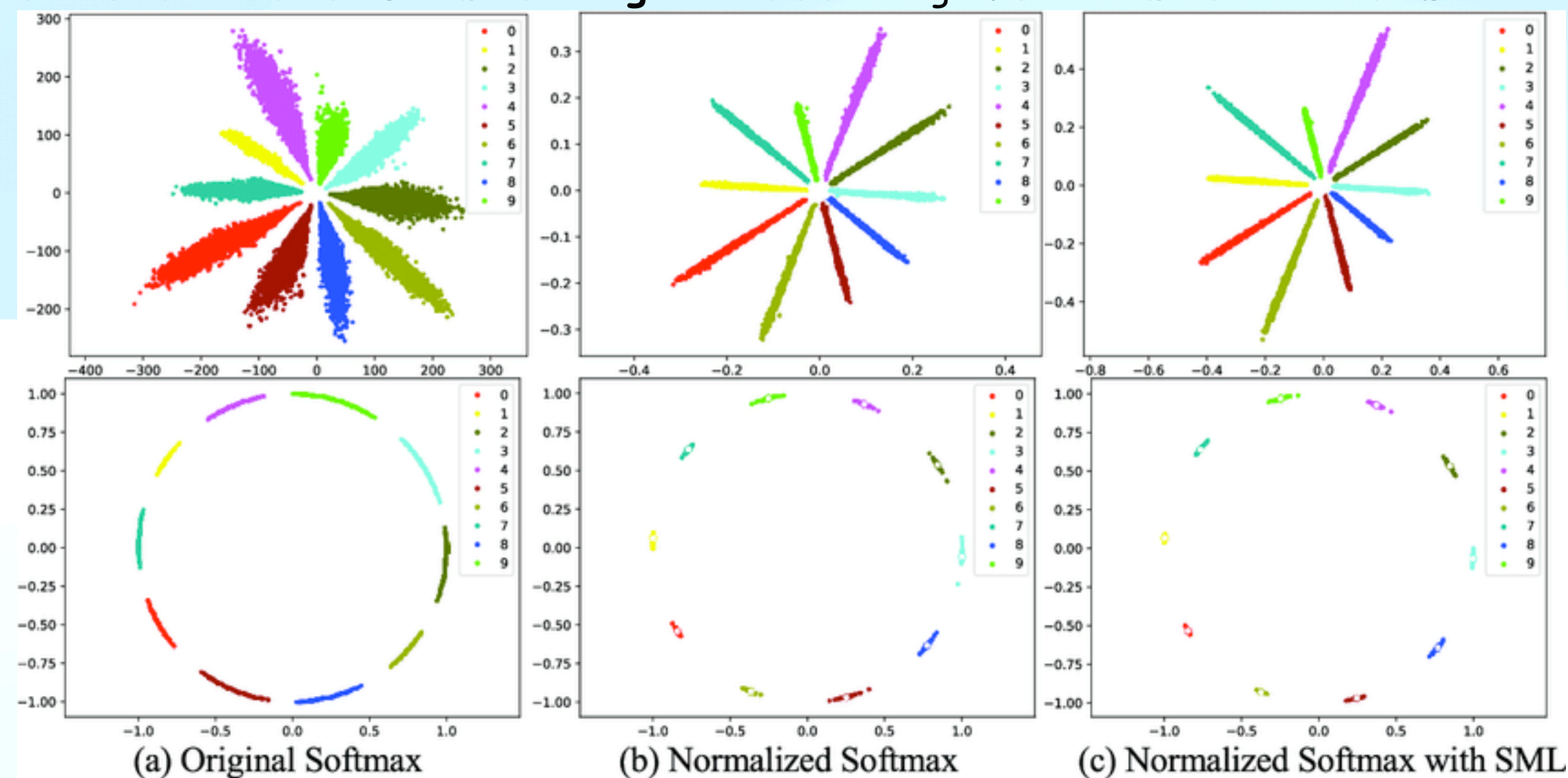
(Un/Semi-)Supervised

Wenn wir zu Daten die passenden Beschreibungen ("**Labels**") haben, spricht man von **Supervised learning**.

Oft hat man aber nur Roh Daten ohne Label.

Auch mit diesen lässt sich lernen, dann spricht man von **unsupervised learning**:

- **Kompression n-input n/x-schicht**
- **Generierung**: Das Modell kann lernen die Daten zu reproduzieren
- automatische **Clustering** Entdeckung von **Mustern und Strukturen** in den Daten.



Mischvarianten nennt man **SemiSupervised**, z.B. wenn **GPT** unvoreingenommen das ganze Internet durchliest (unsupervised) aber dann im letzten Durchgang noch Label bekommt, was 'wahrer' 'hochwertiger' 'sozial akzeptabler' ... Content ist.

Verfahren ohne Ableitung

Pseudo Ableitung

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

wie findet man blind das Minimum

gehe einen Schritt nach links oder rechts

rechts: Wert wird grösser => drehe um gehe doppelt so weit zurück

links: Wert wird kleiner => gehe weiter nach links

random descent

Etwas geplanter (und mathematischer) Minimum finden: Gradient descent

$f'(x)$ wenn die Steigung groß ist gehe vorsichtig kleine Schritte

$f'(x)$ wenn die Steigung klein ist gehe mutig große Schritte

$$y = f(x)$$

$$\text{Fehler} = |y' - y|$$

Forward-Forward-Network (von Hinton)

Nachteil: Funktioniert erstaunlicherweise, aber ist viel langsamer als Gradienten Backpropagation.

Bei Analogon oder Photonen Chips wieder relevant.

Problem: Elektronen und Photonen in einem Chip

AI Grundverständnis

Legendär verständlich sind die online Kurse von Andrew Karpathy (youtube/blogs)

Auf verschiedenem Niveau: Von Anfänger bis Experte (trotzdem recht verständlich)

Koryphäen:

- **Geoff Hinton**
- **Andrew Ng**
- Yann LeCun
- Yoshua Bengio
- Jürgen Schmidhuber (hat alles schon 1990 erfunden)

(es ist es Wert, von jedem das beste Paper / Video anzusehen)

Buch

Anhang B - Autodiff

Differenzierung von Hand

Unterstützung	IDE	Framework
VOLL	PyCharm	PyTorch
Partiell	VSCode	Tensorflow
Kaum	andere	andere

Breiter Output

Idee: reagiere auf verschiedene Eingaben (input) mit mehreren verschiedenen outputs gleichzeitig!

Dann hat man gleichzeitig nicht nur eine Summe, sondern mehrere.

💡 Jede Summe für sich ist genau so einfach wie vorher

Wir schreiben die Summen nur in Notation gemeinsam:

$$y_i = \sum w_{ij} x_j$$

Das hat den Vorteil, dass man die Summen als Matrizen Multiplikation zusammenfassen kann. Rechnerisch ist es aber identisch dazu ein paar Summen per Hand zu rechnen!

$$\underline{y} = \underline{W} * \underline{x}$$

Breiter Output

Wir schreiben die Summen nur in Notation gemeinsam:

$$y_i = \sum w_{ij} x_j$$

Ein Beispiel:

$$\text{input} = [1, 2]$$

$$\text{weights}_1 = [1, 2, 3] \quad \text{Gewichte für output } y_1$$

$$\text{weights}_2 = [2, 3, 4] \quad \text{Gewichte für output } y_2$$

$$y_i = w_i \circ x$$

Breiter Output

Die Aktivierungsfunktion wird auf jeden Eintrag im output Vektor einzeln angewendet:

Einfaches y ohne Aktivierungsfunktion

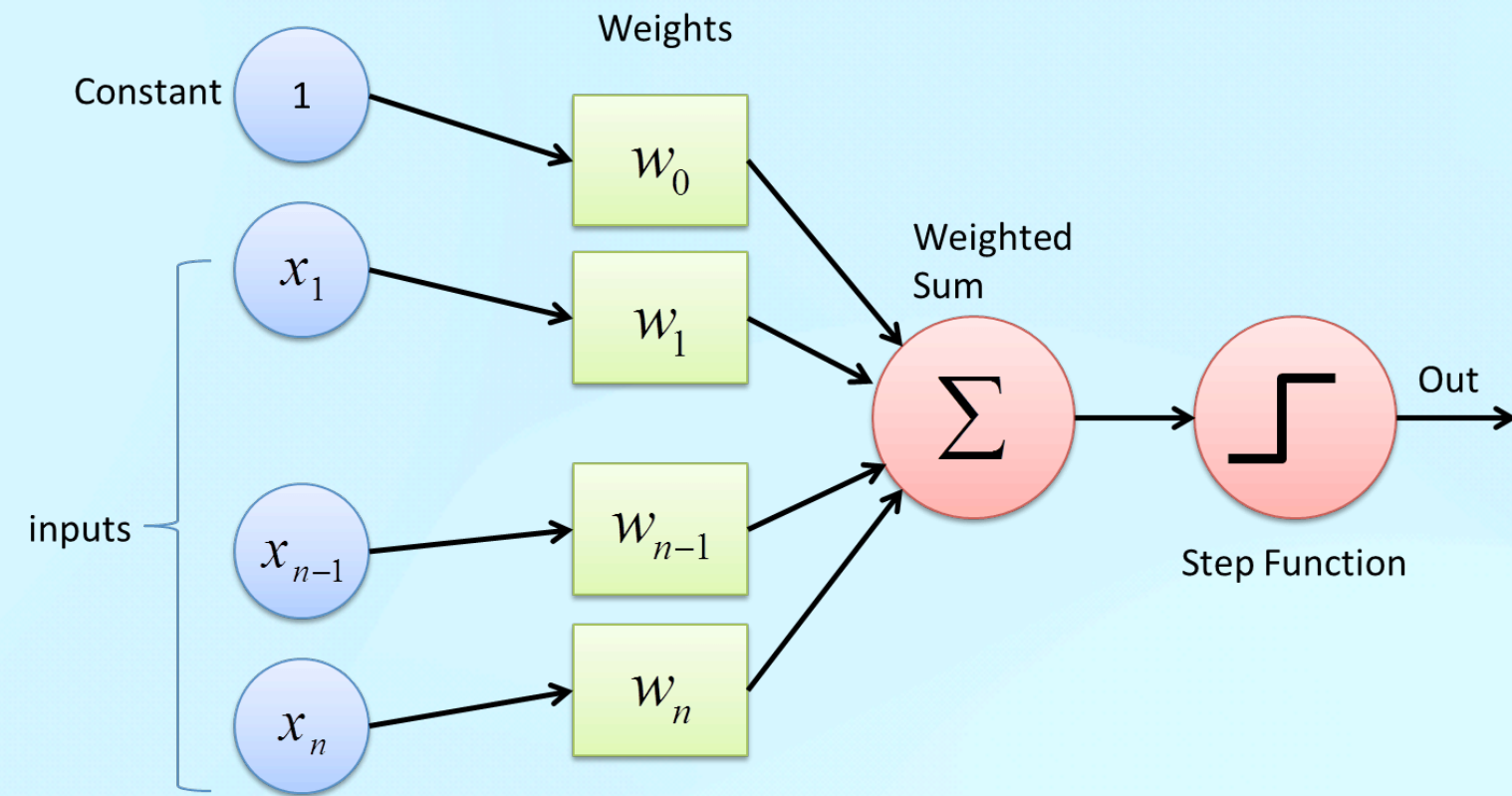
$$y = M * x$$

Besseres y mit Aktivierung nachgeschaltet " y^+ "

$$y^+ = g(M * x) = g(y_i)$$

Breiter Input

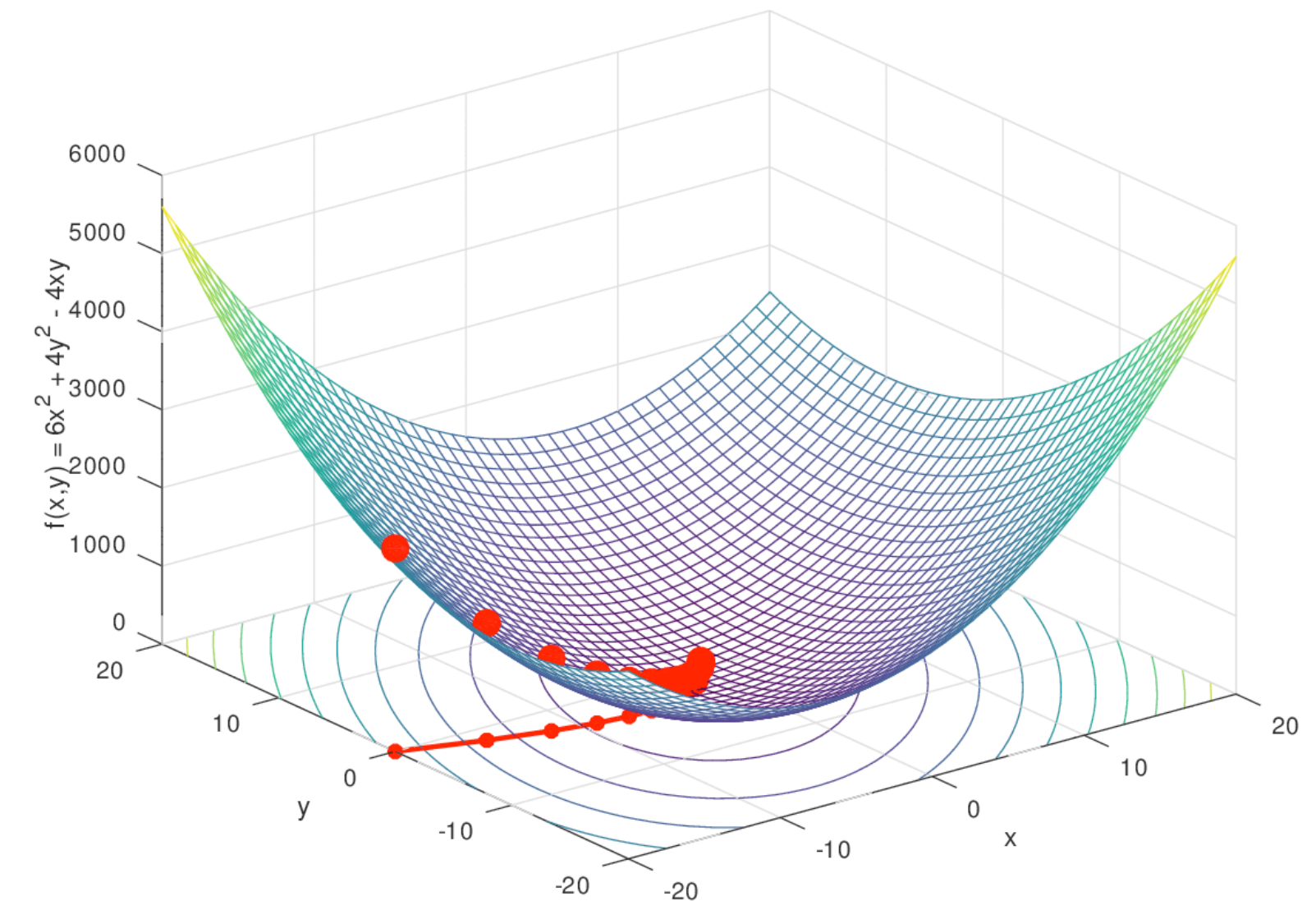
Wie wir gesehen haben bekommt jedes 'Neuron' den Input von mehreren 'Synapsen':



Damit ist unsere Aktivierung Funktion also mathematisch eine Funktion mehrere Variablen.
 $f(x_1, x_2, \dots, x_n) = \dots$

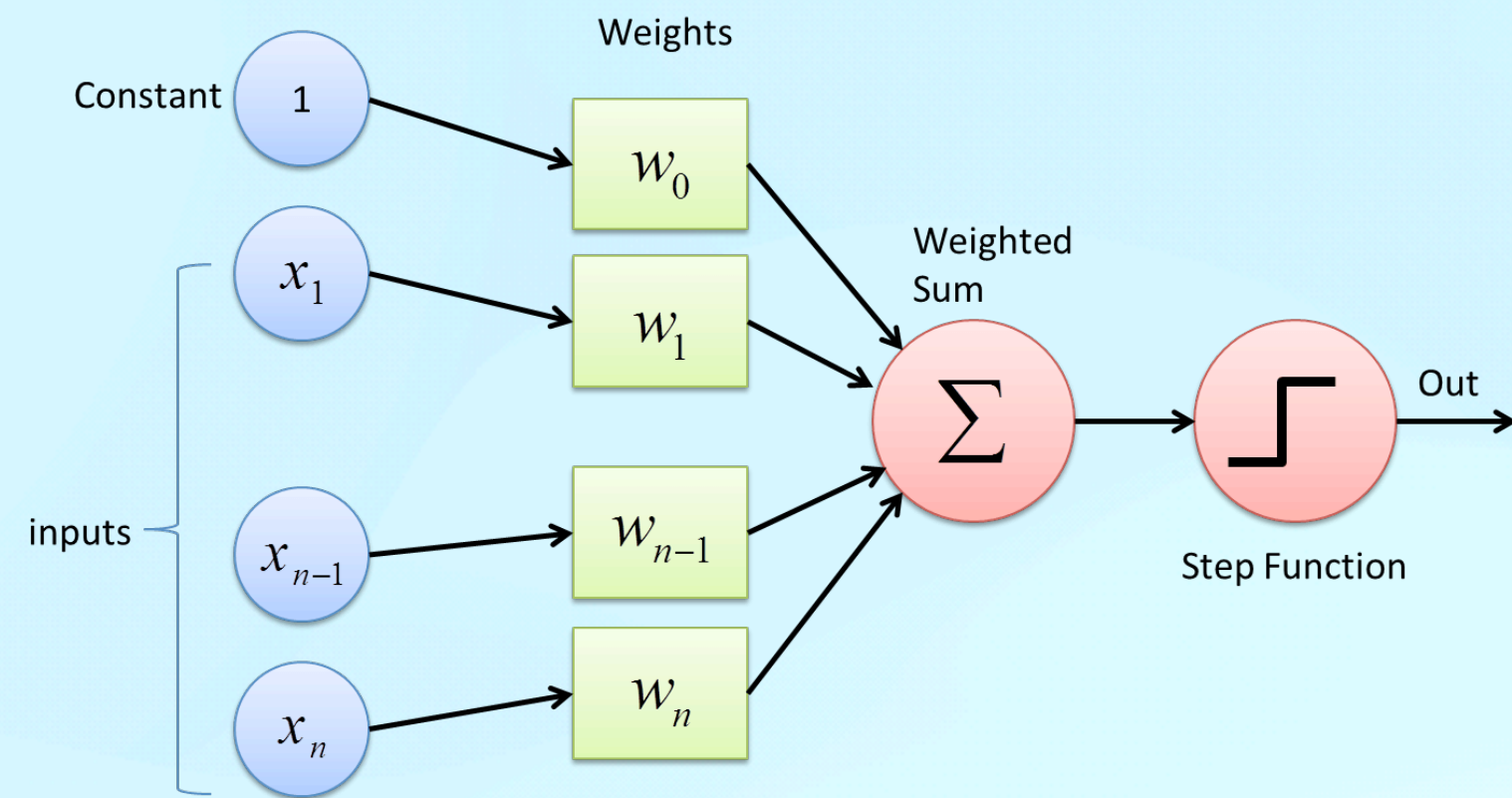
z.B. $f(x_1, x_2) = x_1 * x_2$ oder $f(x_1, x_2) = 2 * x_1 * x_2 * x_2$

Falls das jemand nicht in der Schule oder danach gesehen hat, solche Funktionen kann man sich vorstellen wie Gebirge über der Ebene:



Breiter Input : Batch

Batch: berechne die Aktivierung von mehreren Inputs gleichzeitig!



Stapel von Input Vektoren wird parallel ans Netz gegeben.

Damit $\mathbf{y} = \mathbf{W} * \mathbf{X}$ mit Inputs als Matrix gesammelt (statt $\mathbf{y} = \mathbf{W} * \mathbf{x}$ input vector \mathbf{x})

Gradient

Gradient = Ableitung von Funktion mehrere Variablen

Wenn eine Funktion von mehreren Variablen abhängt kann man sie in verschiedene 'Richtungen' ableiten, also für jede Variable die Ableitung bilden. Diese Ableitungen zusammengefasst als Vektor nennt sich Gradient. Man schreibt üblicherweise ∇f statt f' :

z.B.

$$f(x_1, x_2) = x_1 * x_2$$

$$\nabla f = [df/dx_1, df/dx_2] = [x_2, x_1]$$

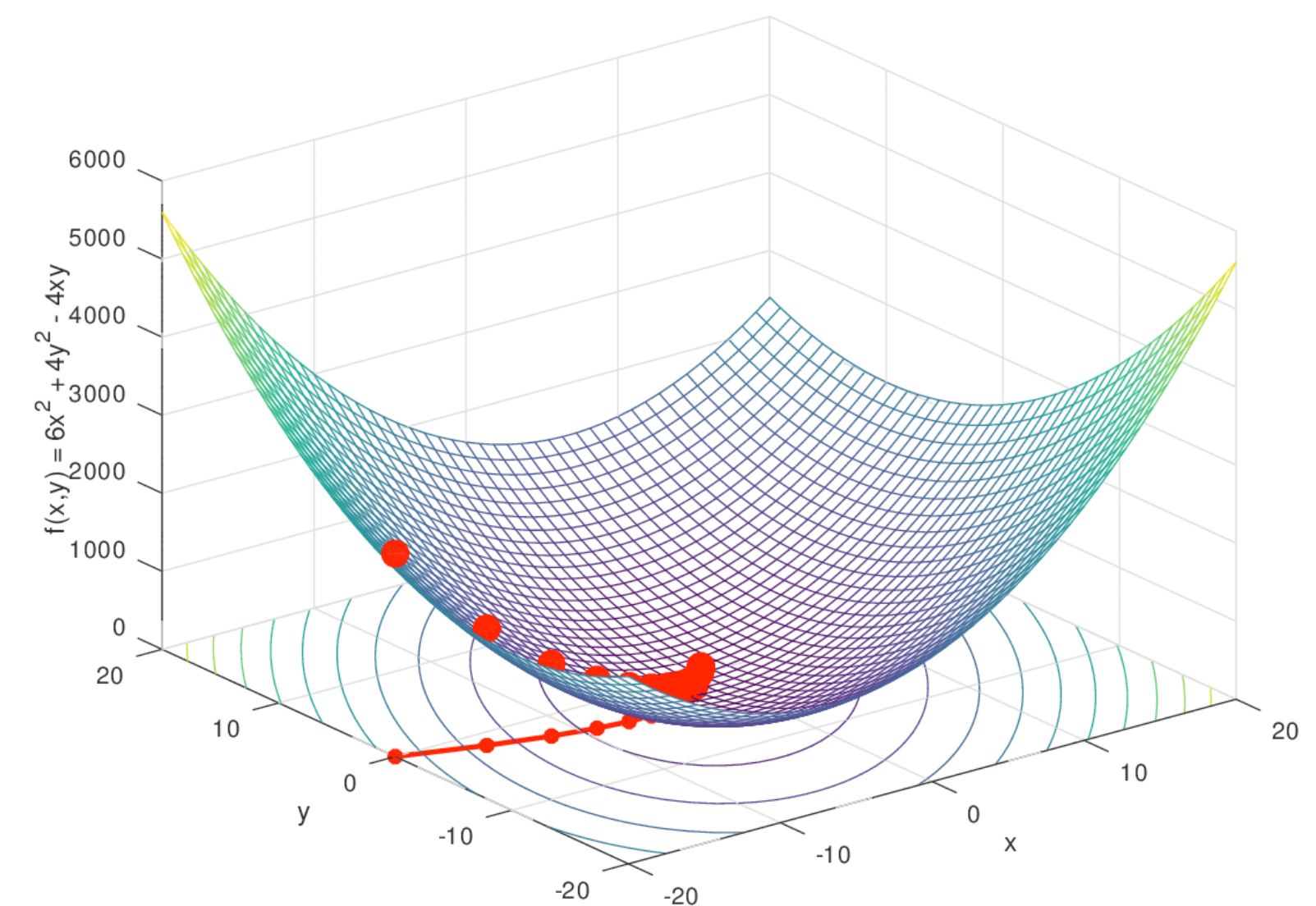
$$f(x_1, x_2) = x_1 + 2 * x_2$$

$$\nabla f = [df/dx_1, df/dx_2] = [1 + 0, 0 + 2] = [1, 2]$$

$$f(x_1, x_2) = 2 * x_1 * (x_2 * x_2)$$

$$\nabla f = [df/dx_1, df/dx_2] = [2 * x_2 * x_2, 2 * x_1 (2 * x_2)]$$

💡 Sowie die Ableitung als tangential Gerade bei einer einfachen Funktion den Grad des Anstiegs oder Abstieg anzeigt, zeigt der Gradient im Gebirge in die Richtung des stärksten Abstiegs beziehungsweise Anstiegs.



Gradient

So wie wir im Newton Verfahren als Beispiel bei der $f(x)=x^2$ Funktion zum Minimum gefunden haben, können wir mit dem Gradienten bei allgemeineren Funktionen die tiefste Stelle im Gebirge finden. Übertragen auf unsere Fehlerfunktion können wir also somit unseren Fehler minimieren.

Das ist die Grundidee von Backpropagation, im Detail erörtern wir das am Donnerstag unter dem Begriff Stochastik Gradient Descent (SGD).

