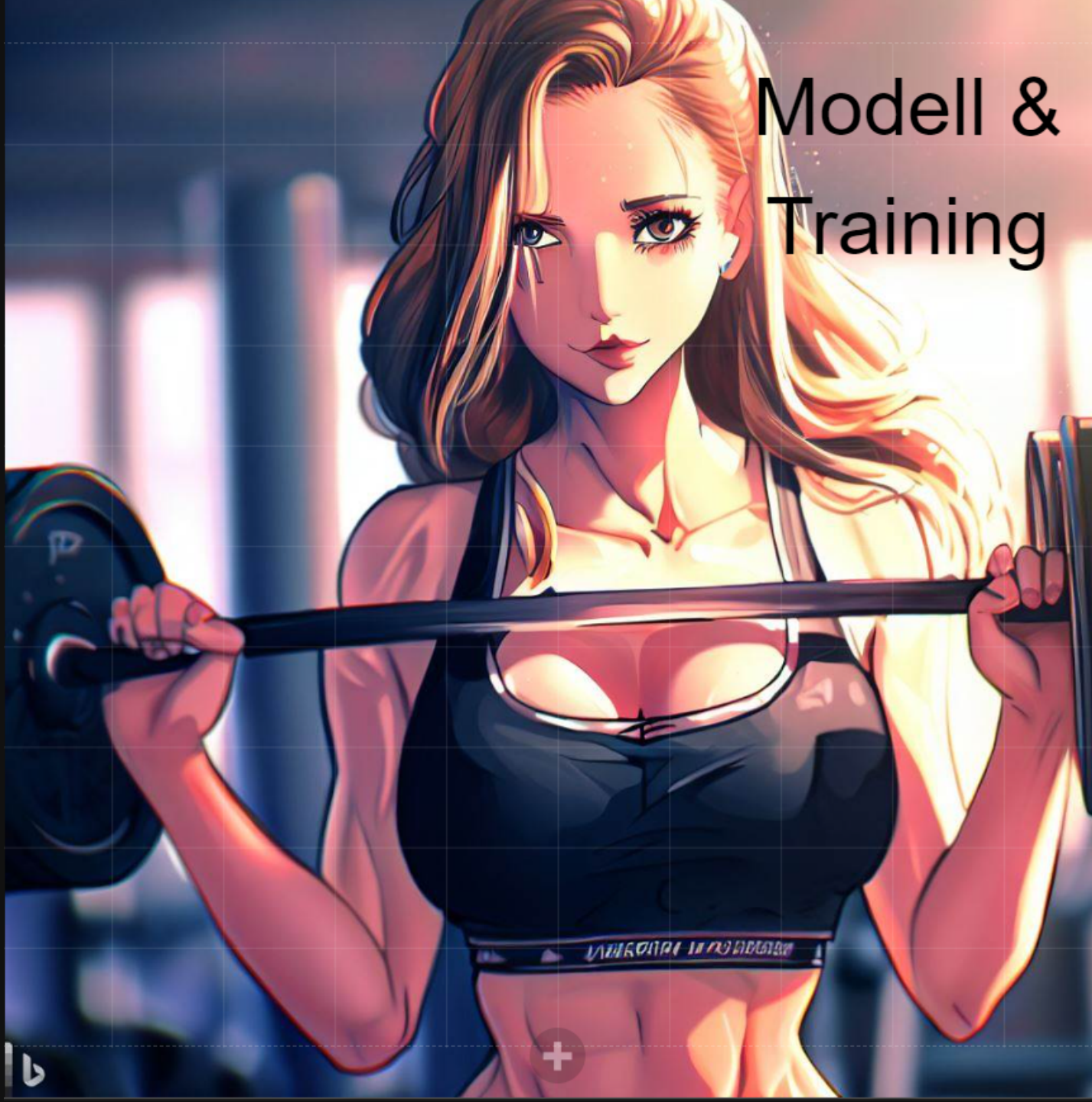


Modell & Training

<https://www.youtube.com/embed/O1ZQtTwPwYk?enablejsapi=1>



Deep Learning Modell

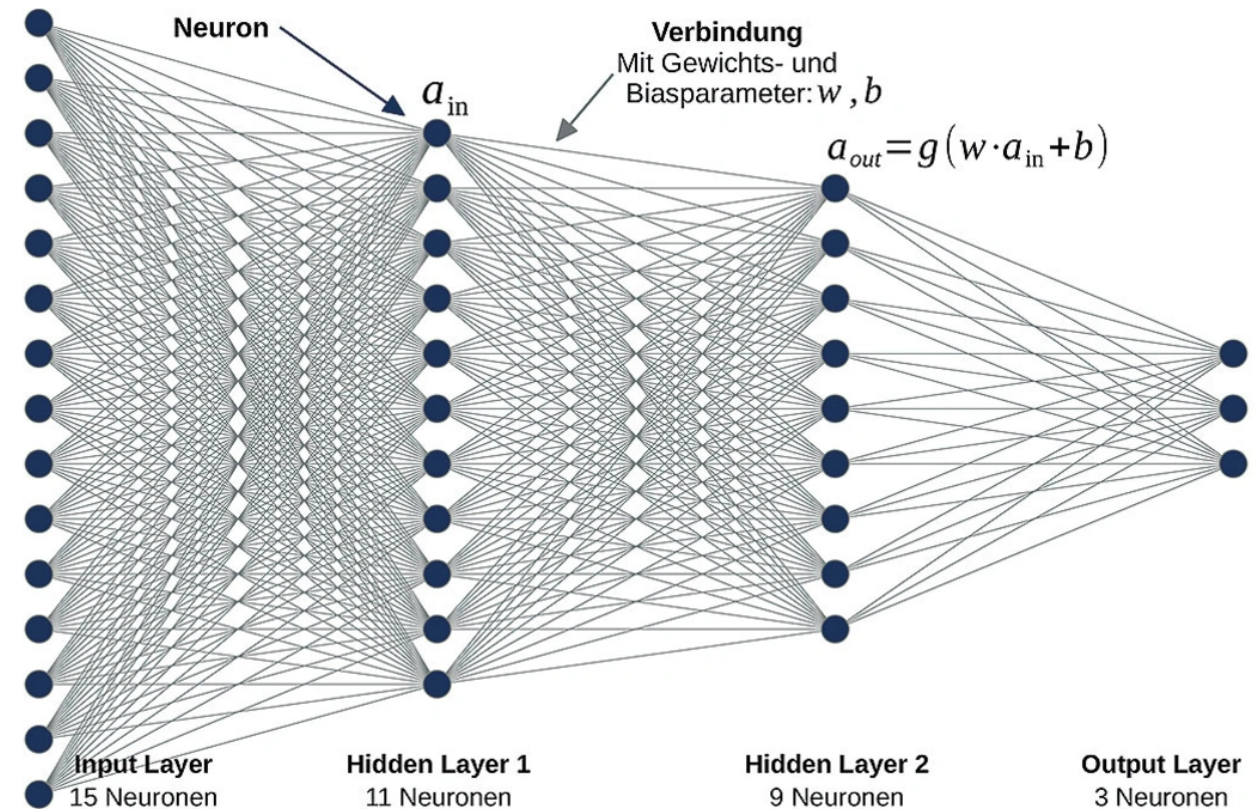
<https://www.youtube.com/embed/qFES8S9D8RM?enablejsapi=1>

Architektur: Neuronale Netzstruktur

Neuronales Netz = einfache Darstellung komplexer Rechnung einfacher Bausteine

LinReg mit Basisfunktionen aus LinReg mit Basis aus Linreg mit Basis aus ...

Modell = Architektur mit Satz von Parametern



mehr zu Neuronalen Netzen

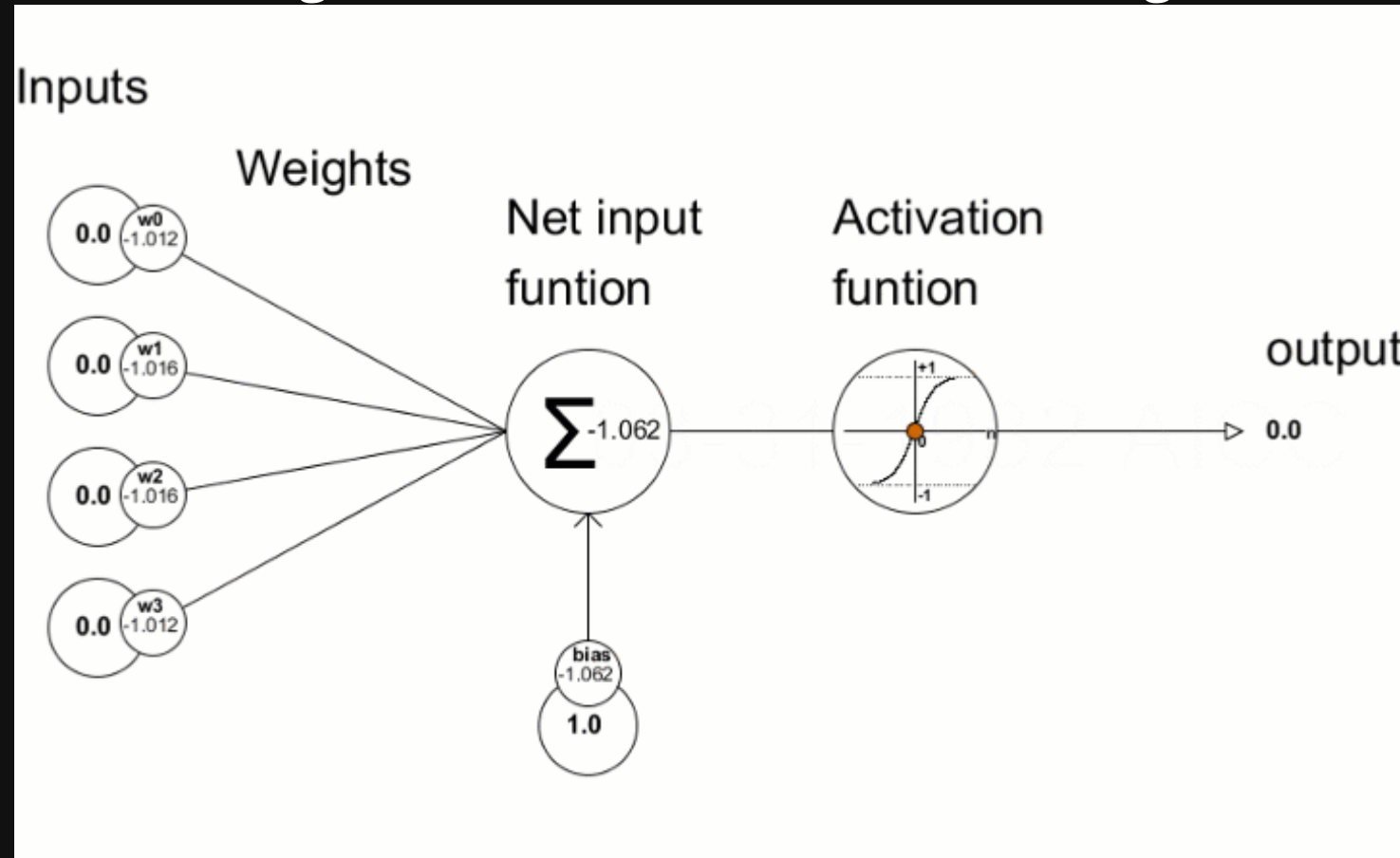
Layer

<https://www.youtube.com/embed/ynFZ3wBqsao?enablejsapi=1>

Layer = Level für Lineare Regression

Mehrere Knoten (Perceptronen)

Knoten = gewichtete Summe & Aktivierungsfunktion



Aktivierungsfunktion

Wofür Aktivierung?

Ohne Aktivierung: N Layer = 1 Layer

Lineare Kombination
von Linearkombinationen
ist eine Linearkombination

Nicht-Lineare Aktivierung
ermöglicht komplexe Modellierung

$$z_{1,j} = \sum_{i=1}^n c_{1,ij} \cdot x_i$$

$$z_2 = \sum_{j=1}^m c_{2,j} \cdot z_{1,j}$$

$$= \sum_{j=1}^m c_{2,j} \cdot \left(\sum_{i=1}^n c_{1,ij} \cdot x_i \right)$$

$$= \sum_{i=1}^n \left(\sum_{j=1}^m c_{2,j} \cdot c_{1,ij} \right) \cdot x_i$$

$$= \sum_{i=1}^n c_{\text{ges},i} \cdot x_i$$

$$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1$$

$$\mathbf{z}_2 = \mathbf{W}_2 \mathbf{z}_1 + \mathbf{b}_2$$

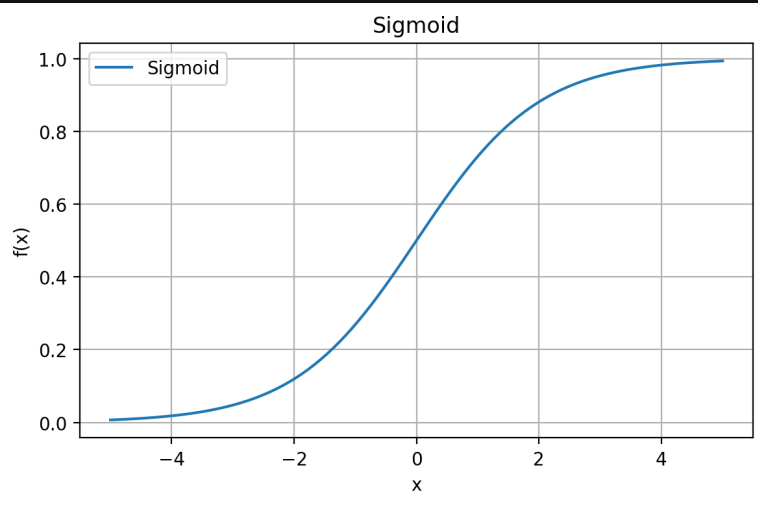
$$= \mathbf{W}_2 (\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$$

$$= \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} + \mathbf{W}_2 \mathbf{b}_1 + \mathbf{b}_2$$

$$= \mathbf{W}_{\text{ges}} \mathbf{x} + \mathbf{b}_{\text{ges}}$$

Aktivierungsfunktion

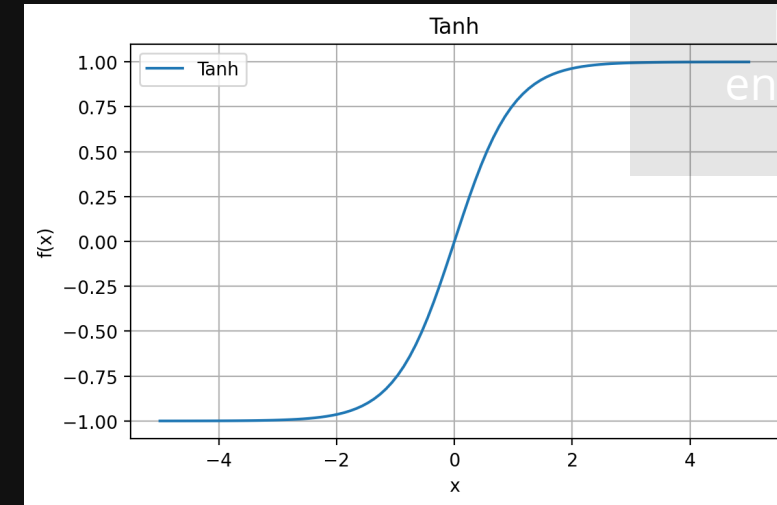
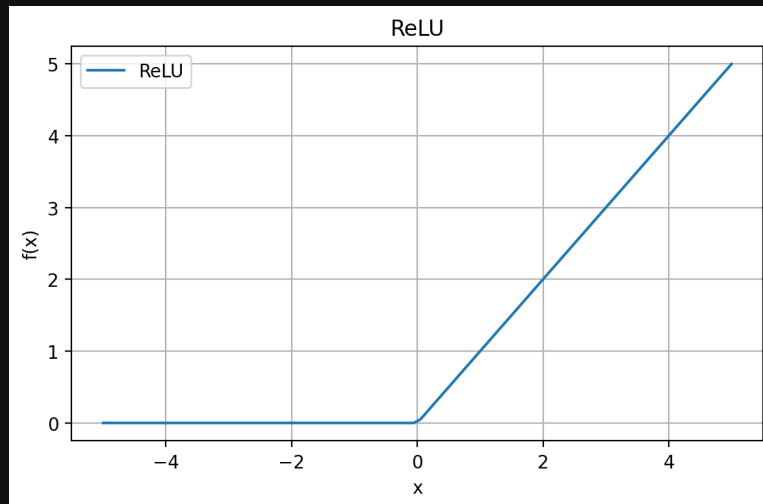
<https://www.youtube.com/embed/IRN7F4E25hk?enablejsapi=1>



Sigmoid

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

$$\text{ReLU}(x) = \max(0, x)$$



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

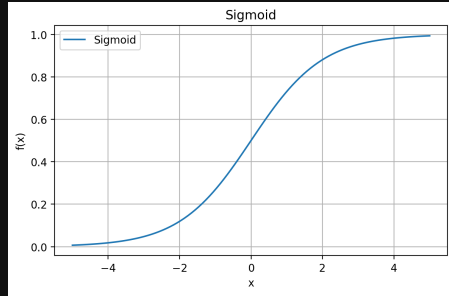
Wahl der Aktivierung:

- Output: Wertebereich Target
- Hidden: Effizienz (?)

Mehr zu Aktivierungen

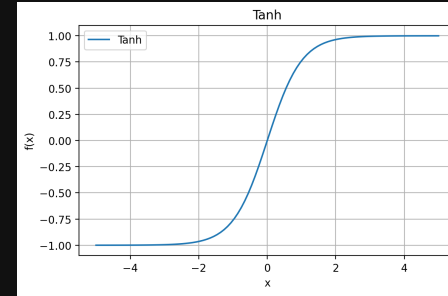
Aktivierungsfunktion

Sigmoid



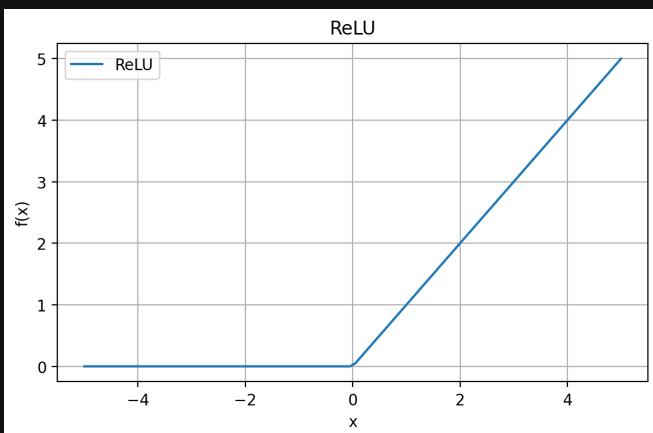
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$



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

$$\tanh'(x) = 1 - \tanh^2(x)$$



$$\text{ReLU}(x) = \max(0, x)$$

$$\text{ReLU}'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Wahl der Aktivierung:

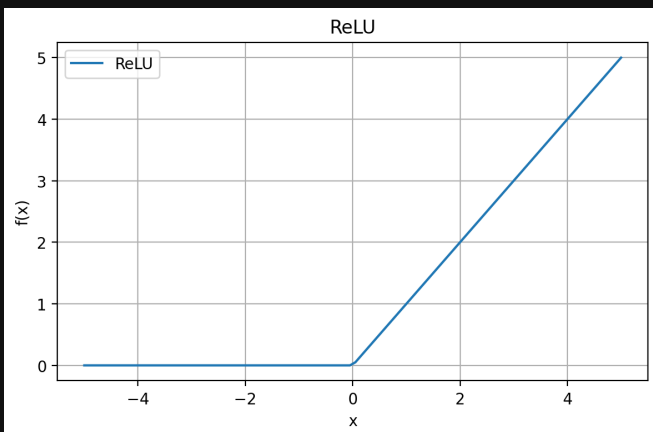
- Output: Wertebereich Target
- Hidden: Effizienz (ReLU)

Mehr zu Aktivierungen

Aktivierungsfunktion

Dying ReLU:

- Permanent $x < 0$
- ReLU & Ableitung = 0
- Perceptron lernt nicht
- Wird Unbrauchbar



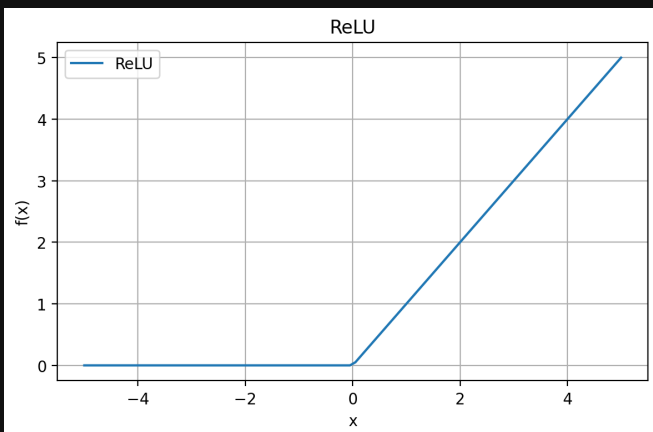
$$\text{ReLU}(x) = \max(0, x)$$

$$\text{ReLU}'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Aktivierungsfunktion

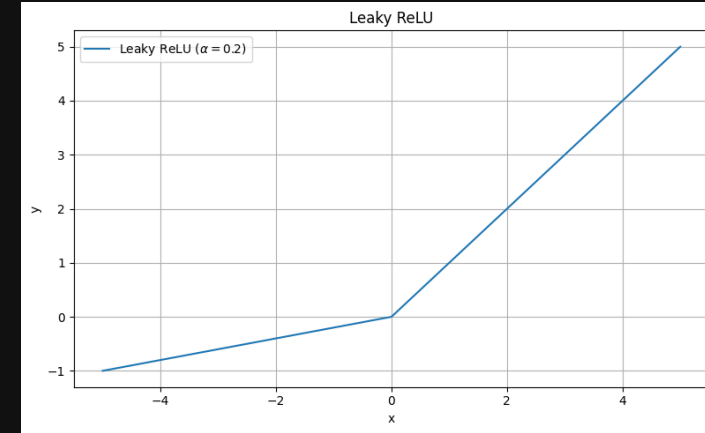
Dying ReLU:

- Permanent $x < 0$
- ReLU & Ableitung = 0
- Perceptron lernt nicht
- Wird Unbrauchbar



$$\text{Leaky ReLU}(x) = \max(\alpha x, x)$$

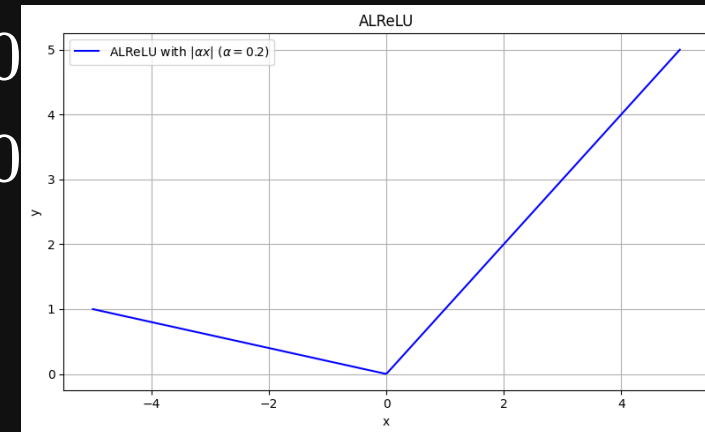
$$\text{Leaky ReLU}'(x) = \begin{cases} \alpha & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$



$$\text{ALReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ |\alpha \cdot x| & \text{if } x < 0 \end{cases}$$

nützlich bei sehr seltener Aktivierung

$$\text{ALReLU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ -\alpha & \text{if } x < 0 \end{cases}$$



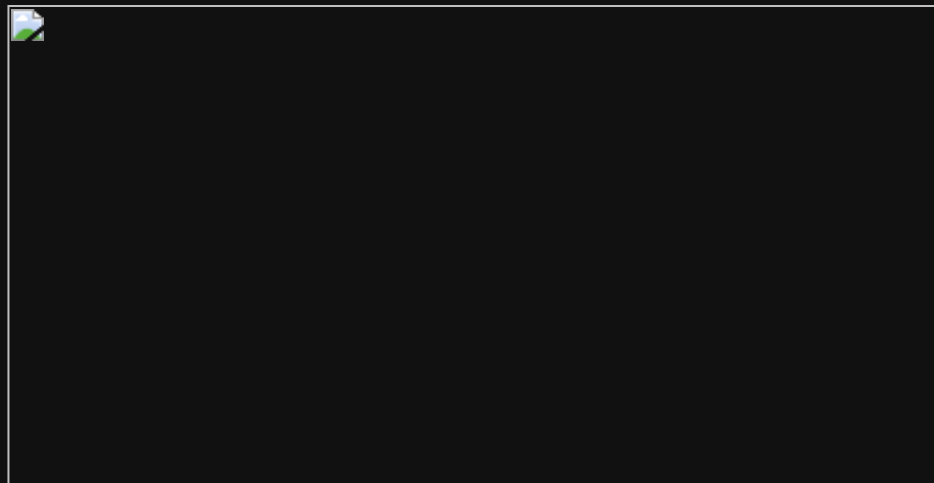
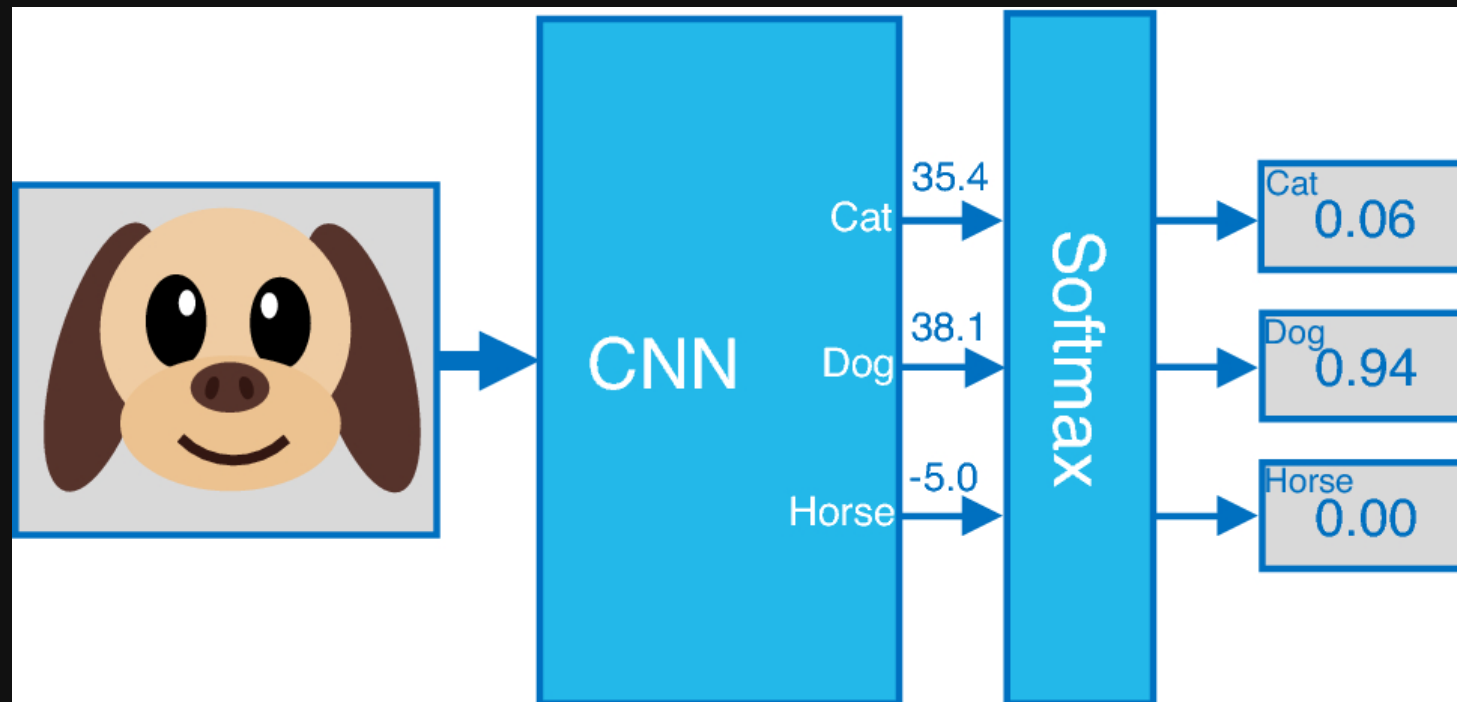
$$\text{ReLU}(x) = \max(0, x)$$

$$\text{ReLU}'(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases}$$



PReLU, RReLU, ELU, SELU, GELU, Swish, Mish, ReLU6, ...

Aktivierungsfunktion

https://www.youtube.com/embed/32bEK-m_30I?enablejsapi=1



Implementation

	Pytorch 	Tensorflow 
Eigenschaften	Pythonic, einfache Syntax schnelleres Training dynamischer Berechnungsgraph höhere Flexibilität	Skalierbar Speichereffizient statisch oder dynamisch
Hauptanwendung	Forschung Prototyping	Grossprojekte Produktion
Community	Forschung	Industrie
Pakete	TorchText, TorchVision, TorchAudio	TF Extended, TF Lite, TF Serving

Beide Frameworks sehr nützlich & weit verbreitet

Mathematik identisch & Aufbau sehr ähnlich

Wahl meist durch Arbeitsumfeld bestimmt

Architektur entwerfen

<https://www.youtube.com/embed/VakbvsvzbD4?enablejsapi=1>

Wenn möglich, bereits existierende Architektur / Modelle verwenden

1. Aufgabe klar definieren (Klassifikation, Regression, Erkennung, ...)
2. Ein- und Ausgabedimension festlegen (MNIST: In: 784; Out: 10)
3. Geeignete Art von Schichten bestimmen (Linear, Convolutional, ...)
4. Anzahl Schichten und Neuronen pro Schicht festlegen
5. Aktivierungsfunktionen festlegen (Hidden & Output)

Implementation: Architektur

- MNIST Classifier

Implementation: Architektur

- MNIST Classifier

```
1 from torch import nn
2 import torch.nn.functional as F
3
4 class Classifier(nn.Module):
5     def __init__(self):
6         super().__init__()
7
8
9
10
11
12     def forward(self, x):
13
14
15
16
17         return x
18
19 model = Classifier()
20 output = model(data)
```

Pytorch
(Meta)

Implementation: Architektur

- MNIST Classifier
- 10 Outputs

```
1 from torch import nn
2 import torch.nn.functional as F
3
4 class Classifier(nn.Module):
5     def __init__(self):
6         super().__init__()
7
8
9
10         self.fc4 = nn.Linear(784, 10)
11
12     def forward(self, x):
13
14
15
16         x = self.fc4(x)
17         return x
18
19 model = Classifier()
20 output = model(data)
```

Pytorch
(Meta)

Implementation: Architektur

- MNIST Classifier
- 10 Outputs
- 3 Hidden Layer

```
1 from torch import nn
2 import torch.nn.functional as F
3
4 class Classifier(nn.Module):
5     def __init__(self):
6         super().__init__()
7         self.fc1 = nn.Linear(784, 64)
8         self.fc2 = nn.Linear(64, 64)
9         self.fc3 = nn.Linear(64, 64)
10        self.fc4 = nn.Linear(64, 10)
11
12    def forward(self, x):
13        x = self.fc1(x)
14        x = self.fc2(x)
15        x = self.fc3(x)
16        x = self.fc4(x)
17        return x
18
19 model = Classifier()
20 output = model(data)
```

Pytorch
(Meta)

Implementation: Architektur

- MNIST Classifier
- 10 Outputs
- 3 Hidden Layer

```
1 from torch import nn
2 import torch.nn.functional as F
3
4 class Classifier(nn.Module):
5     def __init__(self):
6         super().__init__()
7         self.fc1 = nn.Linear(784, 64)
8         self.fc2 = nn.Linear(64, 64)
9         self.fc3 = nn.Linear(64, 64)
10        self.fc4 = nn.Linear(64, 10)
11
12    def forward(self, x):
13        x = F.relu(self.fc1(x))
14        x = F.relu(self.fc2(x))
15        x = F.relu(self.fc3(x))
16        x = self.fc4(x)
17        return x
18
19 model = Classifier()
20 output = model(data)
```

Pytorch
(Meta)

Implementation: Architektur

- MNIST Classifier
- 10 Outputs
- 3 Hidden Layer
- Softmax activation

```
1 from torch import nn
2 import torch.nn.functional as F
3
4 class Classifier(nn.Module):
5     def __init__(self):
6         super().__init__()
7         self.fc1 = nn.Linear(784, 64)
8         self.fc2 = nn.Linear(64, 64)
9         self.fc3 = nn.Linear(64, 64)
10        self.fc4 = nn.Linear(64, 10)
11
12    def forward(self, x):
13        x = F.relu(self.fc1(x))
14        x = F.relu(self.fc2(x))
15        x = F.relu(self.fc3(x))
16        x = F.softmax(self.fc4(x), dim=1)
17        return x
18
19 model = Classifier()
20 output = model(data)
```

Pytorch
(Meta)

Implementation: Architektur

<https://www.youtube.com/embed/rQ6v-xkh4cA?enablejsapi=1>

```
1 import tensorflow as tf
2 from tensorflow.keras import layers
3
4 class Classifier(tf.keras.Model):
5     def __init__(self):
6         super(Classifier, self).__init__()
7         self.fc1 = layers.Dense(64, activation='relu')
8         self.fc2 = layers.Dense(64, activation='relu')
9         self.fc3 = layers.Dense(64, activation='relu')
10        self.fc4 = layers.Dense(10, activation='softmax')
11
12    def call(self, x):
13        x = self.fc1(x)
14        x = self.fc2(x)
15        x = self.fc3(x)
16        x = self.fc4(x)
17        return x
18
19 model = Classifier()
20 model.build((None, 784))
21 model(data)
```

Tensorflow (Google)

```
1 from torch import nn
2 import torch.nn.functional as F
3
4 class Classifier(nn.Module):
5     def __init__(self):
6         super().__init__()
7         self.fc1 = nn.Linear(784, 64)
8         self.fc2 = nn.Linear(64, 64)
9         self.fc3 = nn.Linear(64, 64)
10        self.fc4 = nn.Linear(64, 10)
11
12    def forward(self, x):
13        x = F.relu(self.fc1(x))
14        x = F.relu(self.fc2(x))
15        x = F.relu(self.fc3(x))
16        x = F.softmax(self.fc4(x), dim=1)
17        return x
18
19 model = Classifier()
20 output = model(data)
```

Pytorch
(Meta)

- MNIST Classifier
- 10 Outputs
- 3 Hidden Layer
- Softmax activation

Implementation: Architektur <https://www.youtube.com/embed/OtPB-4Y0iWQ?enablejsapi=1>

```
1 from tensorflow.keras import models
2
3 model = models.Sequential([
4     layers.Dense(64, activation='relu', input_shape=(784,))
5     layers.Dense(64, activation='relu'),
6     layers.Dense(64, activation='relu'),
7     layers.Dense(10, activation='softmax')
8 ])
9
10 model(data)
```

Tensorflow (Google)

```
1 import torch.nn as nn
2
3 model = nn.Sequential(
4     nn.Linear(784, 64),
5     nn.ReLU(),
6     nn.Linear(64, 64),
7     nn.ReLU(),
8     nn.Linear(64, 64),
9     nn.ReLU(),
10    nn.Linear(64, 10),
11    nn.Softmax(dim=1)
12 )
13
14 output = model(data)
```

Pytorch
(Meta)

- MNIST Classifier
- 10 Outputs
- 3 Hidden Layer
- Softmax activation

Trainingsloop

https://www.youtube.com/embed/nWn_5fPEf_E?enablejsapi=1

```
1 trainloader = DataLoader(trainset, batch_size=256, shuffle=True)
```

1. Daten laden (batch)

```
1 for images, labels in trainloader:
```

2. Modell anwenden (forward)

```
1 prediction = model(images)
```

3. Loss berechnen

```
1 loss = criterion(prediction, labels)
```

4. Updates berechnen (backward)

```
1 optimizer.zero_grad()  
2 loss.backward()
```

5. Update durchführen

```
1 optimizer.step()
```

Trainingsloop entwerfen

<https://www.youtube.com/embed/StQ9tHHgxaM?enablejsapi=1>

1. Aufgabe klar definieren
2. Lossfunktion bestimmen
3. Berechnungsschritte definieren

Loss Funktion

<https://www.youtube.com/embed/FVhKvumlyZY?enablejsapi=1>

- Definiert das Ziel des Trainings
- Ziel: Loss minimieren
- erlaubt Vergleich von Modellen
- verschiedene Losses für verschiedene Aufgaben

Loss Funktion

- Mean Squared Error (MSE):
mittlerer quadratische Abweichung

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Loss Funktion

- Mean Squared Error (MSE):
mittlerer quadratische Abweichung

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

- Binäre Cross-Entropy (BCE):
vergleich von Wahrscheinlichkeit einer Klasse

$$BCE = - \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

Loss Funktion

https://www.youtube.com/embed/Y_Sm_xn0Wfss?enablejsapi=1

- Mean Squared Error (MSE):
mittlerer quadratische Abweichung

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

- Binäre Cross-Entropy (BCE):
vergleich von Wahrscheinlichkeit einer Klasse

$$BCE = - \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

- Cross-Entropy (CE):
vergleich von Wahrscheinlichkeiten mehrerer Klassen

$$CE = - \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic})$$

mehr zu
Hintergrund
und
Varianten

Implementation: Training

- Loss: CrossEntropy

Implementation: Training

- Loss: CrossEntropy

```
1 criterion = nn.CrossEntropyLoss()  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13
```

Pytorch

Implementation: Training

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()  
2 optimizer = optim.Adam(model.parameters(), lr=0.003)  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13
```

Pytorch

Implementation: Training

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()  
2 optimizer = optim.Adam(model.parameters(), lr=0.003)  
3  
4 for e in range(epochs):  
5  
6  
7  
8  
9  
10  
11  
12  
13
```

Pytorch

Epoche: alle Daten trainieren

Implementation: Training

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()
2 optimizer = optim.Adam(model.parameters(), lr=0.003)
3
4 for e in range(epochs):
5
6     for images, labels in trainloader:
7
8
9
10
11
12
13
```

Pytorch

Batchwise Input & Target

Implementation: Training

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()
2 optimizer = optim.Adam(model.parameters(), lr=0.003)
3
4 for e in range(epochs):
5
6     for images, labels in trainloader:
7         prediction = model(images)
8         loss = criterion(prediction, labels)
9
10
11
12
13
```

Pytorch

Forward-Pass

Implementation: Training

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()
2 optimizer = optim.Adam(model.parameters(), lr=0.003)
3
4 for e in range(epochs):
5
6     for images, labels in trainloader:
7         prediction = model(images)
8         loss = criterion(prediction, labels)
9
10        optimizer.zero_grad()
11        loss.backward()
12
```

Pytorch

Backward-Pass

Implementation: Training

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()
2 optimizer = optim.Adam(model.parameters(), lr=0.003)
3
4 for e in range(epochs):
5
6     for images, labels in trainloader:
7         prediction = model(images)
8         loss = criterion(prediction, labels)
9
10        optimizer.zero_grad()
11        loss.backward()
12        optimizer.step()
13
```

Pytorch

Update

Implementation: Training

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()
2 optimizer = optim.Adam(model.parameters(), lr=0.003)
3
4 for e in range(epochs):
5     running_loss = 0
6     for images, labels in trainloader:
7         prediction = model(images)
8         loss = criterion(prediction, labels)
9
10        optimizer.zero_grad()
11        loss.backward()
12        optimizer.step()
13
14        running_loss += loss.item()
```

Pytorch

Implementation: Training

<https://www.youtube.com/embed/sjga-FLKP9k?enablejsapi=1>

- Loss: CrossEntropy
- Optimizer: Adam

```
1 criterion = nn.CrossEntropyLoss()
2 optimizer = optim.Adam(model.parameters(), lr=0.003)
3
4 for e in range(epochs):
5     running_loss = 0
6     for images, labels in trainloader:
7         prediction = model(images)
8         loss = criterion(prediction, labels)
9
10        optimizer.zero_grad()
11        loss.backward()
12        optimizer.step()
13
14        running_loss += loss.item()
```

Tensorflow

```
1 model.compile(optimizer=optimizers.Adam(learning_rate=0.003),
2               loss='sparse_categorical_crossentropy',
3               metrics=['accuracy'])
4
5 history = model.fit(train_images, train_labels, epochs=1, batch_size=64)
6
7 print(f'Training loss: {history.history["loss"][0]}')
```

Pytorch

Update

Ableitung der Kosten -> Richtung für Verbesserung -> Update

Gradient Descent $\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$ Lernrate α

Update

Ableitung der Kosten -> Richtung für Verbesserung -> Update

Gradient Descent $\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$ Lernrate α



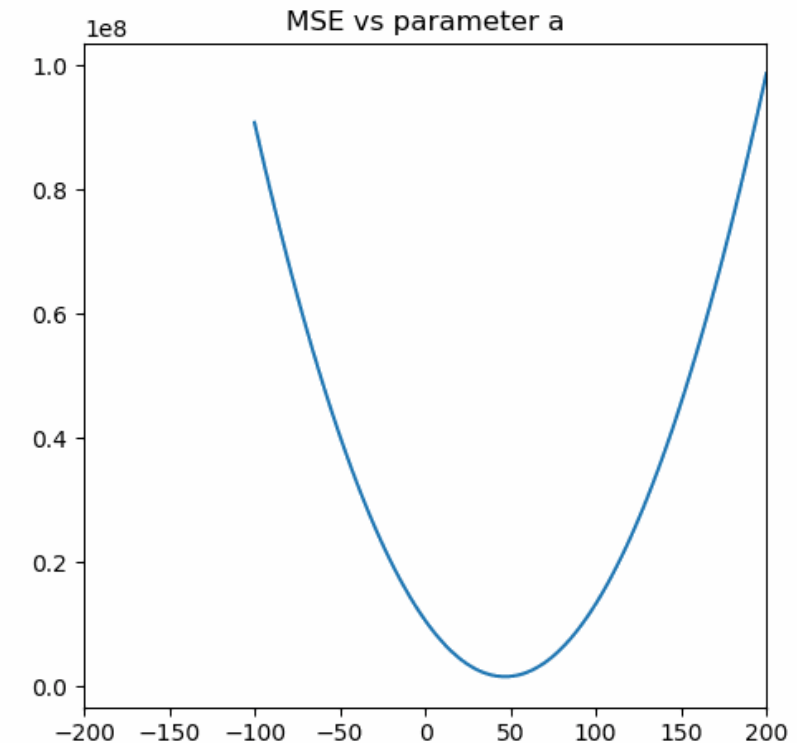
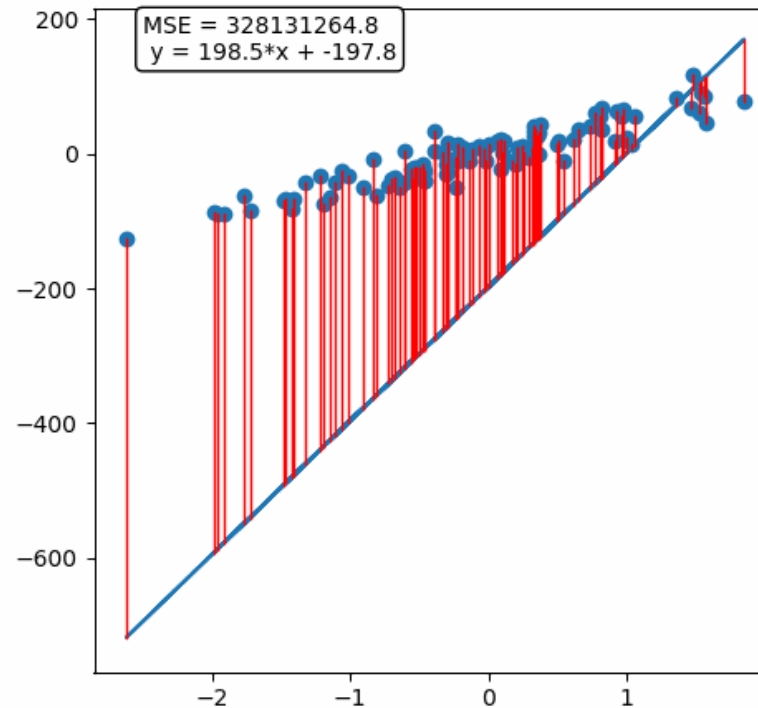
Update

<https://www.youtube.com/embed/0S70j3Lh2Jg?enablejsapi=1>

Ableitung der Kosten -> Richtung für Verbesserung -> Update

Gradient Descent $\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$

Lernrate α



Update

https://www.youtube.com/embed/wyPHES7Sr_U?enablejsapi=1

Ableitung der Kosten -> Richtung für Verbesserung -> Update

Gradient Descent $\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$

Lernrate α

Optimizer

Gradient Descent

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y} = f(x)$$

<https://www.youtube.com/embed/DvTnakOFZos?enablejsapi=1>

x : Input Daten

y : Ziel Daten

\hat{y} : Vorhersage

f : Modell

θ : Parameter

Optimizer

Gradient Descent

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y} = f(x) = a(z) = a(b + w \cdot x)$$

<https://www.youtube.com/embed/QiVluvMSrCs?enablejsapi=1>

b : Bias

w : Gewichte

a : Aktivierung

Optimizer

Gradient Descent

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y} = f(x) = a(z) = a(b + w \cdot x)$$

$$\text{Backward: } \frac{\partial \text{Loss}}{\partial b}$$

Optimizer

Gradient Descent

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y} = f(x) = a(z) = a(b + w \cdot x)$$

$$\text{Backward: } \frac{\partial \text{Loss}}{\partial b} = \frac{\partial \text{Loss}}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b}$$

Optimizer

Gradient Descent

<https://www.youtube.com/embed/hvf555EhBsw?enablejsapi=1>

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y} = f(x) = a(z) = a(b + w \cdot x)$$

$$\text{Backward: } \frac{\partial \text{Loss}}{\partial b} = \frac{\partial \text{Loss}}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b}$$

$$= \sum_i 2(y_i - \hat{y}_i) \cdot 1 \cdot \frac{\partial a}{\partial z} \cdot 1$$

Optimizer

Gradient Descent

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y} = f(x) = a(z) = a(b + w \cdot x)$$

$$\text{Backward: } \frac{\partial \text{Loss}}{\partial b} = \frac{\partial \text{Loss}}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b}$$

$$= \sum_i 2(y_i - \hat{y}_i) \cdot \frac{\partial a}{\partial z}$$

Optimizer

Gradient Descent

<https://www.youtube.com/embed/kPKfXFgClu0?enablejsapi=1>

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y}_i = f(x_i) = a_1(z_1(x_i)) = a_1(b_1 + w_1 \cdot x_i)$$

$$\text{Backward: } \frac{\partial \text{Loss}}{\partial b_1} = \frac{\partial \text{Loss}}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

$$= \sum_i 2(y_i - \hat{y}_i) \frac{\partial a_1}{\partial z_1}$$

Optimizer

Gradient Descent

<https://www.youtube.com/embed/TXWRUcSuars?enablejsapi=1>

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y}_i = f(x_i) = a_2(z_2(a_1(z_1(x_i)))) = a_2(b_2 + w_2 \cdot a_1(b_1 + w_1 \cdot x_i))$$

$$\text{Backward: } \frac{\partial \text{Loss}}{\partial b_1} = \frac{\partial \text{Loss}}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

$$= \sum_i 2(y_i - \hat{y}_i) \cdot \frac{\partial a_2}{\partial z_2} \cdot w_2 \cdot \frac{\partial a_1}{\partial z_1}$$

Optimizer

Gradient Descent

https://www.youtube.com/embed/Fn_U_GPg6il?enablejsapi=1

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\text{Loss} = \sum_i (y_i - \hat{y}_i)^2$$

$$\text{Forward: } \hat{y}_i = f(x_i) = a_3(z_3(a_2(z_2(a_1(z_1(x_i)))))) = a_3(b_3 + w_3 \cdot a_2(b_2 + w_2 \cdot a_1(b_1 + w_1 \cdot x_i)))$$

$$\text{Backward: } \frac{\partial \text{Loss}}{\partial b_1} = \frac{\partial \text{Loss}}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

$$= \sum_i 2(y_i - \hat{y}_i) \cdot \frac{\partial a_3}{\partial z_3} \cdot w_3 \cdot \frac{\partial a_2}{\partial z_2} \cdot w_2 \cdot \frac{\partial a_1}{\partial z_1}$$

Optimizer

Gradient Descent

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

Berechnet Ableitung $\frac{\partial \text{Loss}}{\partial \theta}$
für jedes Update
mit gesamten Datensatz

folgt exakt steilstem Abstieg

Bei Millionen von Daten
sehr rechenintensiv...

Optimizer

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

Gradient Descent

Berechnet Ableitung $\frac{\partial \text{Loss}}{\partial \theta}$
für jedes Update
mit gesamten Datensatz

folgt exakt steilstem Abstieg

Bei Millionen von Daten
sehr rechenintensiv...

Stochastic Gradient Descent (SGD)

Zerlegt Datensatz in Batches
Update nach jeder Batch

folgt Schätzung des steilsten Abstiegs
Fortschritt kann Schwanken

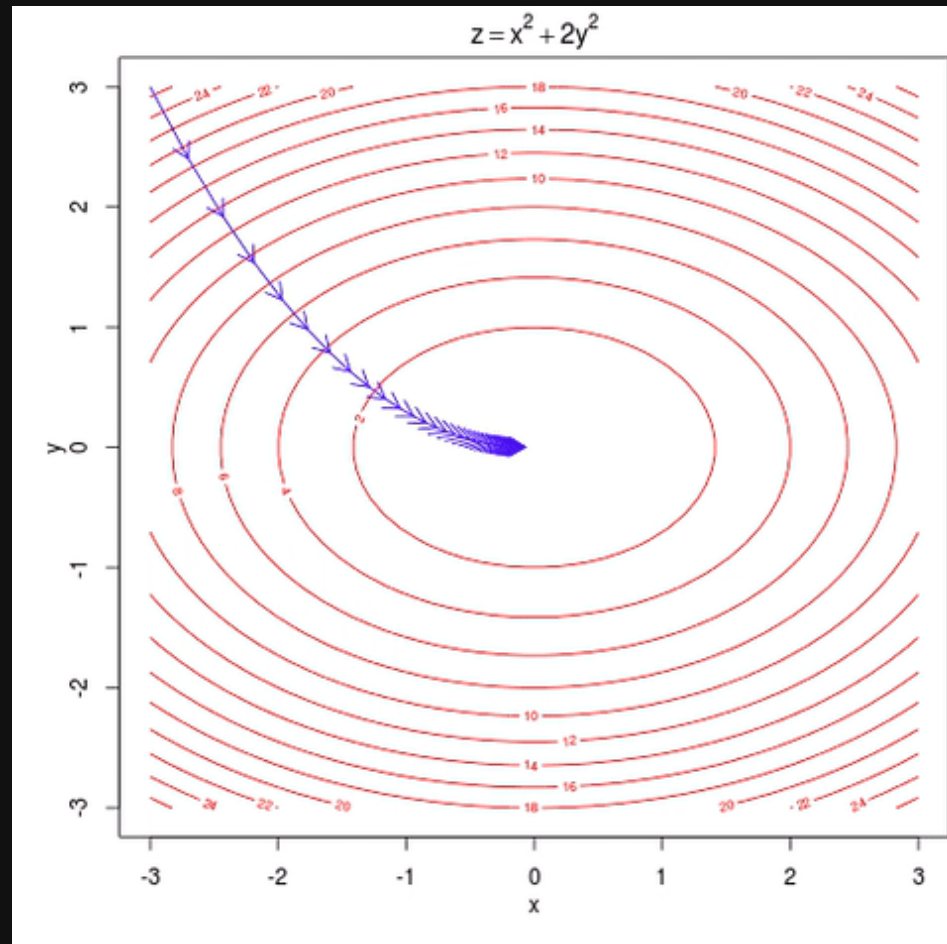
Schnelle Konvergenz
dank schneller Rechnung

Optimizer

Gradient Descent

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

<https://www.youtube.com/embed/WWUapM2JeJ0?enablejsapi=1>



Optimizer

<https://www.youtube.com/embed/213neGdgPjx4?enablejsapi=1>

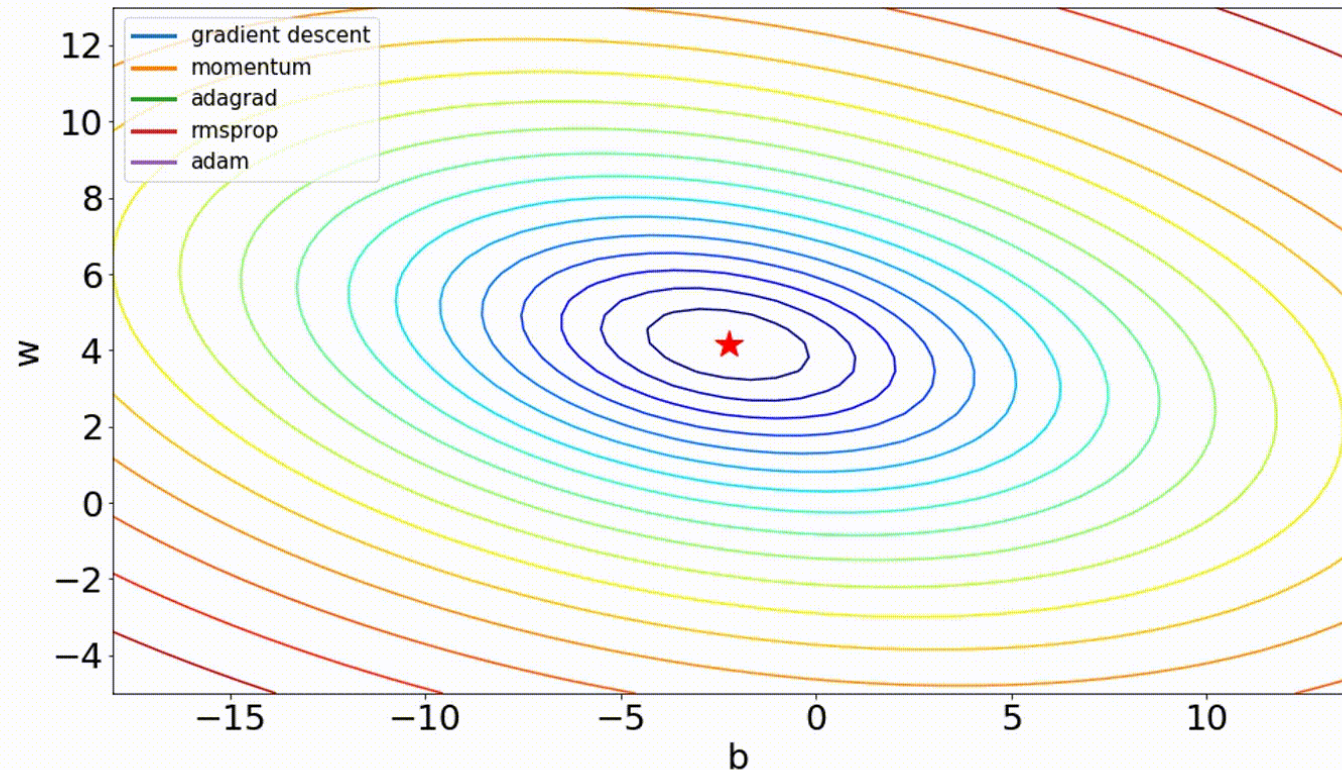
SGD

$$\theta \leftarrow \theta - \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

Momentum

$$v \leftarrow \beta \cdot v + \alpha \cdot \frac{\partial \text{Loss}}{\partial \theta}$$

$$\theta \leftarrow \theta - v$$



mehr zu
Optimizers

Model & Trainingsloop

Hands-On: MNIST Classifier

<https://www.youtube.com/embed/m6pkOvmWXdo?enablejsapi=1>

Bearbeiten Sie [dieses Notebook](#)

- Erstellen Sie einen MNIST Classifier
- Definieren Sie die Trainingsloop
- Testen sie die Trainingsloop ueber zwei Epochen

Die Lösung finden Sie in [diesem Notebook](#)