

# Fundamentals of Neural Networks

---

Mathias Jackermeier

June 10, 2018

Technische Universität München

2018-06-10

## Fundamentals of Neural Networks

Fundamentals of Neural Networks

---

Mathias Jackermeier  
June 10, 2018  
Technische Universität München

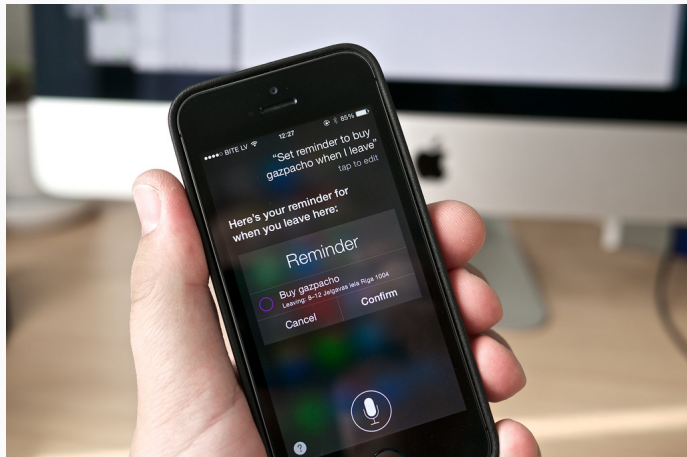


**Figure 1:** A self-driving car.

Credit: Marc van der Chijs / CC BY-ND 2.0



**Figure 1:** A self-driving car.  
Credit: Marc van der Chijs / CC BY-ND 2.0



**Figure 2:** A digital assistant.  
Credit: Kārlis Dambrāns / CC BY 2.0



**Figure 2:** A digital assistant.  
Credit: Kārlis Dambrāns / CC BY 2.0

Jedes intelligente System benutzt neuronale Netze

2018-06-10

Fundamentals of Neural Networks

└ Outline

Outline

---

## Outline

---

Überblick, nicht ins mathematische Detail

2018-06-10

Fundamentals of Neural Networks

└ The Perceptron

The Perceptron

# The Perceptron

---

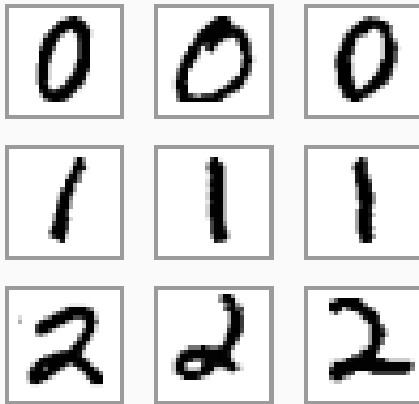
- Predict whether an input image of a handwritten digit shows a zero or another digit

- Predict whether an input image of a handwritten digit shows a zero or another digit



Figure 3: Examples from the MNIST database.  
Credit: Josef Steppan / CC BY-SA 4.0

# MNIST Data Sample



**Figure 3:** Examples from the MNIST database.  
Credit: Josef Steppan / CC BY-SA 4.0

Fundamentals of Neural Networks

└ The Perceptron

└ MNIST Data Sample

2018-06-10

- Predict whether an input image of a handwritten digit shows a zero or another digit
- The image is represented as a flattened vector of pixel intensities  $\mathbf{x} \in \mathbb{R}^{784}$

## Example Task

### Fundamentals of Neural Networks

#### └ The Perceptron

#### └ Example Task

2018-06-10

- Predict whether an input image of a handwritten digit shows a zero or another digit
- The image is represented as a flattened vector of pixel intensities  $\mathbf{x} \in \mathbb{R}^{784}$



- Predict whether an input image of a handwritten digit shows a zero or another digit
- The image is represented as a flattened vector of pixel intensities  $\mathbf{x} \in \mathbb{R}^{784}$
- The output should be 1 if the image shows a zero, otherwise it should be  $-1$

2018-06-10

## Fundamentals of Neural Networks

### └ The Perceptron

### └ Example Task

Example Task

- Predict whether an input image of a handwritten digit shows a zero or another digit
- The image is represented as a flattened vector of pixel intensities  $\mathbf{x} \in \mathbb{R}^{784}$
- The output should be 1 if the image shows a zero, otherwise it should be  $-1$

- Predict whether an input image of a handwritten digit shows a zero or another digit
- The image is represented as a flattened vector of pixel intensities  $\mathbf{x} \in \mathbb{R}^{784}$
- The output should be 1 if the image shows a zero, otherwise it should be  $-1$
- **Idea:** Assign a weight to every input pixel

- Predict whether an input image of a handwritten digit shows a zero or another digit
- The image is represented as a flattened vector of pixel intensities  $\mathbf{x} \in \mathbb{R}^{784}$
- The output should be 1 if the image shows a zero, otherwise it should be  $-1$
- **Idea:** Assign a weight to every input pixel

The perceptron accepts  $n$  input values and computes an output value  $\hat{y}$ :

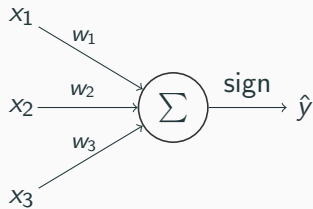
$$\begin{aligned}\hat{y} &= \text{sign} \left( \sum_{i=1}^n w_i x_i \right) \\ &\equiv \hat{y} = \text{sign} \left( \mathbf{w}^\top \mathbf{x} \right)\end{aligned}\quad (1)$$

The perceptron accepts  $n$  input values and computes an output value  $\hat{y}$ :

$$\begin{aligned}\hat{y} &= \text{sign} \left( \sum_{i=1}^n w_i x_i \right) \\ &\equiv \hat{y} = \text{sign} \left( \mathbf{w}^\top \mathbf{x} \right)\end{aligned}\quad (1)$$



Figure 4: A visual representation of the perceptron model.



**Figure 4:** A visual representation of the perceptron model.

1950

Neuron im Gehirn: Nervenzelle Eingaben Ausgaben

- The perceptron is often used in a modified form

2018-06-10

## Fundamentals of Neural Networks

└ The Perceptron

└ Generalizations

- The perceptron is often used in a modified form

f: Aktivierungsfunktion

$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$

# Generalizations

- The perceptron is often used in a modified form
- A scalar bias value can be added to the output computation:

$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$

## Fundamentals of Neural Networks

### └ The Perceptron

#### └ Generalizations

f: Aktivierungsfunktion

- The perceptron is often used in a modified form
- A scalar bias value can be added to the output computation:

$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$

- The sign function can be replaced with a generic function  $f$ :

$$\hat{y} = f(\mathbf{w}^\top \mathbf{x} + b) \quad (3)$$

### └ The Perceptron

#### └ Generalizations

f: Aktivierungsfunktion

- The perceptron is often used in a modified form
- A scalar bias value can be added to the output computation:
$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$
- The sign function can be replaced with a generic function  $f$ :
$$\hat{y} = f(\mathbf{w}^\top \mathbf{x} + b) \quad (3)$$

- The perceptron is often used in a modified form
- A scalar bias value can be added to the output computation:

$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$

- The sign function can be replaced with a generic function  $f$ :

$$\hat{y} = f(\mathbf{w}^\top \mathbf{x} + b) \quad (3)$$

- These modified perceptrons are often called *neurons* or simply *units*

### └ The Perceptron

#### └ Generalizations

f: Aktivierungsfunktion

- The perceptron is often used in a modified form
- A scalar bias value can be added to the output computation:
$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$
- The sign function can be replaced with a generic function  $f$ :
$$\hat{y} = f(\mathbf{w}^\top \mathbf{x} + b) \quad (3)$$
- These modified perceptrons are often called *neurons* or simply *units*



- The perceptron is often used in a modified form
- A scalar bias value can be added to the output computation:

$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$

- The sign function can be replaced with a generic function  $f$ :

$$\hat{y} = f(\mathbf{w}^\top \mathbf{x} + b) \quad (3)$$

- These modified perceptrons are often called *neurons* or simply *units*
- **Notation:** We denote the *weighted input* as

$$z = \mathbf{w}^\top \mathbf{x} + b \quad (4)$$

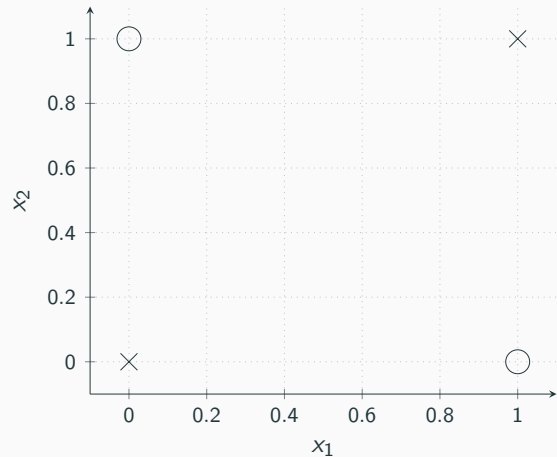
### └ The Perceptron

### └ Generalizations

f: Aktivierungsfunktion

- The perceptron is often used in a modified form
- A scalar bias value can be added to the output computation:
$$\hat{y} = \text{sign}(\mathbf{w}^\top \mathbf{x} + b) \quad (2)$$
- The sign function can be replaced with a generic function  $f$ :
$$\hat{y} = f(\mathbf{w}^\top \mathbf{x} + b) \quad (3)$$
- These modified perceptrons are often called *neurons* or simply *units*
- **Notation:** We denote the *weighted input* as
$$z = \mathbf{w}^\top \mathbf{x} + b \quad (4)$$

# Shortcomings of the Perceptron



**Figure 5:** The perceptron cannot learn the XOR function since the data is not linearly separable.

2018-06-10

## Fundamentals of Neural Networks

### └ The Perceptron

### └ Shortcomings of the Perceptron

Shortcomings of the Perceptron

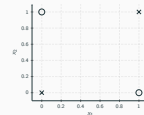


Figure 5: The perceptron cannot learn the XOR function since the data is not linearly separable.

2018-06-10

Fundamentals of Neural Networks  
└─ Feedforward Neural Networks

Feedforward Neural Networks

# Feedforward Neural Networks

---

# Networks of neurons

2018-06-10

Fundamentals of Neural Networks

└ Feedforward Neural Networks

└ Networks of neurons

- **Idea:** A combination of multiple neurons could make much better predictions

- **Idea:** A combination of multiple neurons could make much better predictions

Angelehnt an Gehirn

Auch Multilayer Perzeptron

# Networks of neurons

- **Idea:** A combination of multiple neurons could make much better predictions
- A feedforward neural network is a layered architecture of neurons

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

#### └ Networks of neurons

- **Idea:** A combination of multiple neurons could make much better predictions
- A feedforward neural network is a layered architecture of neurons

Angelehnt an Gehirn

Auch Multilayer Perzeptron

# Networks of neurons

- **Idea:** A combination of multiple neurons could make much better predictions
- A feedforward neural network is a layered architecture of neurons
- The input of a layer is the output of the previous layer

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

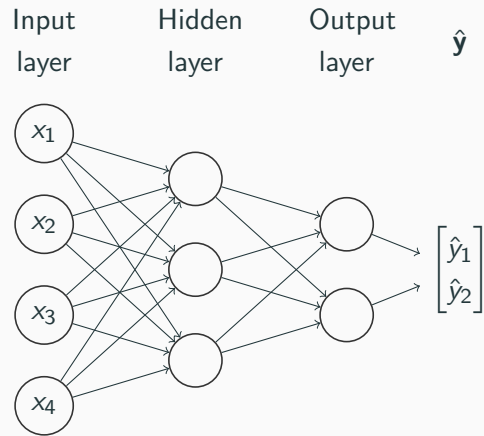
#### └ Networks of neurons

- **Idea:** A combination of multiple neurons could make much better predictions
- A feedforward neural network is a layered architecture of neurons
- The input of a layer is the output of the previous layer

Angelehnt an Gehirn

Auch Multilayer Perzeptron

# Visual Representation



**Figure 6:** A three-layer feedforward neural network.

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Visual Representation

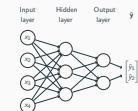


Figure 6: A three-layer feedforward neural network.

Tiefe

Verbindungen nur zwischen Layers

- The design of the output layer depends on the task that we wish to perform

2018-06-10

Fundamentals of Neural Networks

└─ Feedforward Neural Networks

└─ Output Layer

Output Layer

- The design of the output layer depends on the task that we wish to perform



- The design of the output layer depends on the task that we wish to perform
- *Regression*: one single linear neuron

2018-06-10

Fundamentals of Neural Networks

└ Feedforward Neural Networks

└ Output Layer

Output Layer

- The design of the output layer depends on the task that we wish to perform
- *Regression*: one single linear neuron

- The design of the output layer depends on the task that we wish to perform
- *Regression*: one single linear neuron
- *Binary classification*: one single sigmoid neuron

2018-06-10

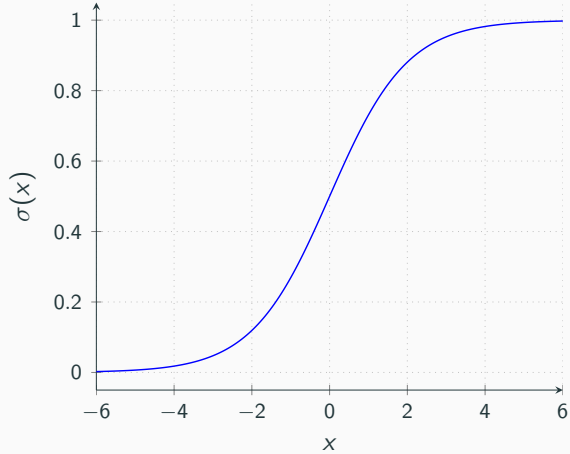
## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

#### └ Output Layer

- The design of the output layer depends on the task that we wish to perform
- *Regression*: one single linear neuron
- *Binary classification*: one single sigmoid neuron

# The Logistic Sigmoid Function



**Figure 7:** The logistic sigmoid function  $\sigma(x) = \frac{1}{1+\exp(-x)}$

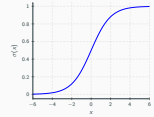
2018-06-10

Fundamentals of Neural Networks

└ Feedforward Neural Networks

└ The Logistic Sigmoid Function

The Logistic Sigmoid Function



**Figure 7:** The logistic sigmoid function  $\sigma(x) = \frac{1}{1+\exp(-x)}$

- The design of the output layer depends on the task that we wish to perform
- *Regression*: one single linear neuron
- *Binary classification*: one single sigmoid neuron
- *Multiclass classification*:  $k$  output units with the softmax function

$$\text{softmax}(x) = \frac{\exp(x)}{\sum_{i=1}^k \exp(z_i)} \quad (5)$$

- The design of the output layer depends on the task that we wish to perform
- *Regression*: one single linear neuron
- *Binary classification*: one single sigmoid neuron
- *Multiclass classification*:  $k$  output units with the softmax function

$$\text{softmax}(x) = \frac{\exp(x)}{\sum_{i=1}^k \exp(z_i)} \quad (5)$$

softmax: Generalisierung, normalisiert

Viele andere Probleme: Auffälliges Verhalten, Bildgenerierung

# Hidden Layers

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Hidden Layers

- The task does not give us any information about how to design the hidden layers

Daher kommt Deep Learning

Verschiedene Repräsentationen (Kanten, Objekte,...)  $\Rightarrow$  Abstraktionen!

Kontinuierlich, nicht linear!

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view

# Hidden Layers

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Hidden Layers

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view

Daher kommt Deep Learning

Verschiedene Repräsentationen (Kanten, Objekte,...)  $\Rightarrow$  Abstraktionen!

Kontinuierlich, nicht linear!

# Hidden Layers

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Hidden Layers

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view
- Deep networks perform almost always better in practice

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view
- Deep networks perform almost always better in practice

Daher kommt Deep Learning

Verschiedene Repräsentationen (Kanten, Objekte,...)  $\Rightarrow$  Abstraktionen!

Kontinuierlich, nicht linear!

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view
- Deep networks perform almost always better in practice
- Activation function: Three common choices are the logistic sigmoid, the tanh, and the rectified linear function

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view
- Deep networks perform almost always better in practice
- Activation function: Three common choices are the logistic sigmoid, the tanh, and the rectified linear function

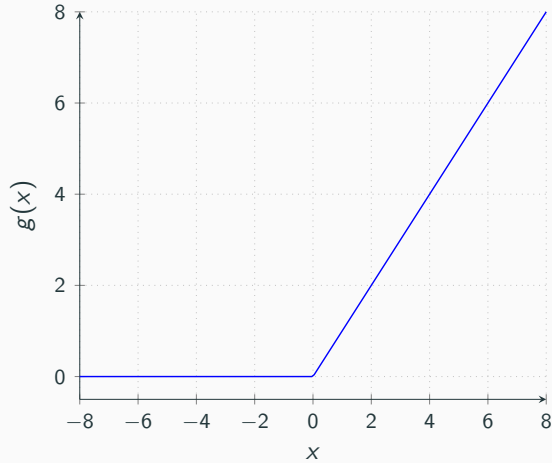
Daher kommt Deep Learning

Verschiedene Repräsentationen (Kanten, Objekte,...)  $\Rightarrow$  Abstraktionen!

Kontinuierlich, nicht linear!



# The Rectified Linear Function



**Figure 8:** The rectified linear function  $g(x) = \max\{0, x\}$

2018-06-10

Fundamentals of Neural Networks

└ Feedforward Neural Networks

└ The Rectified Linear Function

The Rectified Linear Function

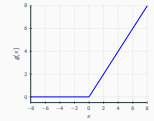


Figure 8: The rectified linear function  $g(x) = \max\{0, x\}$

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view
- Deep networks perform almost always better in practice
- Activation function: Three common choices are the logistic sigmoid, the tanh, and the rectified linear function
- Experimentation and trial & error

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Hidden Layers

- The task does not give us any information about how to design the hidden layers
- Depth: Irrelevant from a theoretical point of view
- Deep networks perform almost always better in practice
- Activation function: Three common choices are the logistic sigmoid, the tanh, and the rectified linear function
- Experimentation and trial & error

- Choose an appropriate input representation

2018-06-10

Fundamentals of Neural Networks

└─ Feedforward Neural Networks

└─ Input Layer

- Choose an appropriate input representation

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$

# Mathematical Formulation

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Mathematical Formulation

$f$  wird komponentenweise angewendet

# Mathematical Formulation

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Mathematical Formulation

f wird komponentenweise angewendet

2018-06-10

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$
- The weight  $w_{ij}^{(l)}$  is the weight from the  $i^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$
- The weight  $w_{ij}^{(l)}$  is the weight from the  $i^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

f wird komponentenweise angewendet

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$
- The weight  $w_{ij}^{(l)}$  is the weight from the  $i^{\text{th}}$  neuron in the  $(l-1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer
- The bias  $b_i^{(l)}$  is the bias of the  $i^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$
- The weight  $w_{ij}^{(l)}$  is the weight from the  $i^{\text{th}}$  neuron in the  $(l-1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer
- The bias  $b_i^{(l)}$  is the bias of the  $i^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

f wird komponentenweise angewendet

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$
- The weight  $w_{ij}^{(l)}$  is the weight from the  $i^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer
- The bias  $b_i^{(l)}$  is the bias of the  $i^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer
- $f^{(l)}$  is the activation function used in the  $l^{\text{th}}$  layer

- We can specify a single neuron with a weight vector  $\mathbf{w}$  and a bias value  $b$
- Since a neural network consists of multiple neurons in a layer, we need weight *matrices*  $\mathbf{W}^{(l)}$  and bias *vectors*  $\mathbf{b}^{(l)}$  to specify the parameters of a layer  $l$
- The weight  $w_{ij}^{(l)}$  is the weight from the  $i^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer
- The bias  $b_i^{(l)}$  is the bias of the  $i^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer
- $f^{(l)}$  is the activation function used in the  $l^{\text{th}}$  layer

f wird komponentenweise angewendet



# Mathematical Formulation

2018-06-10

## Fundamentals of Neural Networks

### └ Feedforward Neural Networks

### └ Mathematical Formulation

- The output at layer  $l$  is then given by

$$\mathbf{a}^{(l)} = f^{(l)} \left( \mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \right) \quad (6)$$

$$\mathbf{a}^0 = \mathbf{x}$$

$$\hat{y} = \mathbf{a}^L$$

- The output at layer  $l$  is then given by

$$\mathbf{a}^{(l)} = f^{(l)} \left( \mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \right) \quad (6)$$

- The output at layer  $l$  is then given by

$$\mathbf{a}^{(l)} = f^{(l)} \left( \mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \right) \quad (6)$$

- The vector of weighted inputs is similarly defined as

$$\mathbf{z}^{(l)} = \mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \quad (7)$$

### └ Feedforward Neural Networks

#### └ Mathematical Formulation

2018-06-10

- The output at layer  $l$  is then given by

$$\mathbf{a}^{(l)} = f^{(l)} \left( \mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \right) \quad (6)$$

- The vector of weighted inputs is similarly defined as

$$\mathbf{z}^{(l)} = \mathbf{W}^{(l)\top} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)} \quad (7)$$

$$\mathbf{a}^0 = \mathbf{x}$$

$$\hat{\mathbf{y}} = \mathbf{a}^L$$

2018-06-10

Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

Training Feedforward Neural  
Networks

---

# Training Feedforward Neural Networks

---

- We have training examples  $\mathbb{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$  with corresponding labels  $\mathbb{Y}$

Zufällig initialisiert

# Training Scenario

2018-06-10

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

#### └ Training Scenario

- We have training examples  $\mathbb{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$  with corresponding labels  $\mathbb{Y}$
- We want to learn a mapping from  $\mathbb{X}$  to  $\mathbb{Y}$

- We have training examples  $\mathbb{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$  with corresponding labels  $\mathbb{Y}$
- We want to learn a mapping from  $\mathbb{X}$  to  $\mathbb{Y}$

Zufällig initialisiert

# Training Scenario

2018-06-10

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

#### └ Training Scenario

- We have training examples  $\mathbb{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$  with corresponding labels  $\mathbb{Y}$
- We want to learn a mapping from  $\mathbb{X}$  to  $\mathbb{Y}$
- **Idea:** Iteratively adjust the parameters of the neural network

- We have training examples  $\mathbb{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$  with corresponding labels  $\mathbb{Y}$
- We want to learn a mapping from  $\mathbb{X}$  to  $\mathbb{Y}$
- **Idea:** Iteratively adjust the parameters of the neural network

Zufällig initialisiert

# Cost Functions

- The cost function  $J(\theta)$  is a measure of how good the network performs

2018-06-10

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

#### └ Cost Functions

Von den Parametern zu einem Skalar

Größer als 0

Auch Loss oder Error

# Cost Functions

- The cost function  $J(\theta)$  is a measure of how good the network performs
- Learning can be framed as minimizing the cost function

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

#### └ Cost Functions

2018-06-10

- The cost function  $J(\theta)$  is a measure of how good the network performs
- Learning can be framed as minimizing the cost function

Von den Parametern zu einem Skalar

Größer als 0

Auch Loss oder Error



- The cost function  $J(\theta)$  is a measure of how good the network performs
- Learning can be framed as minimizing the cost function
- The total cost is a sum over the costs of the individual training examples:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \theta) \quad (8)$$

Von den Parametern zu einem Skalar

Größer als 0

Auch Loss oder Error

- The cost function  $J(\theta)$  is a measure of how good the network performs
- Learning can be framed as minimizing the cost function
- The total cost is a sum over the costs of the individual training examples:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \theta) \quad (8)$$

# Mean squared error

2018-06-10

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

#### └ Mean squared error

- In regression, the per-example loss is commonly

$$\mathcal{L}(\mathbf{x}, y, \theta) = \frac{1}{2}(\hat{y} - y)^2 \quad (9)$$

- In regression, the per-example loss is commonly

$$\mathcal{L}(\mathbf{x}, y, \theta) = \frac{1}{2}(\hat{y} - y)^2 \quad (9)$$

Label: Skalar was wir vorhersagen wollen

Distanz

Erfüllt Bedingungen

- In binary classification, we often use the cross-entropy loss

$$\mathcal{L}(\mathbf{x}, y, \theta) = -y \ln \hat{y} - (1 - y) \ln(1 - \hat{y}) \quad (10)$$

MSE schlecht in Klassifikation

Label: 1 oder 0

- In multiclass classification, the cross-entropy becomes

$$\mathcal{L}(\mathbf{x}, \mathbf{y}, \theta) = -\ln \hat{y}_i \quad (11)$$

Label: one-hot

i-te Klasse

Maximum Likelihood Estimation

# Stochastic Gradient Descent

- (Stochastic) Gradient Descent is the most common algorithm to minimize cost functions in neural networks

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

### └ Stochastic Gradient Descent

⇒ Kleine Änderungen in die entgegengesetzte Richtung des Gradienten  
Learning Rate nicht zu klein und nicht zu groß  
Erweiterungen

2018-06-10

# Stochastic Gradient Descent

- (Stochastic) Gradient Descent is the most common algorithm to minimize cost functions in neural networks
- A change  $\Delta\theta$  in the parameters corresponds roughly to the change

$$\Delta J(\theta) \approx \nabla J(\theta)^\top \Delta\theta \quad (12)$$

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

#### └ Stochastic Gradient Descent

⇒ Kleine Änderungen in die entgegengesetzte Richtung des Gradienten  
Learning Rate nicht zu klein und nicht zu groß  
Erweiterungen

- (Stochastic) Gradient Descent is the most common algorithm to minimize cost functions in neural networks
- A change  $\Delta\theta$  in the parameters corresponds roughly to the change

$$\Delta J(\theta) \approx \nabla J(\theta)^\top \Delta\theta \quad (12)$$

# Stochastic Gradient Descent

- (Stochastic) Gradient Descent is the most common algorithm to minimize cost functions in neural networks
- A change  $\Delta\theta$  in the parameters corresponds roughly to the change

$$\Delta J(\theta) \approx \nabla J(\theta)^\top \Delta\theta \quad (12)$$

- To minimize  $J(\theta)$ , choose

$$\Delta\theta = -\eta \nabla J(\theta), \quad (13)$$

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

#### └ Stochastic Gradient Descent

2018-06-10

- (Stochastic) Gradient Descent is the most common algorithm to minimize cost functions in neural networks
- A change  $\Delta\theta$  in the parameters corresponds roughly to the change

$$\Delta J(\theta) \approx \nabla J(\theta)^\top \Delta\theta \quad (12)$$

- To minimize  $J(\theta)$ , choose

$$\Delta\theta = -\eta \nabla J(\theta), \quad (13)$$

⇒ Kleine Änderungen in die entgegengesetzte Richtung des Gradienten  
Learning Rate nicht zu klein und nicht zu groß  
Erweiterungen

# Stochastic Gradient Descent

- (Stochastic) Gradient Descent is the most common algorithm to minimize cost functions in neural networks
- A change  $\Delta\theta$  in the parameters corresponds roughly to the change

$$\Delta J(\theta) \approx \nabla J(\theta)^\top \Delta\theta \quad (12)$$

- To minimize  $J(\theta)$ , choose

$$\Delta\theta = -\eta \nabla J(\theta), \quad (13)$$

- *Stochastic* Gradient Descent computes only an approximation of the gradient

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

### └ Stochastic Gradient Descent

2018-06-10

- (Stochastic) Gradient Descent is the most common algorithm to minimize cost functions in neural networks

- A change  $\Delta\theta$  in the parameters corresponds roughly to the change

$$\Delta J(\theta) \approx \nabla J(\theta)^\top \Delta\theta \quad (12)$$

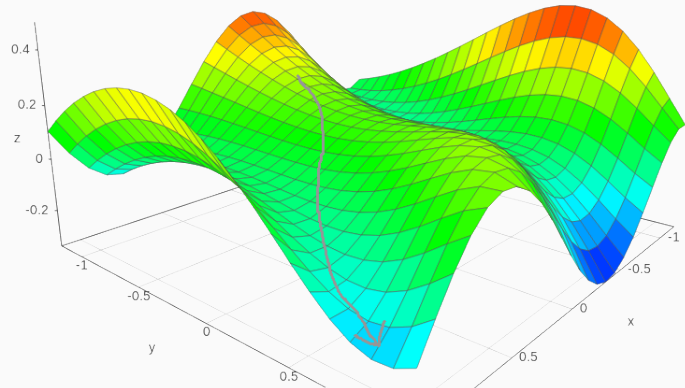
- To minimize  $J(\theta)$ , choose

$$\Delta\theta = -\eta \nabla J(\theta), \quad (13)$$

- *Stochastic* Gradient Descent computes only an approximation of the gradient

⇒ Kleine Änderungen in die entgegengesetzte Richtung des Gradienten  
Learning Rate nicht zu klein und nicht zu groß  
Erweiterungen





**Figure 9:** Stochastic Gradient Descent.

Created with <https://academo.org/demos/3d-surface-plotter/>

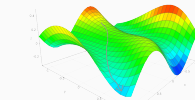
2018-06-10

Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

└ Stochastic Gradient Descent

Stochastic Gradient Descent



**Figure 9:** Stochastic Gradient Descent.  
Created with <https://academo.org/demos/3d-surface-plotter/>

# Back-propagation

2018-06-10

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

### └ Back-propagation

- The back-propagation algorithm efficiently computes the gradient of the cost function

- The back-propagation algorithm efficiently computes the gradient of the cost function

# Back-propagation

2018-06-10

## Fundamentals of Neural Networks

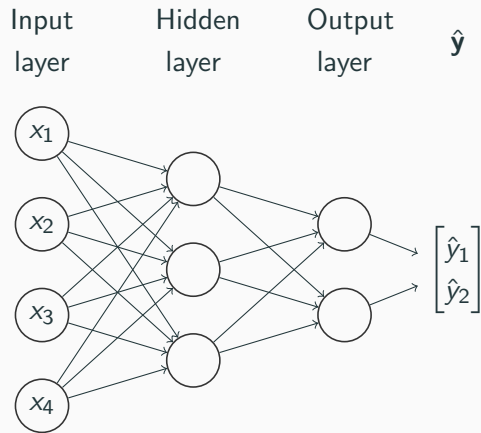
### └ Training Feedforward Neural Networks

#### └ Back-propagation

- The back-propagation algorithm efficiently computes the gradient of the cost function
- It can be derived by recursively applying the chain rule to the layers of the neural network, beginning with the output layer

- The back-propagation algorithm efficiently computes the gradient of the cost function
- It can be derived by recursively applying the chain rule to the layers of the neural network, beginning with the output layer

# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

Fundamentals of Neural Networks  
└ Training Feedforward Neural Networks  
└ Back-propagation

Back-propagation

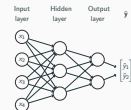
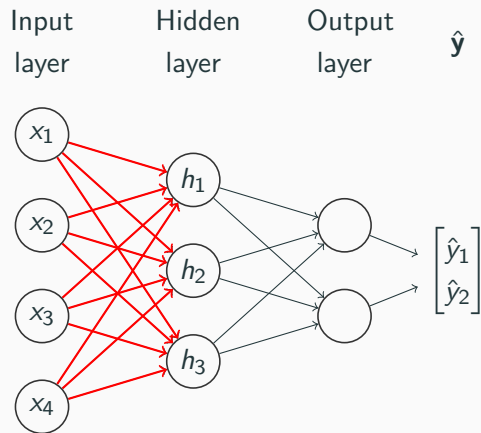


Figure 10: The Back-propagation algorithm.

# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

## Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

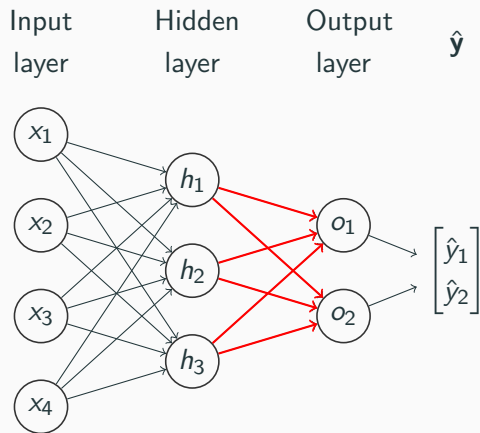
└ Back-propagation

Back-propagation



Figure 10: The Back-propagation algorithm.

# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

## Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

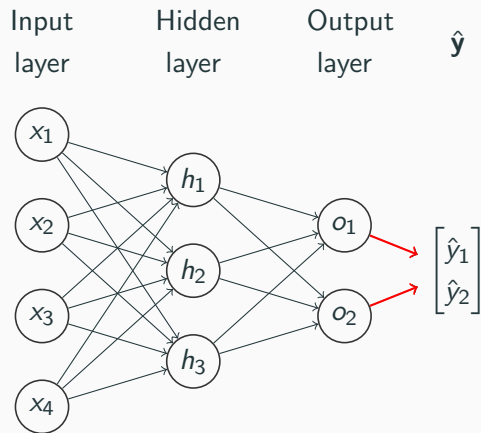
└ Back-propagation

Back propagation



Figure 10: The Back-propagation algorithm.

# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

## Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

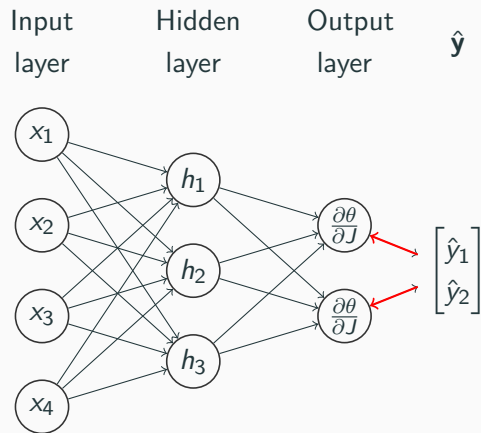
└ Back-propagation

Back-propagation



Figure 10: The Back-propagation algorithm.

# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

## Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

└ Back-propagation

Back-propagation

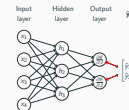
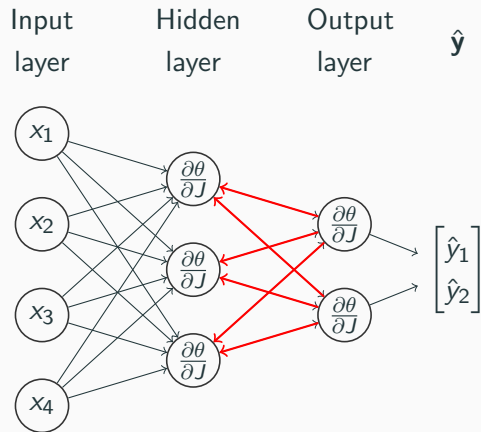


Figure 10: The Back-propagation algorithm.



# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

### └ Back-propagation

Back-propagation

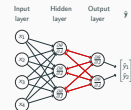
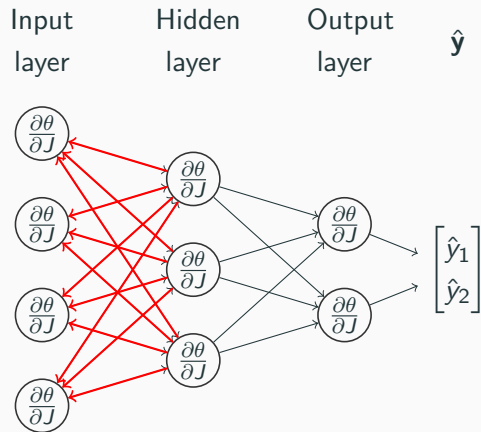


Figure 10: The Back-propagation algorithm.

# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

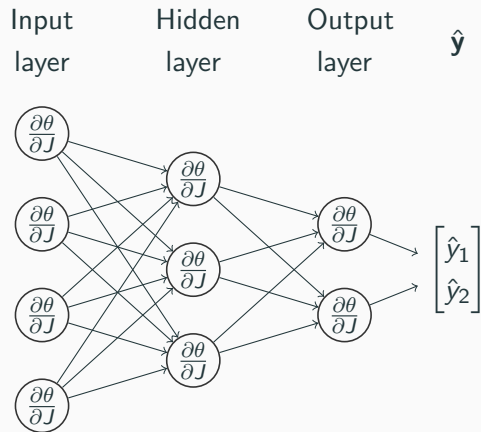
Fundamentals of Neural Networks  
└ Training Feedforward Neural Networks  
└ Back-propagation

Back-propagation



Figure 10: The Back-propagation algorithm.

# Back-propagation



**Figure 10:** The Back-propagation algorithm.

2018-06-10

## Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

└ Back-propagation

Back-propagation

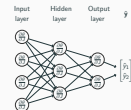


Figure 10: The Back-propagation algorithm.

1. Propagate all training examples of a minibatch forward through the network

2018-06-10

Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

└ The Complete Learning Algorithm

1. Propagate all training examples of a minibatch forward through the network
2. Compute all gradients using back-propagation

2018-06-10

Fundamentals of Neural Networks

└ Training Feedforward Neural Networks

└ The Complete Learning Algorithm

1. Propagate all training examples of a minibatch forward through the network
2. Compute all gradients using back-propagation

1. Propagate all training examples of a minibatch forward through the network
2. Compute all gradients using back-propagation
3. Compute the average gradient

1. Propagate all training examples of a minibatch forward through the network
2. Compute all gradients using back-propagation
3. Compute the average gradient
4. Update the parameters in the negative direction of the gradient

1. Propagate all training examples of a minibatch forward through the network
2. Compute all gradients using back-propagation
3. Compute the average gradient
4. Update the parameters in the negative direction of the gradient

# The Complete Learning Algorithm

1. Propagate all training examples of a minibatch forward through the network
2. Compute all gradients using back-propagation
3. Compute the average gradient
4. Update the parameters in the negative direction of the gradient
5. Repeat until the cost is low enough

2018-06-10

## Fundamentals of Neural Networks

### └ Training Feedforward Neural Networks

### └ The Complete Learning Algorithm

1. Propagate all training examples of a minibatch forward through the network
2. Compute all gradients using back-propagation
3. Compute the average gradient
4. Update the parameters in the negative direction of the gradient
5. Repeat until the cost is low enough



## Conclusion

---

Komplexe Netzwerke einfacher Einheiten

Abstraktionen

Lernen: Kleine Updates der Parameter so dass das Netzwerk besser wird

Überall in Deep Learning

Viele weitere Anwendungen in der Zukunft

**Thank you!**