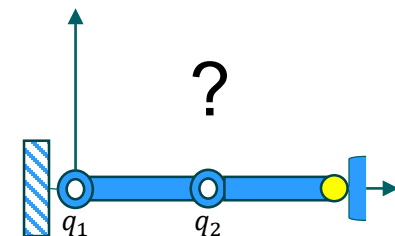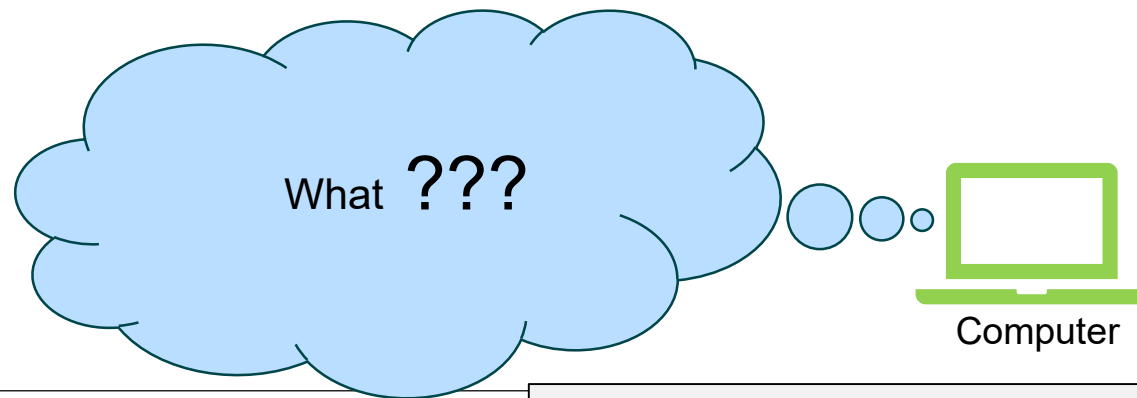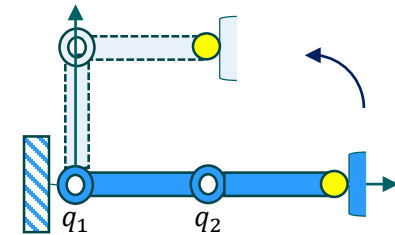# Robotics und Machine Learning (ML)

**Prof. Dr.-Ing. Eric Kaigom**

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

**Learning from experience – What are the coordinates of ⭘ for $q_1 = \frac{\pi}{2}$ and $q_2 = -\frac{\pi}{2}$?**

I know **from experience** that the coordinates of ⭘ are (x = 1, y = 1)!

Human

$q_1$    $q_2$

What **???**

Computer

?

$q_1$    $q_2$

Can a computer (or in general a machine) do the same, **i.e, learn from experience** like a human?

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

**Learning from experience – What are the coordinates of ⦿ for $q_1 = \frac{\pi}{2}$ and $q_2 = -\frac{\pi}{2}$ ?**



I know **from experience** that the coordinates of ⦿ are (x = 1, y = 1)!

Human

I predict **after being trained** that the coordinates of ⦿ are (x = 0.99, y = 1.01)!

Computational methods

Computer

Experience comes from training data

Can a computer (or in general a machine) do the same, **i.e, learn from experience** like a human?? **Yes!**

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

**Learning from experience – What are the coordinates of ○ for $q_1 = \frac{\pi}{4}$ and $q_2 = -\frac{\pi}{4}$ ?**

Well, …. …. I can imagine, but…
…..**actually, I have no idea**!

Human

I **predict** **after training** that the coordinates of ○ are (x = 1.7, y = 0.7)!

Computational methods

Computer

Experience from training data

A trained **machine** can **outperform** human capabilities in e.g. non-nominal cases by **learning** **from data** (**entailing a rich and diversified experience**).

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

# Machine Learning

## Machine Learning – A brief overview

> Machine learning is the **process** of using **algorithms** to learn a **model** that helps **predict** the outcomes of previously **unseen events** (i.e., model input data) by using experiences condensed in **training data.**

What

Challenge and constraints understanding, selection of a learning approach

How

| Data management | Model development | Applications |
|---|---|---|
| • Collection | • Hyperparameterization | • Forward kinematics |
| • Pre-processing | • Optimization | • Backward kinematics |
| • etc | • Validation, etc | • etc |

Our focus

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Robotics and Machine Learning – Model Capture

## Understanding challenges, constraints, and opportunities - example



### Why and what

- Forward kinematics (FK) and backward kinematics (BK) not available in a closed form (e.g., refurbished robots)

- FK and BK prone to uncertainties (e.g., geometry, structural properties, noise, etc.)

- Access to input (e.g., joint pos.) and output (e.g., end-effector pos.) data via measurements and data acquisition

### How

- Instead of a model-based (i.e., analytical) derivation of the input-output or output-input mapping in closed form, data samples are leveraged to capture the mapping model (FK and BK)

- Captured models (FK and BK) can be used in arbitrary (e.g., physical or virtualized) applications and even shared via mail („FK and BK to go!")

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Robotics and Machine Learning – Model Capture

## Selecting a suitable machine learning approach

Our focus is **general**, i.e., encompassing the **non-linear case**!

- • Regression ($x = f(u)$+noise) **Our focus**
- • **Classification** ($x \in C_1 / \cdots / C_n$)
- • **Etc**.

### Supervised learning

- • Input example data known
- • Corresponding/labelled output example data known
- • **Learning goal**: retrieve noised mapping/association function from input and output examples

### Unsupervised learning

- • Input example data known
- • There is no corresponding/labelled output example data available
- • **Learning goal**: retrieve some properties in/from input data

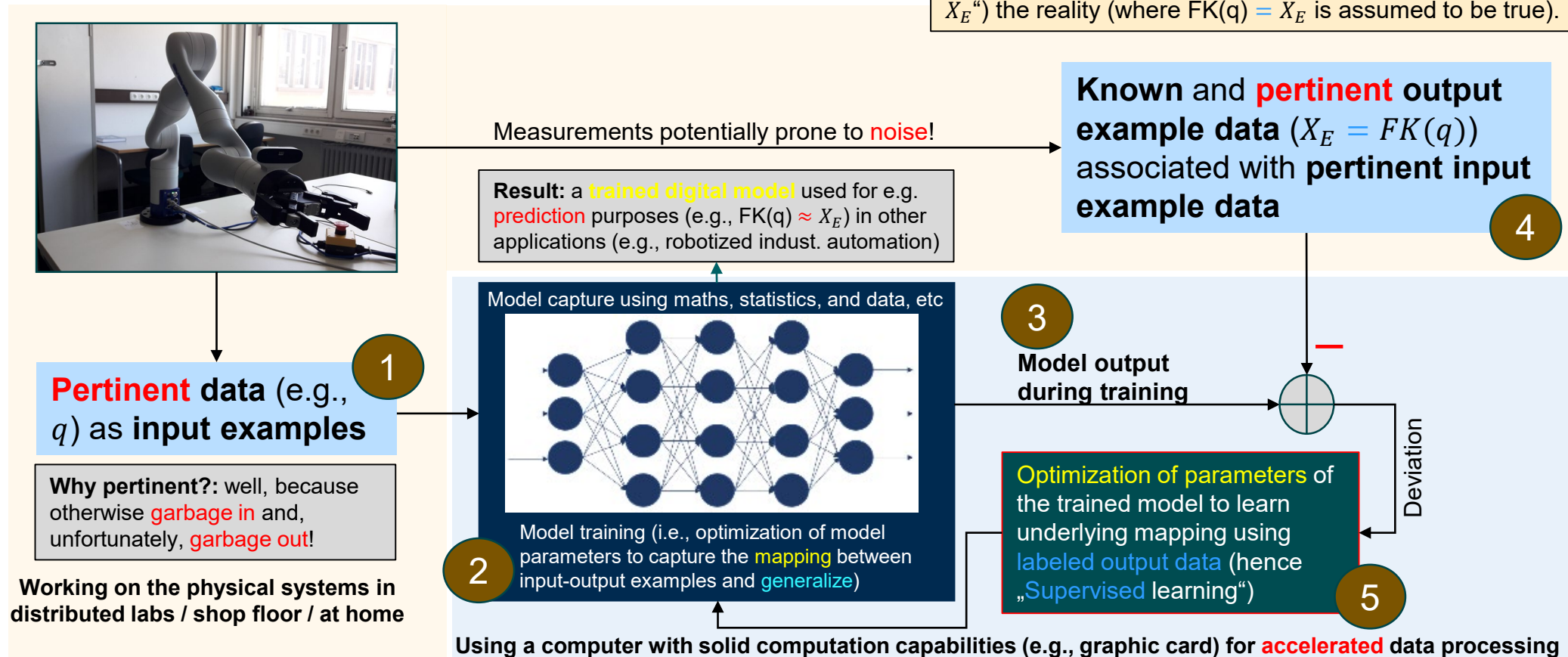- • **Clustering** (captures characteristics from input data to infer properties)
- • **Dimension reduction** (learns (statistically) a representation of input data in a lower dimension space (through e.g. projection) for e.g. efficiency reasons)
- • **Etc**.

### Semi-supervised learning

- • Input example data known
- • Partially labelled output example data known
- • **Learning goal**: retrieve labelling from output examples.

- • **Transductive Super Vector Machines** (harnesses a low density region between class distributions to find decision boundary for separation and generalization)
- • **Generative adversarial networks** (Discriminator generates class labels)
- • **Etc**.

### Reinforcement learning

- • AI system-environment-interaction for the AI system to take action
- • If feedback (environm. to learning system), actions are +/- rewarded
- • **Learning goal**: retrieve a policy to meet a goal (by increasing reward)

- • **Q-Learning** (model-free, set of discrete actions selected via max. Q-value)
- • **Policy gradient** (has a large continuous space of actions and probability-based policies)
- • **Actor-critic** (combines advantages of Q-Learning and Policy gradient)
- • **Etc**.

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Machine Learning – Supervised Learning (Big Picture)

**Hint:** The trained AI model predicts (i.e., approximates, hence „ FK(q) $\approx$ $X_E$") the reality (where FK(q) $=$ $X_E$ is assumed to be true).

**Known** and **pertinent output example data** $(X_E = FK(q))$ associated with **pertinent input example data**

④

Measurements potentially prone to noise!

**Result:** a trained digital model used for e.g. prediction purposes (e.g., FK(q) $\approx$ $X_E$) in other applications (e.g., robotized indust. automation)

Model capture using maths, statistics, and data, etc

③

**Model output during training**



①

**Pertinent data** (e.g., $q$) as **input examples**

Deviation

**Why pertinent?:** well, because otherwise garbage in and, unfortunately, garbage out!

Model training (i.e., optimization of model parameters to capture the mapping between input-output examples and generalize)

②

Optimization of parameters of the trained model to learn underlying mapping using labeled output data (hence „Supervised learning")

⑤

**Working on the physical systems in distributed labs / shop floor / at home**

**Using a computer with solid computation capabilities (e.g., graphic card) for accelerated data processing**

Eric Kaigom
kaigom@fb2.fra-uas.de

**Ethics:** Although compelling from an anthropocentric (human-centered) viewpoint, the related energy and ressource consumption for model training might be (at scale) questionable from an ecocentric perspective (nature with own rights)! How to globally and sustainably address this challenge for the biosphere and a shared well-being?!

Pivotal future research and application avenues YOU might want to streamline and shape!

# Preparing data

# Robotics and Machine Learning – Model Capture

## Why data preparation?
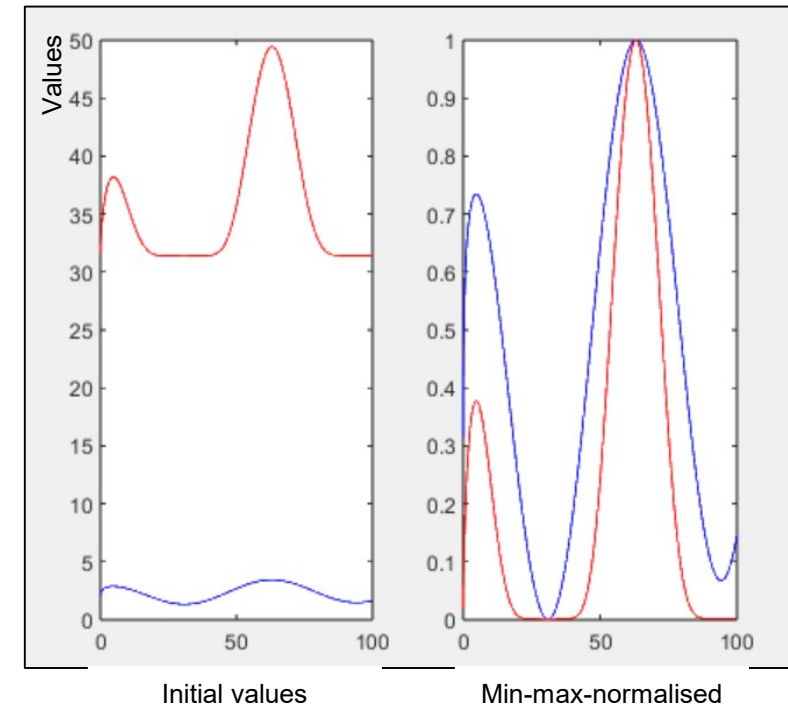
- An important step before data analysis (…in order to learn an unkown mapping or correspondence)

- Pivotal point in Machine Learning: How **features** are **distributed**

- **Implications:**

  - Preserve how the data samples (e.g., joint positions) relate to each other

  - Accommodate the impact of units

  - Avoid effects of relative scale (otherwise: learning performance can drop, i.e., convergence takes longer due to overshoot. More later on.)

- **Challenges:**
  - Impacts of outliers
  - etc

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## **Data preparation - Minimum-Maximum Normalization**

- $x = [x_1 \; x_2 \dots x_n]$ the vector of $n$ scalars

- $x_{min}$ and $x_{max}$ is minimum and maximum value in $x$
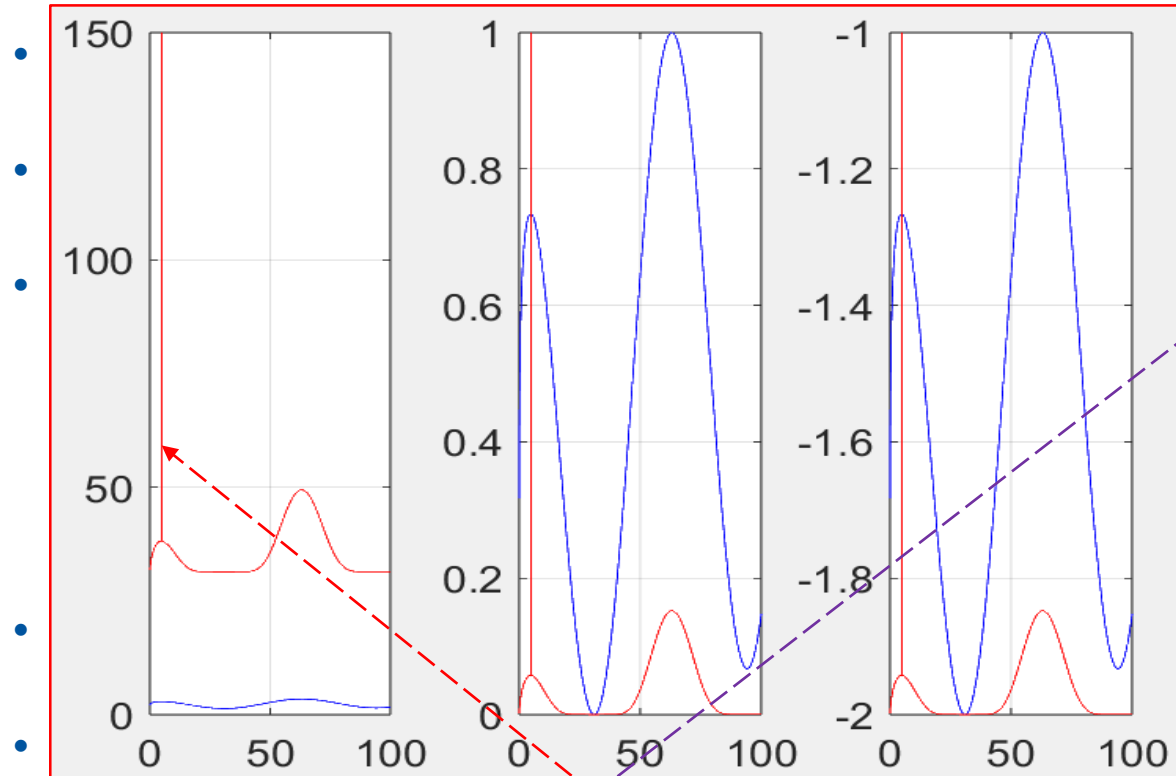
- $\hat{x}$ is the min-max-normalized value of $x$:

$$\hat{x} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

- Values in $\hat{x}$ are squeezed between 0 and 1, i.e., within [0 1]

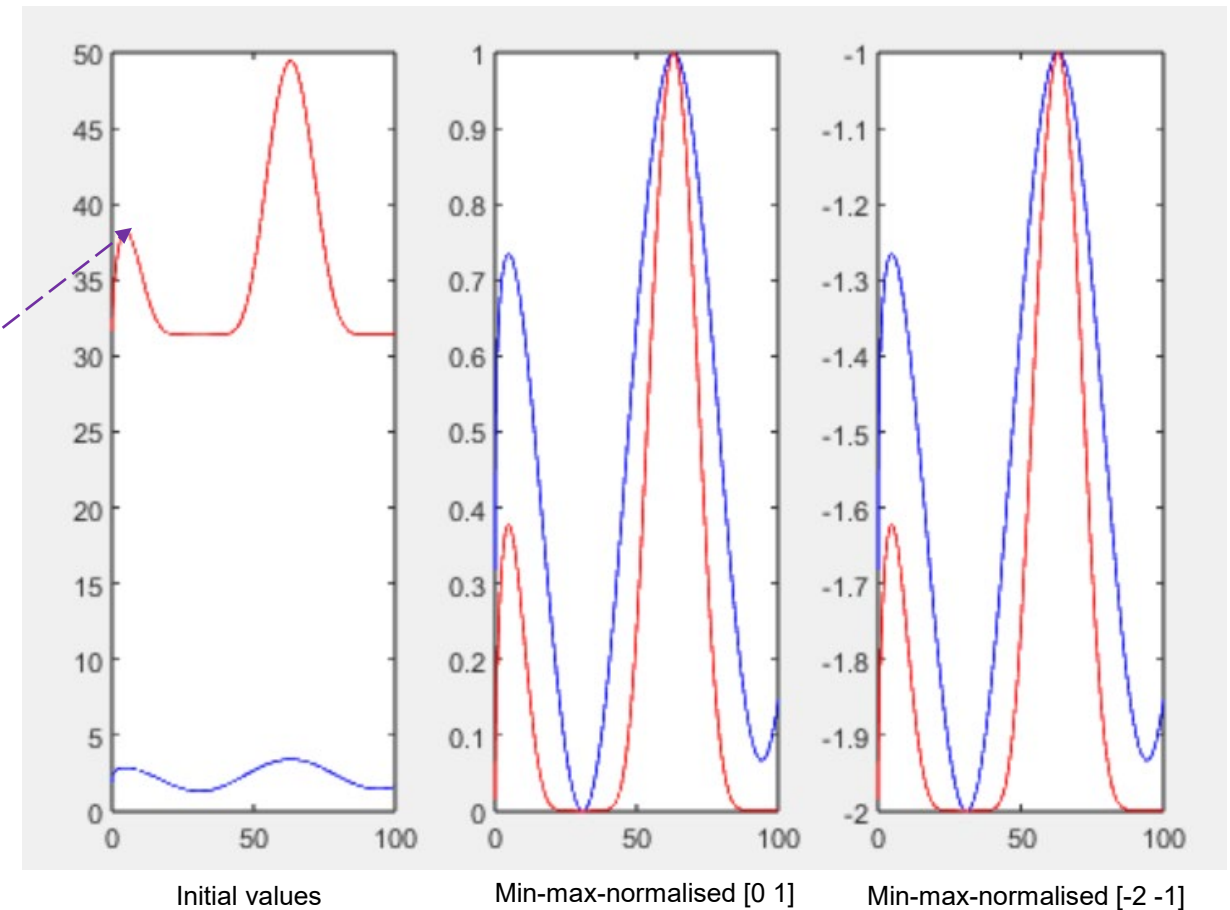- Is it possible to squeeze values between scalars $z_0$ and $z_1$ with $z_0 < z_1$?

Initial values          Min-max-normalised

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Application #1

1) Use Matlab to randomly generate a vector V of 100 scalars between -100 and 100.

2) Carry out the min-max normalization of this vector and store the result in Vn.

3) Set the value at index 20 in V to 500. Store the result in Vo.

4) Carry out the min-max normalization of Vo and store the result in Von.

5) Plot V, Vn and Von in three subplots. What do you observe?

$$\hat{x} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Data preparation -  Minimum-Maximum Normalization

- $x = [x_1 \ x_2 \ ... \ x_n]$ the vector of $n$ scalar data

- $x_{min}$ and $x_{max}$ is minimum and maximum value in $x$

- $\hat{x}$ is the min-max-normalized value of $x$:

$$\hat{x} = z_0 + (z_1 - z_0) \frac{x - x_{min}}{x_{max} - x_{min}}$$

- Values in $\hat{x}$ are squeezed between $z_0$ and $z_1$



Initial values          Min-max-normalised [0 1]          Min-max-normalised [-2 -1]

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Data preparation - Minimum-Maximum Normalization



Initial values | Min-max-normalised [0 1] | Min-max-normalised [-2 -1]

- – Min-max normalization faces outliers issues!
- – Shape (of the **distribution**) is preserved
- – How to influence **mean value** and **standard deviation**?

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Data preparation - z-Score Normalization

- $x = [x_1 \; x_2 \; ... x_n]$ the vector of $n$ scalar data

- Mean value

$$\mu = \frac{1}{n} \sum_{i=1}^{n} x_i$$

- Standard deviation

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \mu)^2}$$

- Normalization:

$$\hat{x} = \frac{x - \mu}{\sigma}$$

- **Hints**:
  - The feature $\hat{x}$ hat mean $\hat{\mu} = 0$ and standard dev. $\hat{\sigma} = 1$
  - Data are **not** necessary confined within [0 1] (during training), which is helpful for robustness (w.r.t **new** test data outside [-1 1])!



Initial values

Normalised

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Data preparation -  z-Score Normalization

- $x = [x_1 \; x_2 \ldots x_n]$ the vector of $n$ scalar data

- Mean value

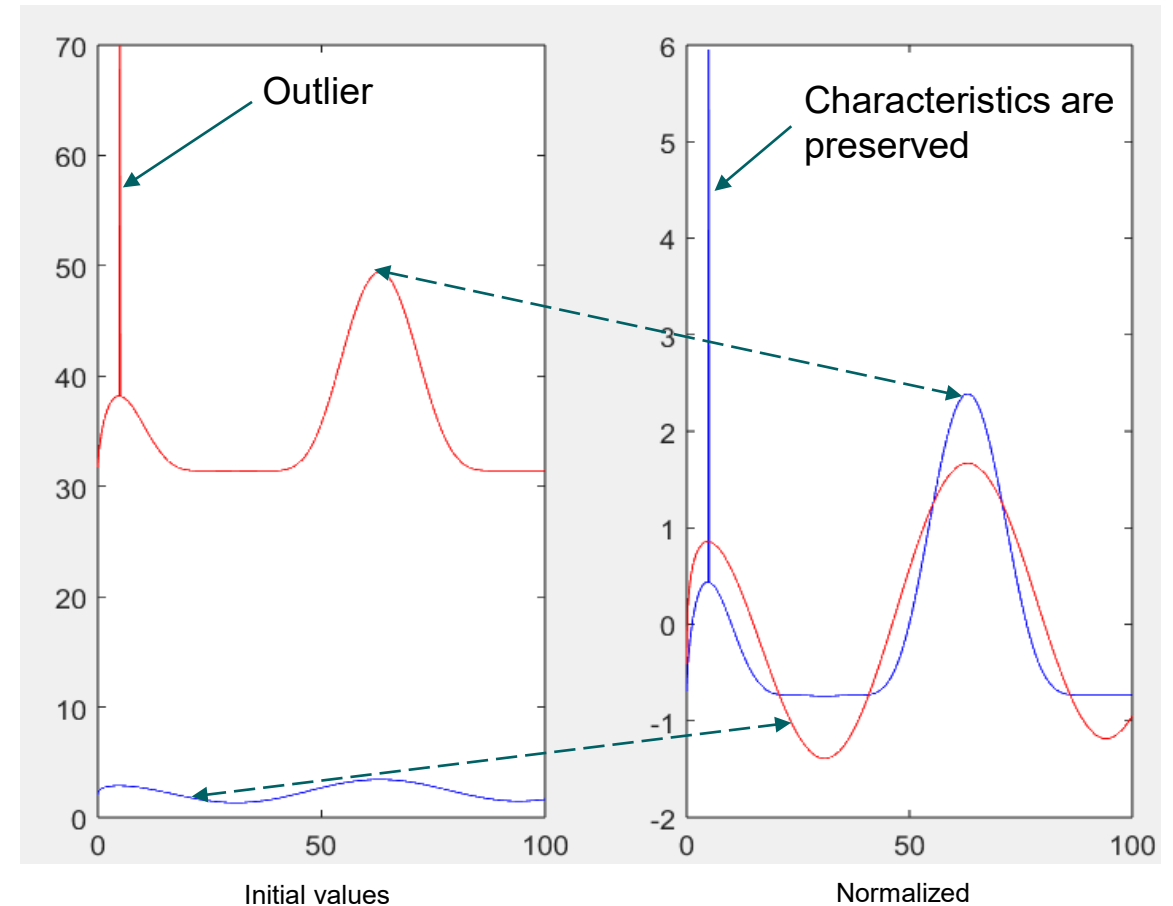$$\mu = \frac{1}{n}\sum_{i=1}^{n} x_i$$

- Standard deviation

$$\sigma = \sqrt{\frac{1}{n-1}\sum_{i=1}^{n}(x_i - \mu)^2}$$

- Normalisation:

$$\hat{x} = \frac{x - \mu}{\sigma}$$

- **Hints**:
  - Useful properties about outliers are kept!
  - Impact of outliers on normalized data is mitigated (hence robustness!)



Outlier

Characteristics are preserved

Initial values

Normalized

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Application #2:

- Assume that $u = [x_1, \cdots, x_n]^T$ is a dataset

- Calculate Min-Max(z-Score(x))

## Solution:

$$\hat{x} = \frac{x - x_{min}}{x_{max} - x_{min}} = \frac{\dfrac{x - \mu}{\sigma} - \dfrac{x_{min} - \mu}{\sigma}}{\dfrac{x_{max} - \mu}{\sigma} - \dfrac{x_{min} - \mu}{\sigma}} = \frac{x - x_{min}}{x_{max} - x_{min}}$$
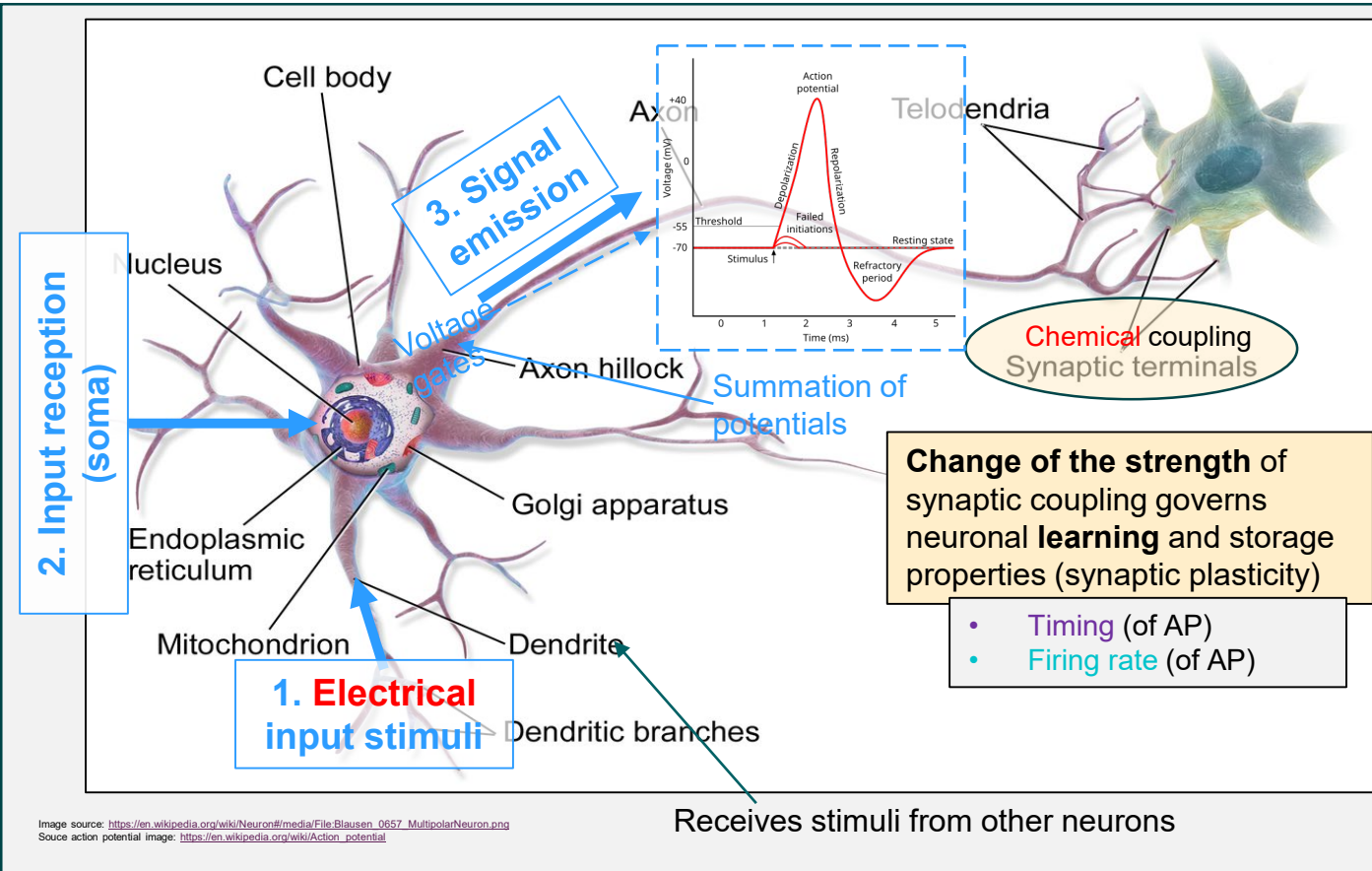
```
>> n=10; a=rand(n,1)'; zm = normalize(a,'range')- normalize(normalize(a),'range')

zm =

   1.0e-15 *

   -0.1110   -0.0555        0        0        0        0        0        0        0   -0.1110
```
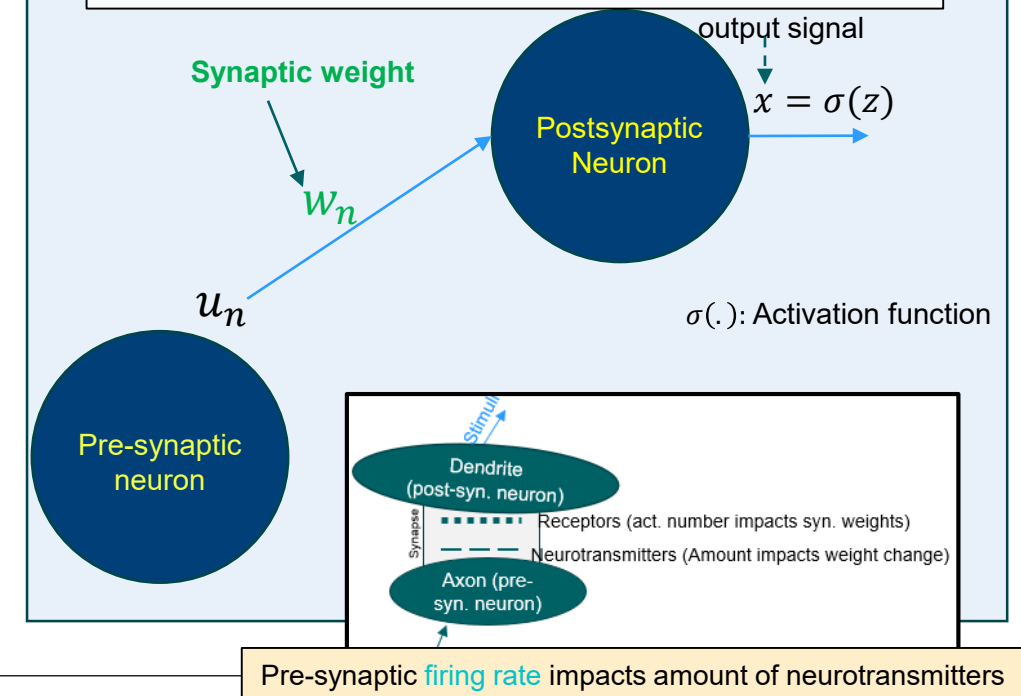
Z-score

Min-max
Between [0 1]

Between [0 1]

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Using neurons to reason

## Neuron model – Synapse dynamics



Cell body

Action potential

Axon

Telodendria

+40

Voltage (mV)

0

Threshold

-55

-70

Depolarization

Repolarization

Failed initiations

Stimulus ↑

Resting state

Refractory period

0 1 2 3 4 5
Time (ms)

**3. Signal emission**

Voltage gates

Summation of potentials

Nucleus

Axon hillock

Chemical coupling
Synaptic terminals

**2. Input reception (soma)**

Endoplasmic reticulum

Golgi apparatus

Mitochondrion

Dendrite

**1. Electrical input stimuli**

Dendritic branches

**Change of the strength** of synaptic coupling governs neuronal **learning** and storage properties (synaptic plasticity)

- Timing (of AP)
- Firing rate (of AP)

Receives stimuli from other neurons

Image source: https://en.wikipedia.org/wiki/Neuron#/media/File:Blausen_0657_MultipolarNeuron.png
Souce action potential image: https://en.wikipedia.org/wiki/Action_potential

### Spike-timing-dependent plasticity

- Process associated with the modification (optimization) of the synaptic weights
- Depends upon temporal dynamics (e.g., latency) between action potential (AP) in pre- and postsynaptic neurons (+/- strength of long-term potentiation resp. depression)
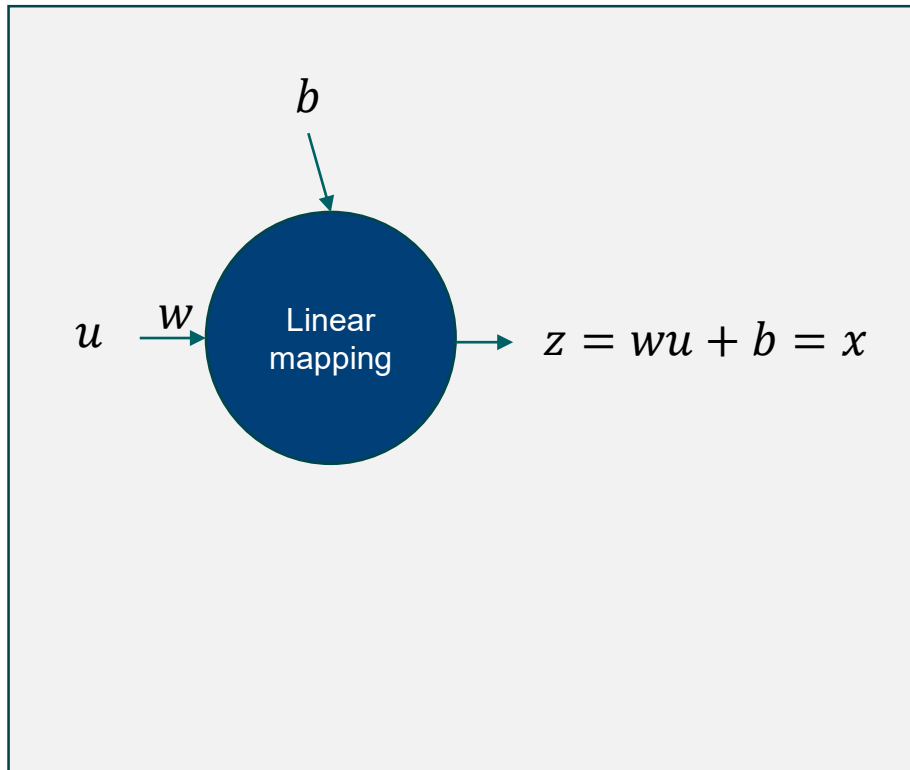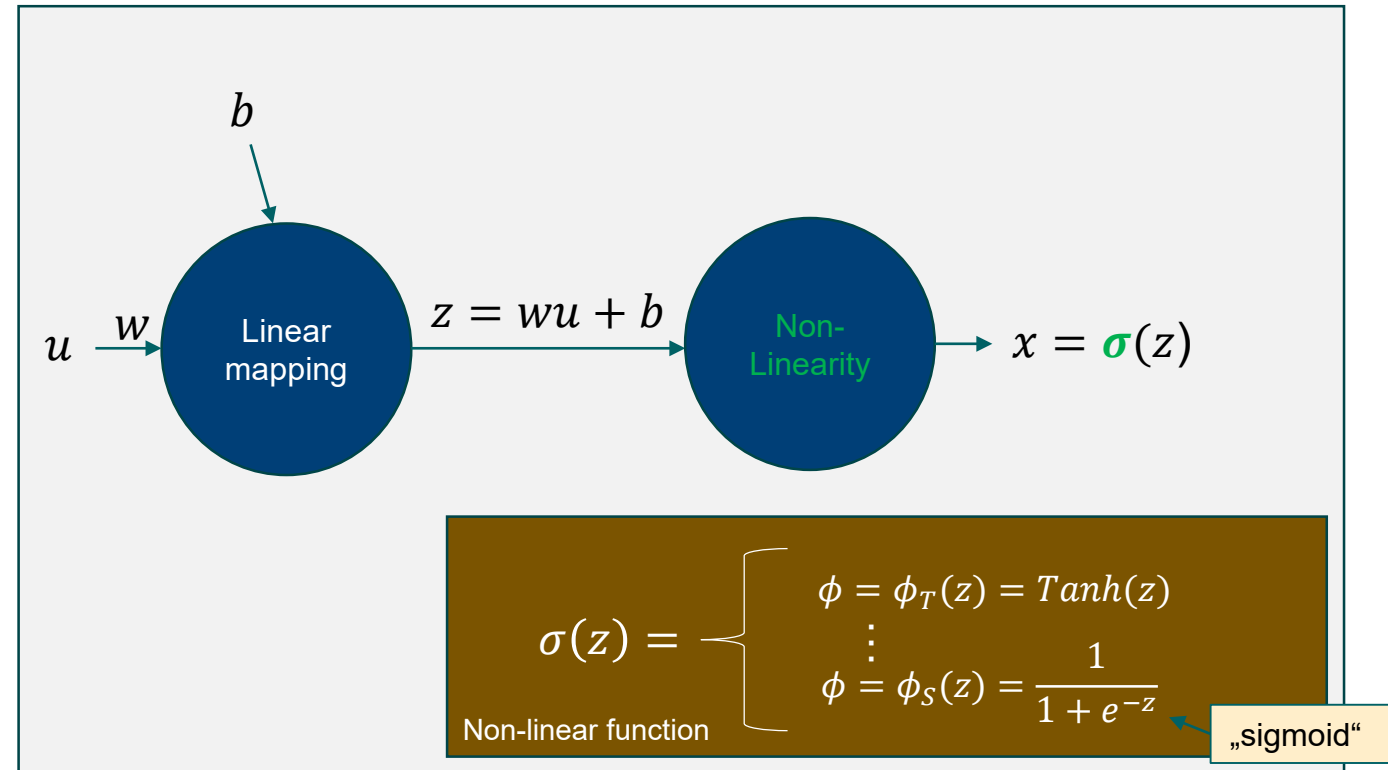
Synaptic weight

$w_n$

$u_n$

output signal

Postsynaptic Neuron

$x = \sigma(z)$

$\sigma(.)$: Activation function

Pre-synaptic neuron

Stimulus

Dendrite (post-syn. neuron)

Receptors (act. number impacts syn. weights)

Neurotransmitters (Amount impacts weight change)

Synapse

Axon (pre-syn. neuron)

Pre-synaptic firing rate impacts amount of neurotransmitters

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Neuron model – Synapse dynamics



2. Input reception (soma)

3. Signal emission

1. **Electrical input stimuli**

Voltage gates

Summation of potentials

Chemical coupling
Synaptic terminals

**Change of the strength** of synaptic coupling governs neuronal **learn**ing and storage properties (synaptic plasticity)

- Time coding
- Rate coding

Image source: https://en.wikipedia.org/wiki/Neuron#/media/File:Blausen_0657_MultipolarNeuron.png
Souce action potential image: https://en.wikipedia.org/wiki/Action_potential

Receives stimuli from other neurons



$b$ = (Bias)

Synaptic weight (used to **learn**!)

**Output signal**

$u_1$  $w_1$

$u_2$  $w_2$

$w_n$

$u_n$

$$z = b + \sum_{k=1}^{n} x_k w_k$$

$x = \sigma(z)$

$\sigma(.)$: Activation function

**Input signal**

**Neuron model**

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Artificial neuron



$$z = wu + b = x$$

Identity activation

$$z = wu + b$$

$$x = \boldsymbol{\sigma}(z)$$

Non-linear function

$$\sigma(z) = \begin{cases} \phi = \phi_T(z) = Tanh(z) \\ \vdots \\ \phi = \phi_S(z) = \dfrac{1}{1 + e^{-z}} \end{cases}$$

„sigmoid"

From identity activation to **non-linear** activation

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Data preparation and activation function

Saturation

**Input** data

**Output** data

Data preparation
• Min-max
• Z-score
• …etc

Learning system
(neuron/neural network)

$\phi$

Activation function $\phi$

$$\phi = \phi_T(u) = Tanh(u)$$

$$\phi = \phi_S(u) = \frac{1}{1 + e^{-u}}$$

Range [0 1] aligns with the notion of „**probability**" (for classification purposes)

**Data preparation helps**
• keep the input training data in the **range** of the activation function $\phi$
• **accelerate** the training process

Called „Sigmoid"

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Activation function - Saturation



Data preparation
• Min-max
• Z-score
• …etc

Learning system
(neuron/neural network)

**Output** data

$\phi$

Activation function $\phi$

Saturation

$$\phi = \phi_T(u) = Tanh(u)$$

Range [0 1] aligns with the notion of „**probability**" (for classification purposes)

$$\phi = \phi_S(u) = \frac{1}{1 + e^{-u}}$$

Called „Sigmoid"

### When saturated
• **Output** of the activation function $\phi$ does not respond to variations of its input
• learning capabilities drop, as potentially captured non-linear mappings are limited

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Activity function – Product of derivatives



"Sigmoid"

$$\phi = \phi_S(u) = \frac{1}{1 + e^{-u}}$$

"Sigmoid derivative"

$$\phi_S'(u) = (1 - \phi_S(u))\phi_S(u)$$



$\phi_S'(u)$

$\phi_S'(u)^* \phi_S'(u)$

$\phi_S'(u)^* \phi_S'(u) * \phi_S'(u)$

**Observation**: The amplitude of the product of the derivative $\phi_S'(u)$ by itself decreases with the increasing number of product terms.

Application#3: verify above-mentioned behavior ($\phi_S'(u) * \cdots * \phi_S'(u)$) of the sigmoid by writing and executing a matlab code for three product terms.

## Activation function – Fixing saturation issues

$$f(u) = u, \text{if } u \geq 0$$
$$f(u) = 0, \text{otherw.}$$

Rectified Linear Unit  (**ReLu**) Activation

$$f(u) = u, \text{if } u \geq 0$$
$$f(u) = a \cdot u, \text{otherw.}, a \in \mathbb{R}$$

Here $a = 0.5$. However, usually, $a \in [0.01, 0.1]$

Leaky Rectified Linear Unit   (**LReLu**)

- Not sensitive to negative inputs (might lead to dead neuron!)
- Unbounded for positive inputs

- Customized response (via $a$) to negative inputs (helps control the amplitude of derivative for negative inputs)
- Unbounded and no saturation hazard (thus: enhances learning rate during (optimization-base**d)** learning)

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Activation function – Fixing derivative issues

$\phi_T(u)$ provides a larger derivative (i.e., "gradient", more in short)



$\phi = \phi_T(u) = Tanh(u)$

Data preparation
• Min-max
• Z-score
• …etc

Learning system
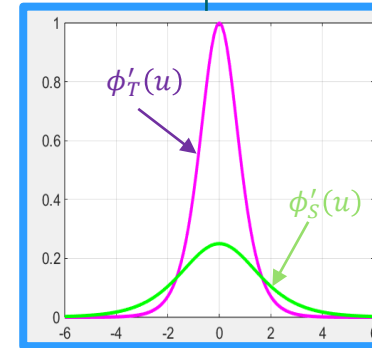(via parameter optimization)

$\phi$

Activation function $\phi$

Range [0 1] aligns with the notion of „probability"

$\phi = \phi_S(u) = \dfrac{1}{1 + e^{-u}}$

**Derivative of $\phi_T$ and $\phi_S$:**

$$\phi_T'(u) = 1 - \tanh^2(u)$$

„Sigmoid derivative"

$$\phi_S'(u) = \big(1 - \phi_S(u)\big)\phi_S(u)$$

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## What is an „optimization" and how does it work?

**Goal**: looking for a location with „maximum" or „minimum" function value



$f(u)$

Descent toward the obvious minimum

$u$

**Optimization goal**: find $u^*$ that e.g. minimizes the 1D (or n-D) loss function $f(u)$.
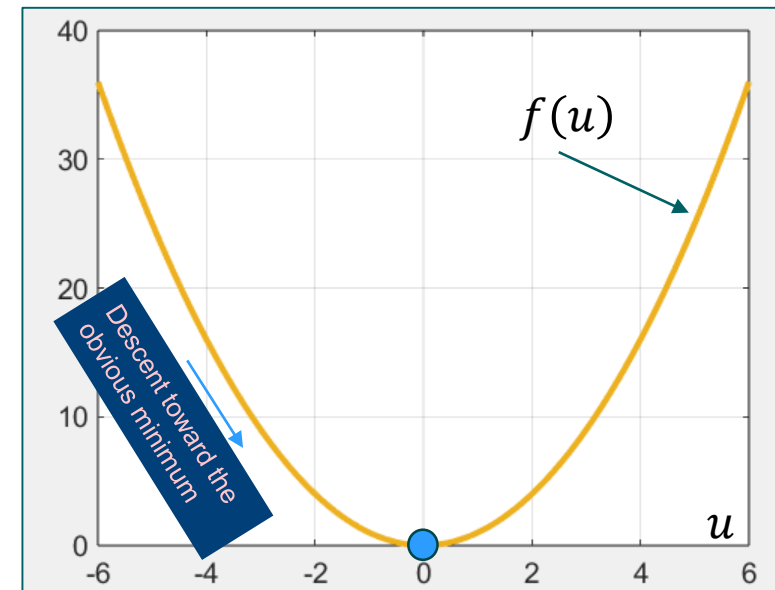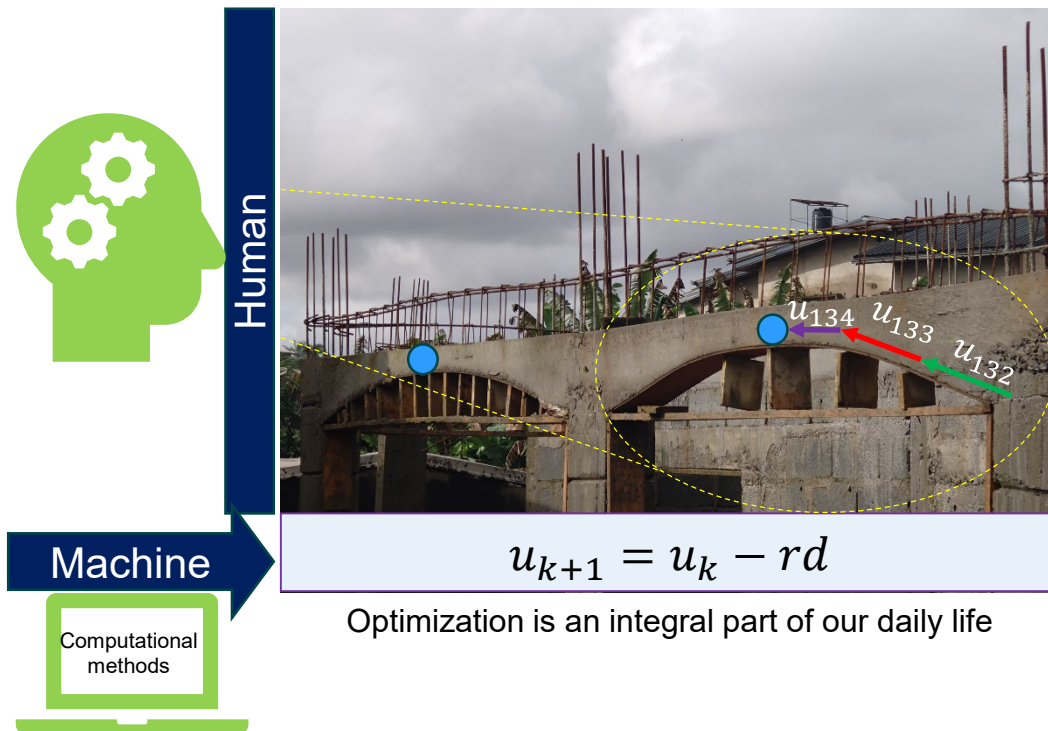


$u_{130}$ $u_{131}$ $u_{132}$

$$u_{k+1} = u_k - rd$$

Optimization is an integral part of our daily life

Whereas humans tend to perceive optimum in some cases from experiences (see, e.g., r.h.s), the gradient can help machines do so in a systematic way!

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## What is an „optimization" and how does it work?



$$u_{k+1} = u_k - rd$$

Optimization is an integral part of our daily life



$f(u)$

Descent toward the obvious minimum

**Optimization goal**: find $u^*$ that e.g. minimizes the 1D (or n-D) loss function $f(u)$.

While humans tend to perceive an optimum in some cases (see, e.g., l.h.s), the gradient $d$ can help machines do so in a systematic way with rate $r$!

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## What is a gradient?

- Preliminary: $x = f(u)$ is a real function of a real variable $u$.

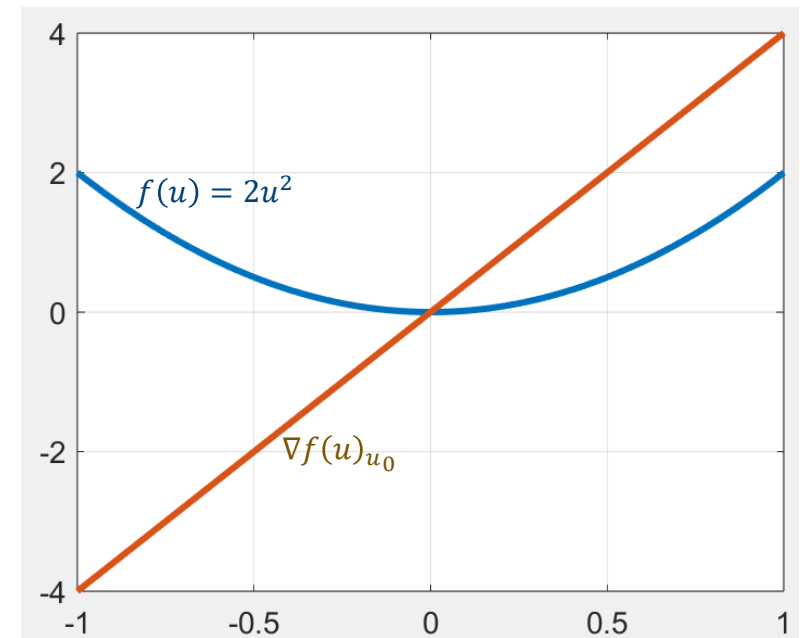- Gradient $\nabla f(u)_{u_0}$ of $f(u)$ w.r.t $u$ at position $u = u_o$ is given by

$$\nabla f(u)_{u_0} = \left.\frac{\partial f(u)}{\partial u}\right|_{u=u_0}$$

> $f$ is differentiable at $u_0$

1.Order derivative w.r.t $u$ at $u = u_o$

**In the *n*-dimensional space: $f : \mathbb{R}^n \rightarrow \mathbb{R}$**

- $u = (u_1, \cdots, u_n)$

- $\nabla f(u)_{u_0} = \left[\frac{\partial f(u)}{\partial u_1}, \cdots, \frac{\partial f(u)}{\partial u_n}\right]^T_{u=u_0}$



$f(u) = 2u^2$

$\nabla f(u)_{u_0}$

Gradient example

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## What is a gradient?

- Preliminary: $x = f(u)$ is a real function of a real variable $u$.

- Gradient $\nabla f(u)_{u_0}$ of $f(u)$ w.r.t $u$ at position $u = u_o$ is given by

$$\nabla f(u)_{u_0} = \frac{\partial f(u)}{\partial u}\bigg|_{u=u_0}$$

1.Order derivative w.r.t $u$ at $u = u_o$

**Application #4:**

- $f(u) = u^2 \rightarrow \nabla f(u)_{u0=1} = \frac{\partial f(u)}{\partial u}\big|_{u=1} = \frac{\partial u^2}{\partial u}\big|_{u=1} = 2u\big|_{u=1} = 2 \cdot 1 = 2$

- $f(u) = \frac{1}{u-1} \rightarrow \nabla f(u)_{u0=0} = \frac{\partial f(u)}{\partial u}\big|_{u=0} = \frac{0-1}{(u-1)^2}\big|_{u=0} = -1$

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Application #5

**Challenge**:

$$x = f(u_1, u_2) = 2u_1^3 + \ln(u_2)$$

Calculate the gradient of $f(u_1, u_2)$ at $u = (u_1, u_2)$

**Solution:**

$$x = f(u_1, \dots, u_n) \rightarrow \nabla f = \begin{bmatrix} \dfrac{\partial f}{\partial u_1} \\ \vdots \\ \dfrac{\partial f}{\partial u_n} \end{bmatrix} \qquad \Longrightarrow \qquad \nabla f = \begin{bmatrix} 6u_1^2 \\ \dfrac{1}{u_2} \end{bmatrix}$$

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Application #6

### Challenge:

$$x = f(u_1, u_2) = -u_2^2 + u_2 + 17$$

Calculate the gradient of $f(u_1, u_2)$ at $u = (u_1, u_2)$

### Solution:

$$x = f(u_1, \dots, u_n) \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial u_1} \\ \vdots \\ \frac{\partial f}{\partial u_n} \end{bmatrix}$$

$$\nabla f = \begin{bmatrix} 0 \\ -2u_2 + 1 \end{bmatrix}$$

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Application #7

### Challenge:

$$x = f(u_1, u_2) = -u_1 u_2^2 + u_1 u_2 + 17 u_2$$

Calculate the gradient of $f(u_1, u_2)$ at $u = (u_1, u_2)$

### Solution:

$$x = f(u_1, \ldots, u_n) \rightarrow \nabla f = \begin{bmatrix} \dfrac{\partial f}{\partial u_1} \\ \vdots \\ \dfrac{\partial f}{\partial u_n} \end{bmatrix}$$

$$\nabla f = \begin{bmatrix} -u_2^2 + u_2 \\ -2u_1 u_2 + u_1 + 17 \end{bmatrix}$$

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Application #8

### Challenge:

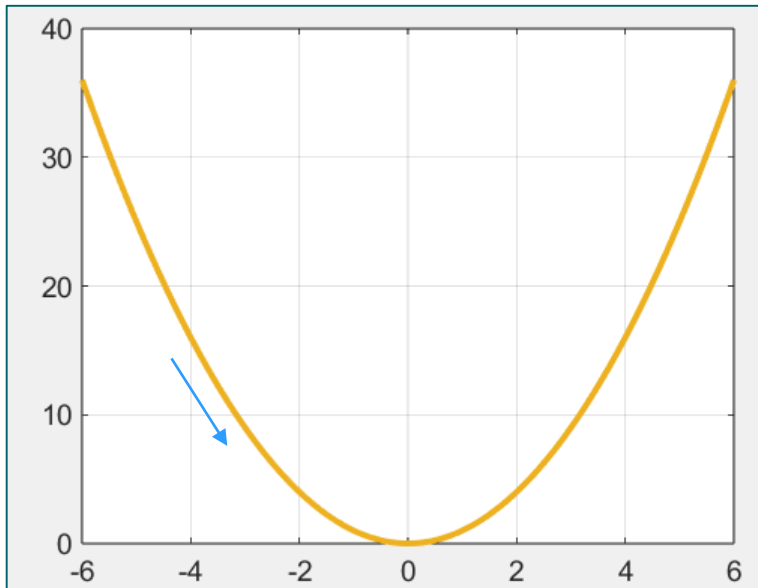$$x = f(u_1, u_2, u_3) = -u_1 u_2^2 u_3 + u_1 u_2 + 17 u_2 + u_3$$

Calculate the gradient of $f(u_1, u_2, u_3)$ at $u = (u_1, u_2, u_3)$

### Solution:

$$x = f(u_1, \dots, u_n) \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial u_1} \\ \vdots \\ \frac{\partial f}{\partial u_n} \end{bmatrix} \qquad \Longrightarrow \qquad \nabla f = \begin{bmatrix} -u_2^2 u_3 + u_2 \\ -2 u_1 u_3 u_2 + u_1 + 17 \\ -u_1 u_2^2 + 1 \end{bmatrix}$$

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Why is the gradient useful?

- Suppose that we have a function $x = f(u) = u^2$



How to find $u^*$ that minimize $f(u)$?

**Approach**:

$$u_{k+1} = u_k - \underbrace{r \nabla f(u)_{u_k}}_{= d} = u_k - r \left. \frac{\partial f(u)}{\partial u} \right|_{u=u_k}$$

- $r$ is a **hyperparameter** that reflects the **learning rate**
- The **fastest decay** of $f(u)$ is the **opposite direction** of the **gradient** $\nabla f(u)$ of $f(u)$
- Hence, $-r \nabla f(u)$ steers the series until the convergence $u_{k+1} \to u^*$

Source: https://cdn.britannica.com/32/124632-004-1C08C796/craters-Mount-Cameroon.jpg



**Caution: Many local optima might exist**

## Optimization routine based upon the gradient

Goal: find minimum of $x = f(u) = u^2$
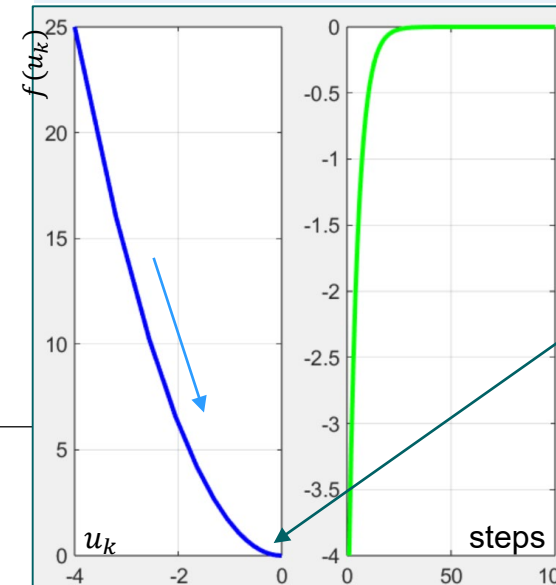
```
1    clear;
2    syms u                  %symbolic variable
3    f =u^2;                 %function to be optimized
4    c = gradient(f,u); %symbolic gardient
5    uk = -5;            % initial uk value
6    r = 0.1;           % learning rate
7    interval=1:100; % 100 steps
8    U=[]; functionValue=[];%container
9    for i=interval
10     u = uk;
11     uk = uk - r*subs(c); %main iteration loop
12     U=[U uk];          % for plotting
13     functionValue = [functionValue (subs(f))];
14   end
15   subplot(1,2,1); plot(U,functionValue,'b');
16   subplot(1,2,2); plot(interval,U,'g');
```

You are welcome to optimize and share the code in our lecture forum!

**Approach**:

$$u_{k+1} = u_k - r\nabla f(u)_{u_k} = u_k - r\frac{\partial f(u)}{\partial u}\bigg|_{u=u_k}$$

- $r$ is a **hyperparameter** that reflects the learning rate
- The **fastest decay** of $f(u)$ is the **opposite direction** of the **gradient** $\nabla f(u)$ of $f(u)$
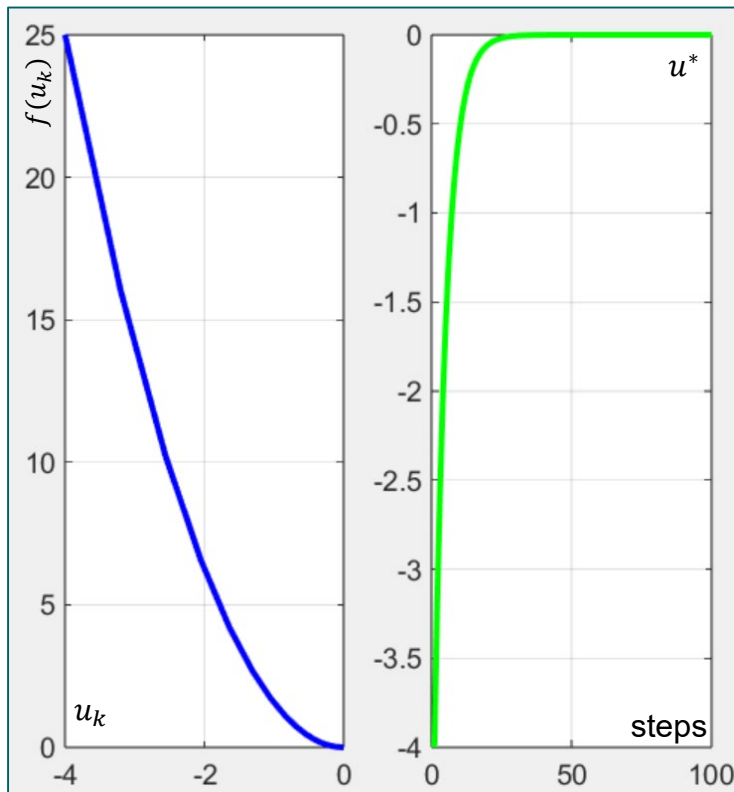- Hence, $-r\nabla f(u)$ steers the series until the convergence $u_{k+1} \rightarrow u^*$

**Observations:**
- Goal-oriented decay toward minimum value
- Minimum position $u^*$ is attained
- However …

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

**We are aware of the gradient usefulness. The path toward usefulness can be struggeling! Why?**

- Minimizing $x = f(u) = u^2$ to find minimum $u = u^*$



Learning rate $r = 0.1$. Optimization remains smooth and fast.

Learning rate $r = 0.97$. Optimization takes longer to converge while oscillating!

Eric Kaig
kaigom@fb2.fra-uas.de

**Observation #1:** Values of hyperparameters $r$ **impact** optimization processes!

FRANKFURT
UNIVERSITY
PPLIED SCIENCES

# Linear regression

## Objective of regression – Find a function that captures the set of points



- Input: $u = [-5.3\ -2.2\ 0.5\ 1.7\ 2.5\ 4.2\ 6.9\ 9.2]^T$;

- Output: $x = [-9.8\ -3.6\ 1.70\ 4.1\ 5.7\ 9.1\ 14.5\ 19.1\ ]^T$;

- Points A (u,x) potentially line up with a line L

- If so, then approxim. of L should read: $\hat{x} = wu + b$

- $w$ (a scalar) is the slope, given by the **gradient**

- $b$ (also a scalar) is the „bias"

## Objective of regression – Find a function that captures the set of points



- $u = [-5.3\ -2.2\ 0.5\ 1.7\ 2.5\ 4.2\ 6.9\ 9.2\ ]^T;$

- $x = [-5.6\ 0.6\ 6.0\ 8.4\ 10.0\ 13.4\ 18.8\ 23.4\ ]^T;$

$\dfrac{\partial f}{\partial u_1}$ is growth rate of $f([u_1 \cdots u_n])$ in the $u_1$ direction

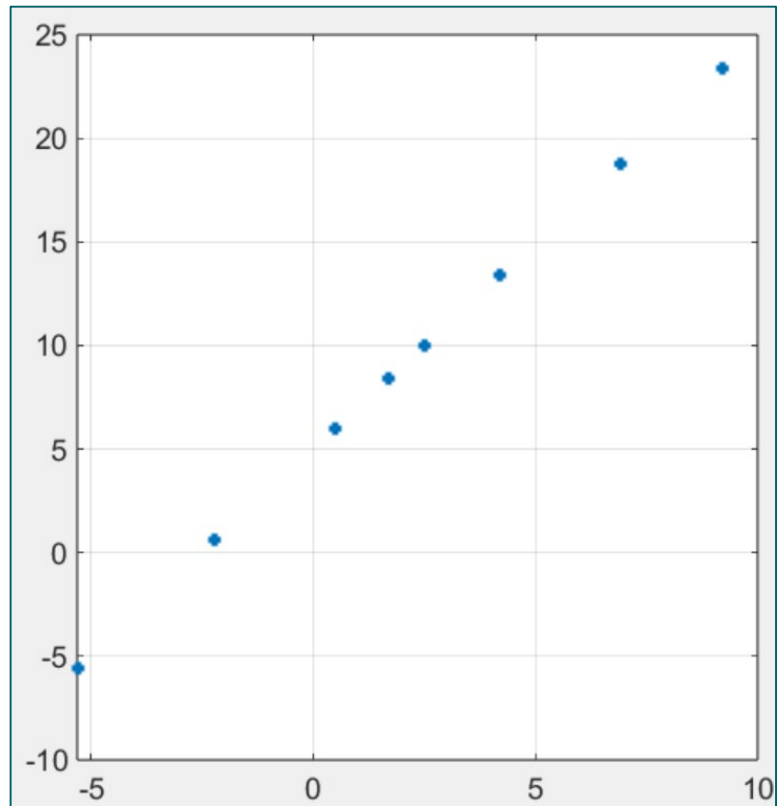prediction

**Hint:**

$x = f(u_1, \dots, u_n) \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial u_1} \\ \vdots \\ \frac{\partial f}{\partial u_n} \end{bmatrix}$ is a vector of $\mathbb{R}^n$!

- Line $\hat{x} = wu + b$

- Loss $L(w,b) = \sum_{i=1}^n (\hat{x}_i - x_i)^2$

**Goal:** learning $w$ and $b$ that min. $L(w,b)$

$\nabla f$ is the direction of the steepest **ascent** (hence, $-\nabla f$ is the direction of steepest **descent**) at $u$.

- Gradient $\nabla L(w,b) = \begin{bmatrix} \frac{\partial L(w,b)}{\partial w} \\ \frac{\partial L(w,b)}{\partial b} \end{bmatrix} = \sum_{i=1}^n \begin{bmatrix} 2u(\hat{x}_i - x_i) \\ 2(\hat{x}_i - x_i) \end{bmatrix}$

- $u_{k+1} = u_k - r\nabla L(u)_{u_k}, u_k = [w_k\ b_k]^T$

**Hint:**

$\dfrac{\partial L}{\partial w} = \dfrac{\partial L}{\partial \hat{x}} \dfrac{\partial \hat{x}}{\partial w}$

$\dfrac{\partial L}{\partial b} = \dfrac{\partial L}{\partial \hat{x}} \dfrac{\partial \hat{x}}{\partial b}$

Eric Kaigom
kaigom@fb2.fra-uas.de

OF APPLIED SCIENCES

## Optimization pitfall



Source: https://cdn.britannica.com/32/124632-004-1C08C796/craters-Mount-Cameroon.jpg

Eric Kaigom
kaigom@fb2.fra-uas.de

# Robotics and Machine Learning – Model Capture

## When to update the parameters?

- Line $\hat{x} = wu + b$

- Loss $L(w,b) = \sum_{i=1}^{n}(\hat{x}_i - x_i)^2$

- Gradient $\nabla L(w,b) = \begin{bmatrix} \frac{\partial L(w,b)}{\partial w} \\ \frac{\partial L(w,b)}{\partial b} \end{bmatrix} = \sum_{i=1}^{n} \begin{bmatrix} 2u(\hat{x}_i - x_i) \\ 2(\hat{x}_i - x_i) \end{bmatrix}$

- $u_{k+1} = u_k - r\nabla L(u)_{u_k}, u_k = [w_k \ b_k]^T$

**Hint:**

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{x}}\frac{\partial \hat{x}}{\partial w}$$
$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{x}}\frac{\partial \hat{x}}{\partial b}$$

---

**Gradient descent:**

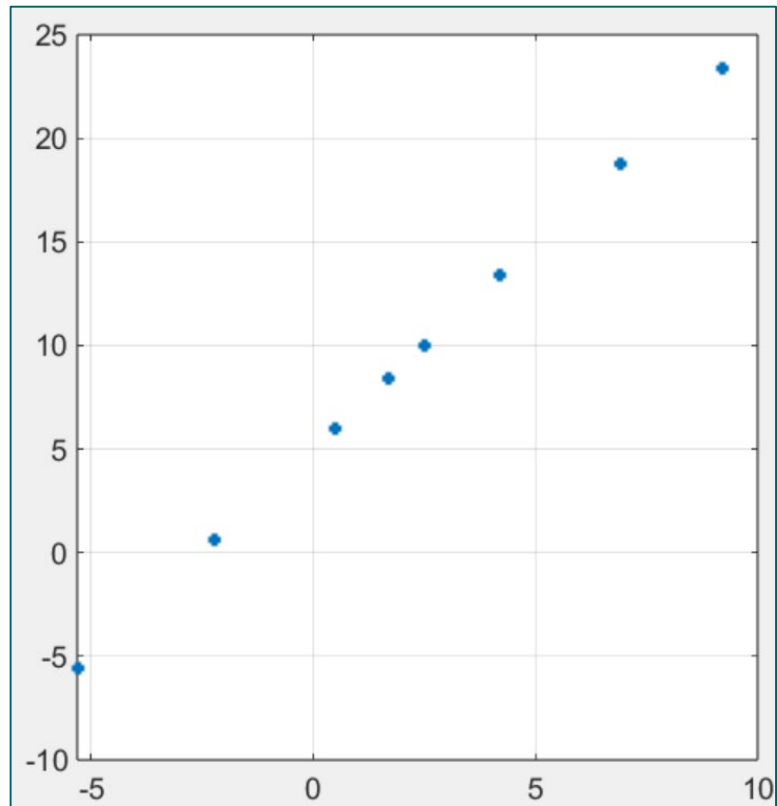Update after a loss evalation related to all training data samples (after each **epoch**)

---

**Stochastic gradient descent:**

Update after the loss related to a training data sample has been computed

---

**Mini-batch stochastic gradient descent:**

Update after the loss related to a number of training data samples has been computed

---

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Application #9: Write and execute a Matlab code that learns w and b!



- u = [-5.3 -2.2 0.5 1.7 2.5 4.2 6.9 9.2 ]$^T$;

- x = [-5.6 0.6 6.0 8.4 10.0 13.4 18.8 23.4 ]$^T$;

$\frac{\partial f}{\partial u_1}$ is growth rate of $f([u_1 \cdots u_n])$ in the $u_1$ direction

prediction

**Hint:**

$x = f(u_1, \ldots, u_n) \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial u_1} \\ \vdots \\ \frac{\partial f}{\partial u_n} \end{bmatrix}$ is a vector of $\mathbb{R}^n$!

- Line $\hat{x} = wu + b$

- Loss $L(w,b) = \sum_{i=1}^{n}(\hat{x}_i - x_i)^2$

$\nabla f$ is the direction of the steepest **ascent** (hence, $-\nabla f$ is the direction of steepest **descent**) at $u$.

**Goal:** learning $w$ and $b$ that min. $L(w,b)$

- Gradient $\nabla L(w,b) = \begin{bmatrix} \frac{\partial L(w,b)}{\partial w} \\ \frac{\partial L(w,b)}{\partial b} \end{bmatrix} = \sum_{i=1}^{n} \begin{bmatrix} 2u(\hat{x}_i - x_i) \\ 2(\hat{x}_i - x_i) \end{bmatrix}$

- $u_{k+1} = u_k - r\nabla L(u)_{u_k}, u_k = [w_k\ b_k]^T$

**Hint:**

$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial w}$

$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{x}} \frac{\partial \hat{x}}{\partial b}$

Eric Kaigom
kaigom@fb2.fra-uas.de

OF APPLIED SCIENCES

# Robotics and Machine Learning – Model Capture

## Objective of regression – Find a function that captures the set of points



Ground truth

The gradient-based routine has captured (i.e. **learned (!)**) the line using the set of input and output points!

Once we have the weights $w$ and $b$, we have **captured** the mapping between input and output data!

## Objective of regression – Find a function that captures the set of points



```
1   clear; clc;
2   input = [-5.3 -2.2 0.5 1.7 2.5 4.2 6.9 9.2];
3   output=[-5.6 0.6 6.0 8.4 10.0 13.4 18.8 23.4]; %ou = 2*in+5
4   uk=[0;0];
5   r = 0.001; U = []; w = 0; b = 0;
6   for i = 1:10000
7   uk = uk -r*gradientloss(uk,input,output);
8   U = [U uk]; w = uk(1); b = uk(2);
9   end
10
11  subplot(1,2,1); plot(U(1,:),U(2,:));
12  subplot(1,2,2); plot(input,output,'+'); hold on; plot(input,w*input+b,'o');
13
14  function [r] =gradientloss (uk, in, out)
15  r1 = sum(2*(uk(1)*in+uk(2)-out).*in);
16  r2 = sum(2*(uk(1)*in+uk(2)-out)    );
17  r=[r1;r2];
18  end
```

The gradient-based routine has captured (i.e. learned) the line using the set of input and output points!

**Observation #2:** **Gradient** helps **LEARN** the mapping (Here: $\hat{x} = wu + b$) using input & output data! **However…**

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Objective – Find a function that captures the set of points



- u = [-5.3 -2.2 0.5 1.7 2.5 4.2 6.9 9.2];

- x = [-4.0 1.0 1.9 9.0 9.90 12.9 17.0 22.9];

---

- Still a line $\hat{x} = wu + b$ ?!

- Loss $L(w,b) = \sum_{i=1}^{n}(\hat{x}_i - x_i)^2$

- Gradient $\nabla L(w,b) = \begin{bmatrix} \frac{\partial L(w,b)}{\partial w} \\ \frac{\partial L(w,b)}{\partial b} \end{bmatrix} = \sum_{i=1}^{n}\begin{bmatrix} 2u(\hat{x}_i - x_i) \\ 2(\hat{x}_i - x_i) \end{bmatrix}$

- $u_{k+1} = u_k - r\nabla L(u)_{u_k}, u_k = [w_k\ b_k]^T$

**Hint:**

$$x = f(u_1, \ldots, u_n) \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial u_1} \\ \vdots \\ \frac{\partial f}{\partial u_n} \end{bmatrix}$$

**…how good** does this approach with **linearity assumption** perform?!

Obviously demanding non-linearities, as they happen in REAL LIFE!

**Recall**: the real world (and thus application data…) is (mostly) non-linear.

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Objective – Find a function that captures the set of points



Dynamics of optimization of weights

Predicted data
(i.e. from $\hat{x} = wu + b$)

Considerable deviation!!

Training data
(i.e., used to compute the weights w and b using gradient)

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Objective – Find a function that captures the set of points



- u = [-1.5967  -0.9356  -0.3599  -0.1040   0.0666   0.4292   1.0049   1.4954];
- x = [-0.0990   3.2492   4.7410   4.9784   4.9911   4.6316   2.9802   0.5275];

- Obviously not a line $\hat{x} = wu + b$!

**Hint:**
$$x = f(u_1, \ldots, u_n) \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial u_1} \\ \vdots \\ \frac{\partial f}{\partial u_n} \end{bmatrix}$$

- Loss $L(w, b) = \sum_{i=1}^{n}(\hat{x}_i - x_i)^2$

- Gradient $\nabla L(w, b) = \begin{bmatrix} \frac{\partial L(w,b)}{\partial w} \\ \frac{\partial L(w,b)}{\partial b} \end{bmatrix} = \sum_{i=1}^{n} \begin{bmatrix} 2u(\hat{x}_i - x_i) \\ 2(\hat{x}_i - x_i) \end{bmatrix}$

- $u_{k+1} = u_k - r\nabla L(u)_{u_k}, u_k = [w_k \; b_k]^T$

**…how good** does this approach with **linearity assumption** perform?!

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

**Objective – Find a function that captures the set of points**



Dynamics of optimization of weights

Loss

Training and prediction data

Ground truth

Optimized linear model

Obviously, the linear model is no longer suitable! We need to upskill our model… But how?! Any idea?!

**Observation #3:** The performance of $\hat{x} = wu + b$ drops as non-linearities and outliers are involved in training data!

UNIVERSITY OF APPLIED SCIENCES

# Learning using a neuron

## 1 neuron case : Optimization of the weights using (Input data, Output data) and gradient

$$\text{Input } u = \begin{bmatrix} 1.887 \\ \vdots \\ -0.288 \end{bmatrix}$$

$u_1$, $u_n$

$b$

$w$

Neuron

$$\text{Output } x = \begin{bmatrix} -0.774 \\ \vdots \\ 0.887 \end{bmatrix}$$

$x_1$, $x_n$

Label data

Training set of data examples = $(u_i, x_i)$, $i \in \{1, \cdots, n\}$

**Objectives**: How to find suitable $w$ and $b$, two scalars, that (learn to) **associate (or map)** multiple input examples in $u$ with corresponding output examples (i.e., labels) in $x$ ?

.. and with good generalization!

**Approach:**

$u$ → neuron → $x$   (a)

(b)
- Loss gradient computation using input/output/state data
- Gradient descent **optimization**

Computed using the training examples $(u, x)$

Values of a $w$ and a $b$ that minimize Loss function

(c) Trained neuronal model (Mapping)!

(d) new $u$ not seen before → Learned mapping as the AI model → Accurately predicted $x$

Eric Kaigom
kaigom@fb2.fra-uas.de

52

Note: Similar distribution usually assumed (realistic?). However, out of the scope of this lecture!

(Use in tests/applications testdata **normalized** with same normalization parameters as training)

T UNIVERSITY

## Neuron-related gradient computation

**Objective:** How to compute $\frac{\partial L}{\partial w}$ ?

Artificial non-linear neuron

Training label data
(„ground truth")

$x_T$

$b$

Neuron weight

$u$   $w$   Linear mapping   $z = wu + b$   Activation $\sigma$   $x = \boldsymbol{\sigma}(z)$

Neuron input

Neuron output.
**Nonlinearity**
(**important**) is
added via $\sigma$ to
capture more
features

e.g., $L = 0.5(x - x_T)^2$

$$\frac{\partial L}{\partial w} = ?$$

Observe that: $\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x} \frac{\partial x}{\partial z} \frac{\partial z}{\partial w}$

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Neuron-related gradient computation

**A** Forward step: Computing the loss function $L(x(u), x_m)$

**Objective:** How to compute $\frac{\partial L}{\partial w}$ ?

**Rationale:** We need the derivative (i.e., $\frac{\partial L}{\partial w}$) to be aware of the **direction of fastest decrease** (given by $-\frac{\partial L}{\partial w}$) of $L$ as a function of $w$ to update $w_{k+1}$. Recall that we strive to find
$$w_{k+1} = w_k - r\frac{\partial L}{\partial w} \rightarrow w^*$$
that minimizes $L$!

$b$

Neuron training output

Measured training data (the labels)

$u$ $w$ Linear mapping $z = wu + b$ Activa-tion $\sigma(z)$ $x = \sigma(z) \longrightarrow$ Loss $= L(x(u), x_m)$

Neuron training data

A neuron made up of its linear $(z)$ and non-linear $(\sigma(z))$ parts

This has implications:

**How? Observe that:**
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z}\frac{\partial z}{\partial w}$$

$\frac{\partial L}{\partial w}$ $\longleftarrow$ $\frac{\partial z}{\partial w}$ $\longleftarrow$ $\frac{\partial x}{\partial z}$ $\longleftarrow$ $\frac{\partial L}{\partial x}$

Backward step: Computing the **gradient** of loss function $L(x(u), x_m)$ **B**

Eric kaigo

**Observation**: Gradient $\frac{\partial L}{\partial w}$ computed (to optimize $w$) through a **backward propagation**!

## Neuron-related gradient computation

Forward step: Computing the loss function $L(x(u), x_m)$

**Objective:** How to compute $\frac{\partial L}{\partial w}$ ?

**Application: a**ssume that
- $L(x, x_m) = \frac{1}{2}(x - x_m)^2$
- $\sigma(z) = \frac{1}{1+e^{-z}}$

**Rationale:** We need the derivative (i.e., $\frac{\partial L}{\partial w}$) to be aware of the **direction of fastest decrease** (given by $-\frac{\partial L}{\partial w}$) of $L$ as a function of $w$ to update $w_{k+1}$. Recall that we strive to find
$$w_{k+1} = w_k - r\frac{\partial L}{\partial w} \to w^*$$
that minimizes $L$!

$b$

$u$ $\xrightarrow{w}$ Linear mapping $\xrightarrow{z = wu + b}$ Activation $\sigma(z)$ $\rightarrow$ $x = \sigma(z) \longrightarrow$ Loss $= L(x(u), x_m)$

**How? Observe that:**
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z}\frac{\partial z}{\partial w}$$
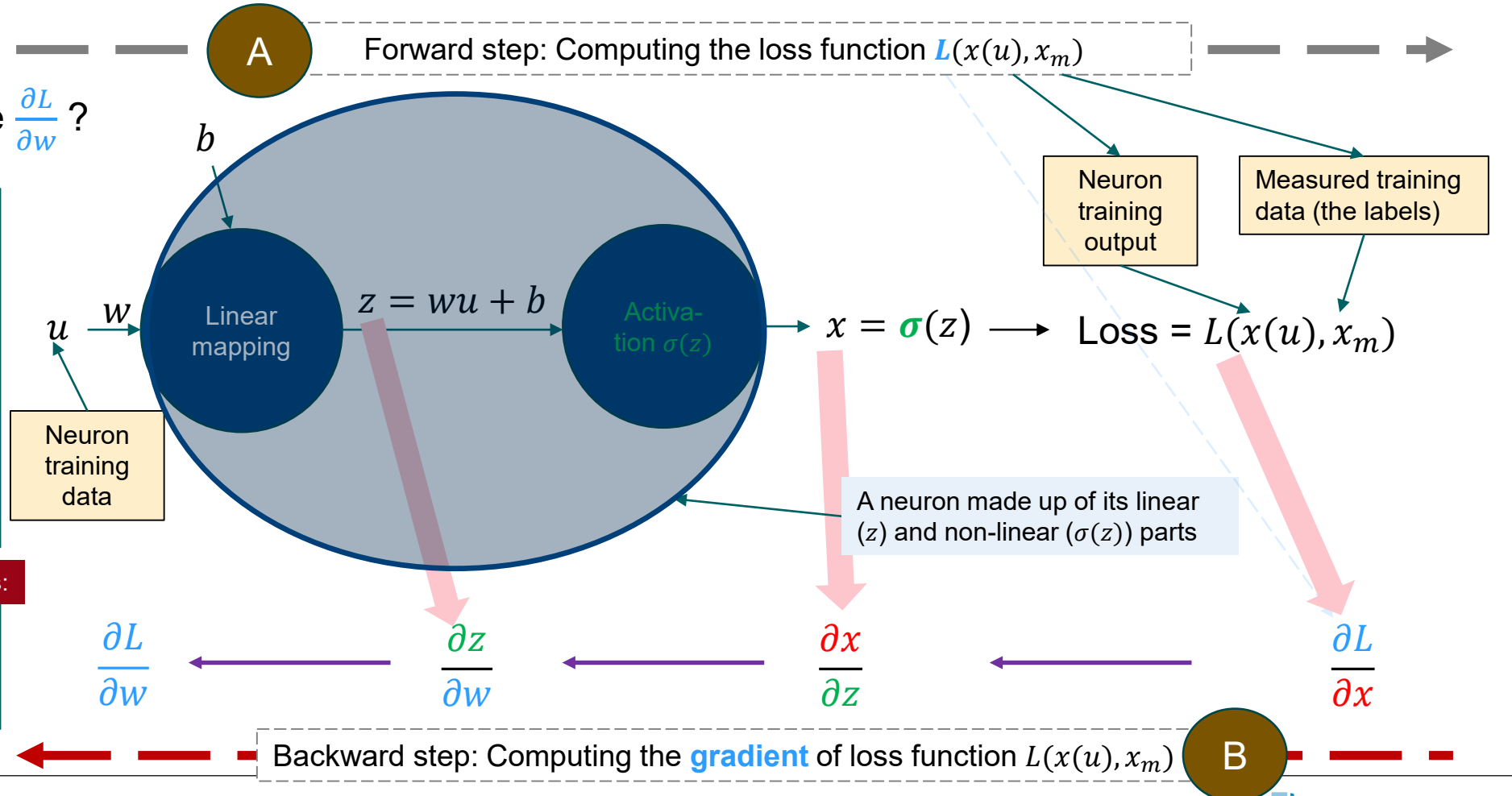
$\frac{\partial L}{\partial w} \longleftarrow \frac{\partial z}{\partial w} \longleftarrow \frac{\partial x}{\partial z} \longleftarrow \frac{\partial L}{\partial x}$

Backward step: Computing the **gradient** of loss function $L(x(u), x_m)$

Eric K
kaigor

$$\frac{\partial L}{\partial w} = (x - x_m)\sigma(z)(1 - \sigma(z))u$$

$$\frac{\partial z}{\partial w} = u$$

$$\frac{\partial x}{\partial z} = \sigma(z)(1 - \sigma(z))$$

$$\frac{\partial L}{\partial x} = (x - x_m)$$

OF APPLIED SCIENCES

## Neuron-related gradient computation

**Objective:** How to compute $\frac{\partial L}{\partial w}$ ?

**Application:** assume that
- $L(x, x_m) = (x - x_m)^2$
- $\sigma(z) = \tanh(z)$

**Rationale:** We need the derivative (i.e., $\frac{\partial L}{\partial w}$) to be aware of the **direction of fastest decrease** (given by $-\frac{\partial L}{\partial w}$) of $L$ as a function of $w$ to update $w_{k+1}$.

Recall that we strive to find
$$w_{k+1} = w_k - r\frac{\partial L}{\partial w} \to w^*$$
that minimizes the loss $L$!

$b$

$u \xrightarrow{w}$ Linear mapping $\xrightarrow{z = wu + b}$ Activation $\sigma(z)$ $\longrightarrow$ $x = \sigma(z) \longrightarrow$ Loss $= L(x(u), x_m)$

**How? Observe that:**
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z}\frac{\partial z}{\partial w}$$
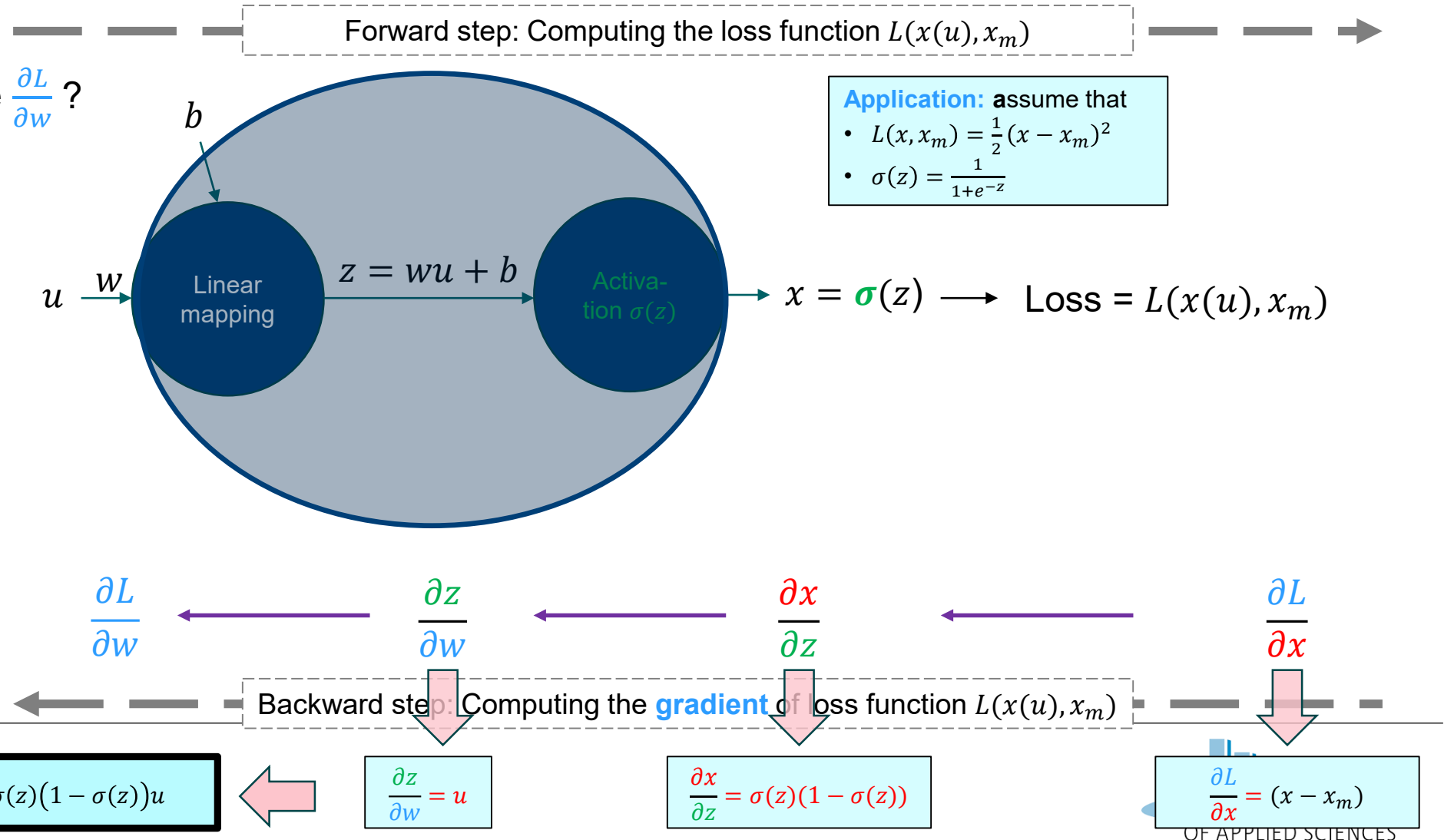
$\frac{\partial L}{\partial w} \longleftarrow \frac{\partial z}{\partial w} \longleftarrow \frac{\partial x}{\partial z} \longleftarrow \frac{\partial L}{\partial x}$

Backward step: Computing the **gradient** of loss function $L(x(u), x_m)$

Eric K
kaigor

$$\frac{\partial L}{\partial w} = 2(x - x_m)(1 - \sigma^2(z))u$$

$$\frac{\partial z}{\partial w} = u$$

$$\frac{\partial x}{\partial z} = 1 - \sigma^2(z)$$

$$\frac{\partial L}{\partial x} = 2(x - x_m)$$

OF APPLIED SCIENCES

## Neuron-related gradient computation



**Known activation function $\sigma(z)$**

Assume that
- $L(x, x_m) = (x - x_m)^2$
- $\sigma(z) = \tanh(z)$

$$\frac{\partial L}{\partial w} \longleftarrow \frac{\partial z}{\partial w} \longleftarrow \frac{\partial x}{\partial z} \longleftarrow \frac{\partial L}{\partial x}$$

Backward step: Computing the **gradient** of loss function $L(x(u), x_m)$

$$\boxed{\frac{\partial L}{\partial w} = 2(x - x_m)(1 - \sigma^2(z))u}$$

$$\frac{\partial z}{\partial w} = u$$

$$\frac{\partial x}{\partial z} = 1 - \sigma^2(z)$$

$$\frac{\partial L}{\partial x} = 2(x - x_m)$$

Assume that
- $L(x, x_m) = \frac{1}{2}(x - x_m)^2$
- $\sigma(z) = \frac{1}{1 + e^{-z}}$

$$\frac{\partial L}{\partial w} \longleftarrow \frac{\partial z}{\partial w} \longleftarrow \frac{\partial x}{\partial z} \longleftarrow \frac{\partial L}{\partial x}$$

Backward step: Computing the **gradient** of loss function $L(x(u), x_m)$

$$\boxed{\frac{\partial L}{\partial w} = (x - x_m)\sigma(z)(1 - \sigma(z))u}$$

$$\frac{\partial z}{\partial w} = u$$

$$\frac{\partial x}{\partial z} = \sigma(z)(1 - \sigma(z))$$

$$\frac{\partial L}{\partial x} = (x - x_m)$$

**Observe that:**

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z}\frac{\partial z}{\partial w} = p\frac{\partial z}{\partial w}$$

with $p = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z} = \frac{\partial L}{\partial z}$

Impact of the activation function on $p$ is kept hidden!

**Unknown activation function $\sigma(z)$**

## Application

1) Generate the **set $u$** of **input training data** made up of **10000 random numbers** between -1 and 1

2) Generate the set $y = x_m$ of label data with $y = x_m =$ **243.5∗$u$ − 162**

3) Implement a neuron with **input $u$**, **weight $w$, bias $b$, state $z$, identity output $x$** as activation (i.e., $x(z) = \sigma(z) = z$) in Matlab. The neuron maps the input signal $u$ to an ouput signal $x$

4) Write a loop (the maximum number of iterations is $10^4$) that uses the neuron to run the
   a) **Forward propagation** (from $u$ to the current neural network output $x$)
   b) **Compute the gradient of the loss function $L$** (what is a **useful** one?) by **using $u, x,$ and $x_m$**
   c) **Update the weights $(w, b)$** of the neuron

5) Did you retrieve ($w = 243.5, b = 162$)? What do you observe while increasing/decrasing the learning rate $r$?

6) Repeat from 2) with $x_m = 2u^2$+1. What did you observe?

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Learning with **several fully connected** neurons

## Neural network

$b_{in,2}$

$b_{1,2}$

$b_{out,2}$

$u_{in,2}$ $\xrightarrow{\quad w_{in,2} \quad}$ Non-Linearity (neuron input) $\xrightarrow{\quad w_{1,2} \quad}$ Non-Linearity (neuron 12) $\xrightarrow{\quad w_{out,2} \quad}$ Non-Linearity (neuron output) $\xrightarrow{\quad} x_{out,2}$

$x_{in,2} = \sigma_{in,2}(w_{in,2} u_{in,2} + b_{in,2})$

$x_{1,2} = \sigma_{1,2}(w_{1,2} x_{in,2} + b_{1,2})$

$x_{out,2} = \sigma_{out,2}(w_{out,2} x_{1,2} + b_{out,2})$

Input layer        Hidden layer        Output layer

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Fully connected neural network (NN)



$u_{in}$

Non-Linearity (neuron in)

Non-Linearity (neuron 11)

Non-Linearity (neuron 12)

Non-Linearity (neuron 13)

Non-Linearity (neuron 21)

Non-Linearity (neuron 22)

Non-Linearity (neuron 23)

Non-Linearity (neuron out)

$u_{out}$

Input Layer

1. Layer

2. Layer

Output Layer

Eric Kaigom
kaigom@fb2.fra-uas.de

MiMo (Multiple Input and Multiple Output) **are usual** (e.g., **2 joint positions as input** and **2 Cartesian positions as output…**).

OF APPLIED SCIENCES

## Neural network



$u_{in}$

$b_1$

Neuron input

Neuron 1

Neuron 2

Neuron 3

Neuron 1

Neuron 2

Neuron 3

Neuron output

$u_{out}$

$w_{11}^{in}$

$w_{12}^{in}$

$w_{13}^{in}$

$w_{11}^1$

$w_{12}^1$

$w_{13}^1$

$w_{21}^1$

$w_{22}^1$

$w_{23}^1$

$w_{31}^1 w_{32}^1$

$w_{32}^1$

$w_{1out}^2$

$w_{2out}^2$

$w_{3out}^2$

Input Layer

**1. Layer**

**2. Layer**

Output Layer

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Gradient via backpropagation

**Recall that:**

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z}\frac{\partial z}{\partial w}$$

$$\frac{\partial L}{\partial w^4} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial w^4}$$

$$\frac{\partial L}{\partial w^3} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial x^2}\frac{\partial x^2}{\partial z^2}\frac{\partial z^2}{\partial w^3}$$

$$\frac{\partial L}{\partial w^2} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial x^2}\frac{\partial x^2}{\partial z^2}\frac{\partial z^2}{\partial x^1}\frac{\partial x^1}{\partial z^1}\frac{\partial z^1}{\partial w^2}$$

$w^n$ means weights $w$ w.r.t $n$-th layer

**1** — Forward input propagation … $\quad \overrightarrow{w^1, z^1, x^2}\ \overrightarrow{w^2, z^2, x^3}\ \overrightarrow{w^3, z^3, x^4}$



$u_{in1}$, $u_{in2}$, $u_{in3}$

$w^{in}$ $\quad w^1 \quad w^2 \quad w^{out}$

Input Layer | 1. hidden Layer | 2. hidden Layer | 3. hidden Layer | Output Layer

**2** — **Backward gradient propagation …** $\quad \dfrac{\partial L}{\partial w^1} \leftarrow \dfrac{\partial L}{\partial w^2} \leftarrow \dfrac{\partial L}{\partial w^{out}}$

**3** — Weight update… $\qquad w_{k+1}^i = w_k^i - r\dfrac{\partial L}{\partial w^i}$

**4** — Feeding the neural network with training data…

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKF
UNIVERSI
OF APPLIED SCIENCES

## Gradient via backpropagation

**Recall that:**

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z}\frac{\partial z}{\partial w}$$
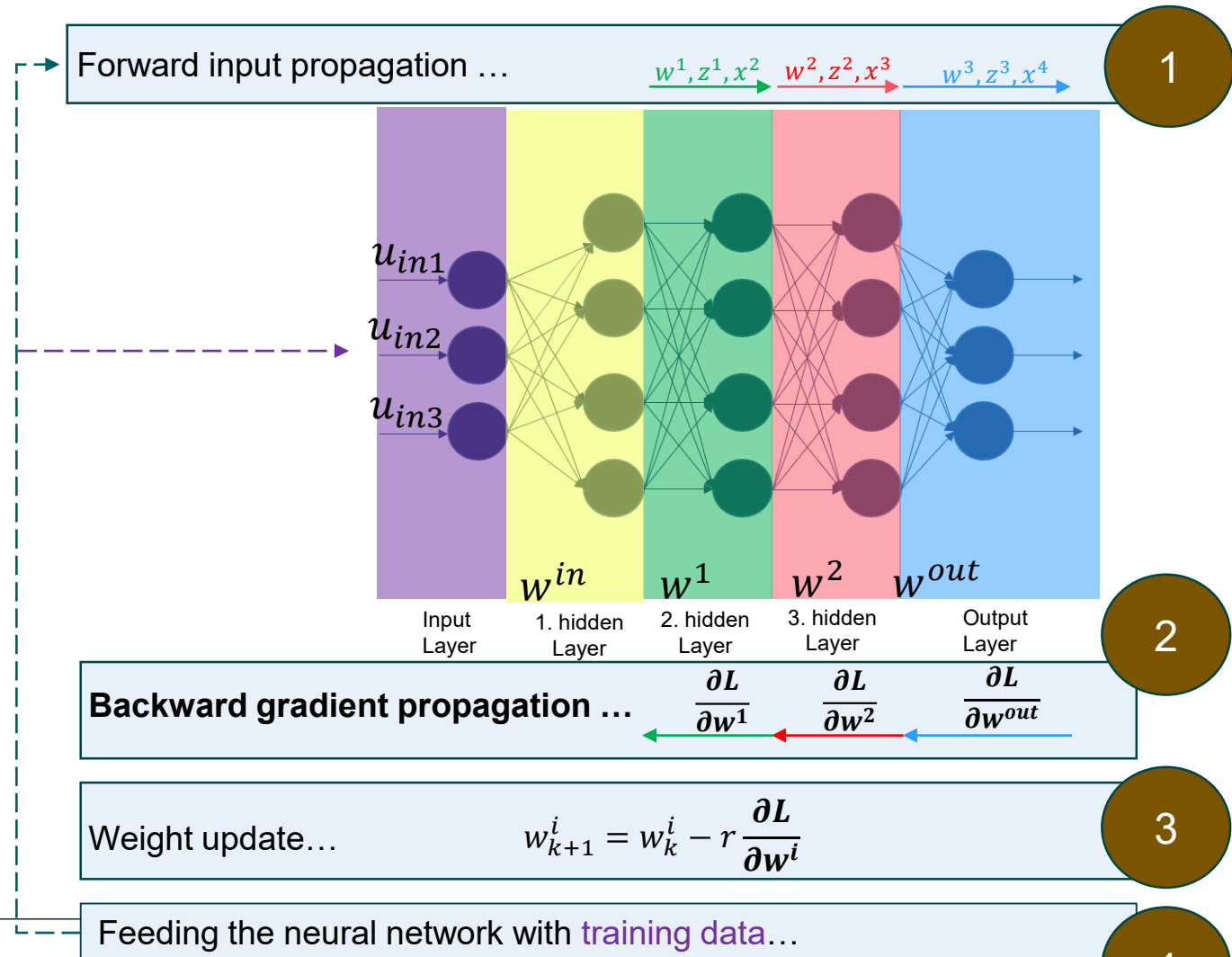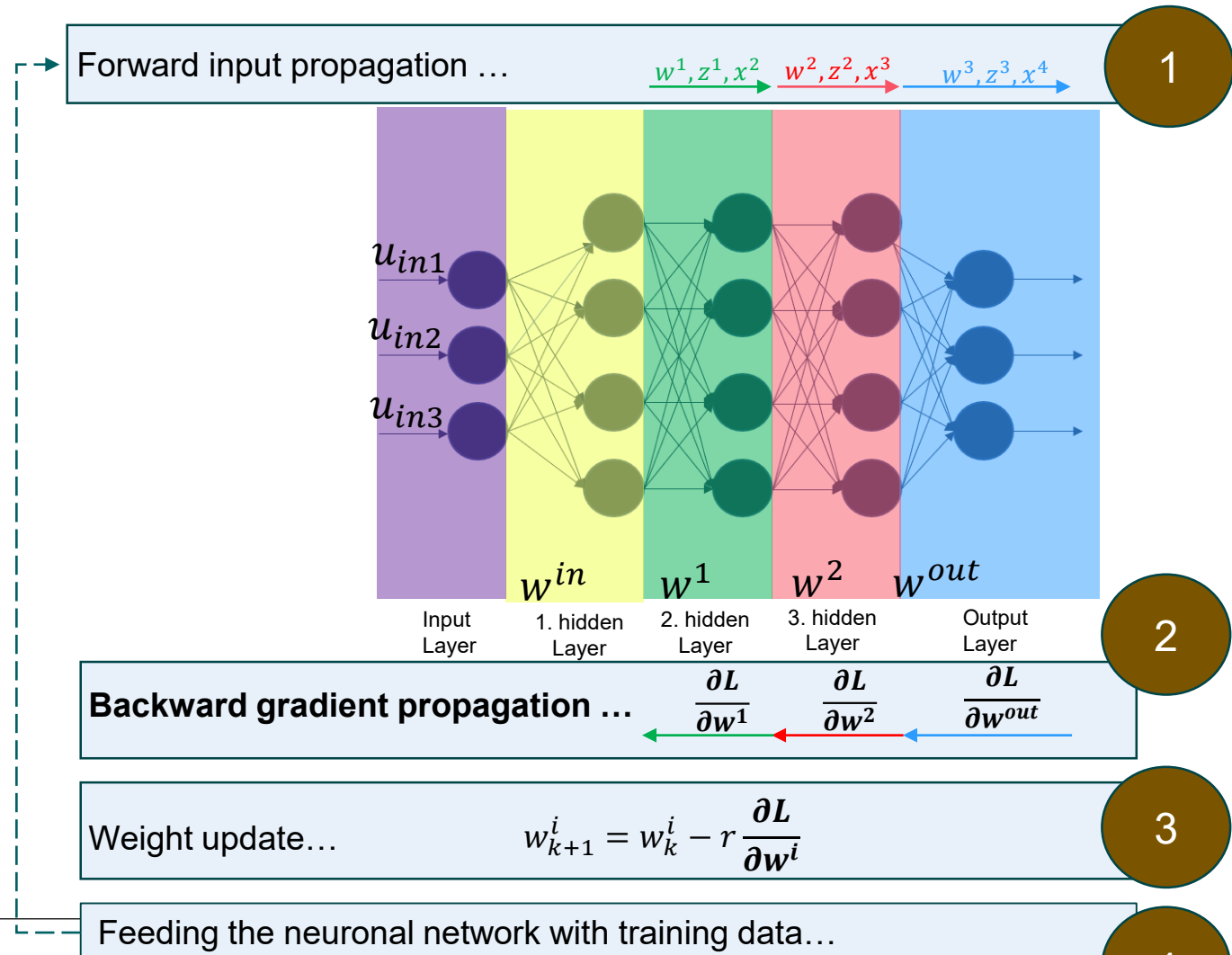
$$\frac{\partial L}{\partial w^4} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial w^4}$$

$$\frac{\partial L}{\partial w^3} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial x^2}\frac{\partial x^2}{\partial z^2}\frac{\partial z^2}{\partial w^3}$$

$$\frac{\partial L}{\partial w^2} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial x^2}\frac{\partial x^2}{\partial z^2}\frac{\partial z^2}{\partial x^1}\frac{\partial x^1}{\partial z^1}\frac{\partial z^1}{\partial w^2}$$

**Observations:**
Some **terms** appear **many times** in the gradient-related chain**!**

**1** Forward input propagation … $\quad \underrightarrow{w^1, z^1, x^2} \; \underrightarrow{w^2, z^2, x^3} \; \underrightarrow{w^3, z^3, x^4}$



$u_{in1}$
$u_{in2}$
$u_{in3}$

$w^{in} \qquad w^1 \qquad w^2 \qquad w^{out}$

Input Layer | 1. hidden Layer | 2. hidden Layer | 3. hidden Layer | Output Layer

**2** **Backward gradient propagation …** $\quad \underleftarrow{\frac{\partial L}{\partial w^1}} \; \underleftarrow{\frac{\partial L}{\partial w^2}} \; \underleftarrow{\frac{\partial L}{\partial w^{out}}}$

**3** Weight update… $\qquad w_{k+1}^i = w_k^i - r\frac{\partial L}{\partial w^i}$

**4** Feeding the neuronal network with training data…

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKF
UNIVERSI
OF APPLIED SCIENCES

## Gradient via backpropagation

**1**

Forward input propagation …   $w^1, z^1, x^2$   $w^2, z^2, x^3$   $w^3, z^3, x^4$

**Recall that:**

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial x}\frac{\partial x}{\partial z}\frac{\partial z}{\partial w}$$



$u_{in1}$   $u_{in2}$   $u_{in3}$

$w^{in}$   $w^1$   $w^2$   $w^{out}$

Input Layer | 1. hidden Layer | 2. hidden Layer | 3. hidden Layer | Output Layer

$$\frac{\partial L}{\partial w^4} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial w^4}$$

$$\frac{\partial L}{\partial w^3} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial x^2}\frac{\partial x^2}{\partial z^2}\frac{\partial z^2}{\partial w^3}$$

$$\frac{\partial L}{\partial w^2} = \frac{\partial L}{\partial x^3}\frac{\partial x^3}{\partial z^3}\frac{\partial z^3}{\partial x^2}\frac{\partial x^2}{\partial z^2}\frac{\partial z^2}{\partial x^1}\frac{\partial x^1}{\partial z^1}\frac{\partial z^1}{\partial w^2}$$
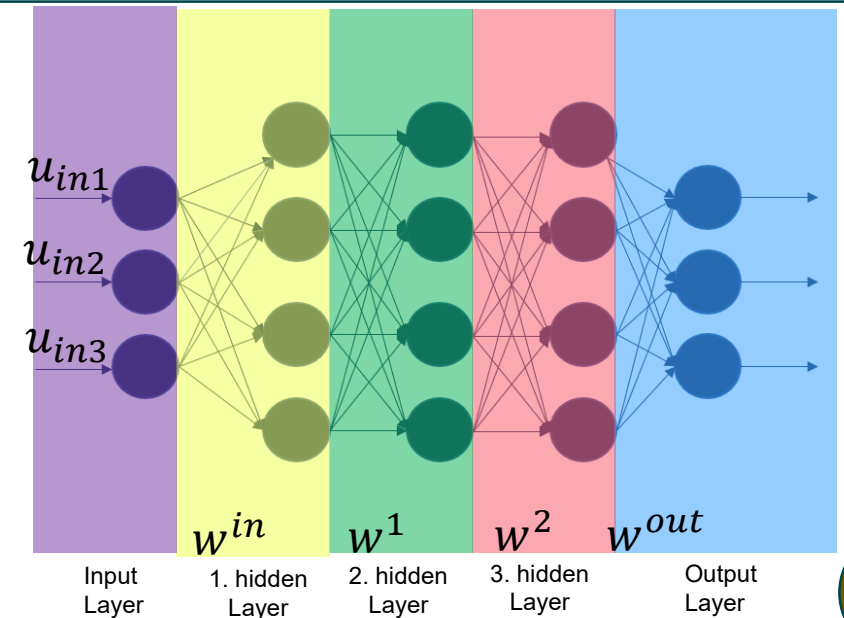
**2**

**Backward gradient propagation …**   $\frac{\partial L}{\partial w^1}$   $\frac{\partial L}{\partial w^2}$   $\frac{\partial L}{\partial w^{out}}$

**Observation**:

If some of the $\frac{\partial ()}{\partial ()}$ (e.g., $\frac{\partial x^2}{\partial z^2}$) in the chain are too small or simply vanish, the **gradient** is likely to **disappear**!

What does it implicate?

**3**

Weight update…   $w_{k+1}^i = w_k^i - r\frac{\partial L}{\partial w^i}$

Big challenge faced by **deep neural networks**!

Feeding the neuronal network with training data…

**4**

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKF
UNIVERSI
OF APPLIED SCIENCES

## Neural network

## Application

1) Implement a matlab function **[x1, x2] = forward(q1, q2, w)** that reflects a fully connected neural network made up of **two hidden layers** with **two neurons** for each of them. The function receives the following parameters
   a) Two real inputs **q1** and **q2**
   b) A vector **w** of weights for all layers
   and returns two real outputs **x1** and **x2**

2) Implement a loop that runs $10^6$ times, feeds **forward(q1, q2, w)** with randomly generated **q1, q2**, and **w**.

3) Set **q1** and **q2** as a vector of 1000 random values between -1 and 1. Furthermore, assume (ground truth) that
   • **X1 = q1 + q2**
   • **X2 = q1 – q2**
4) Run the loop while feeding the network (i.e., **forward (…)**) with **q1** and **q2** value pairwise and compute a meaningful loss function for each epoch using **X1** and **X2**.

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Machine Learning – Updating the weights of the neural network in practice



```
% Compute the contribution to the gradient of the loss for each weight
for i = 1:length(trainedWeights)
    % Perturb the current weight using espsilonValue
    trainedWeightsPerturbedRight = trainedWeights;
    trainedWeightsPerturbedRight(i) = trainedWeightsPerturbedRight(i) + espsilonValue;
```

Measurements potentially prone to noise!

**Data label**

(4)

**Result:** a trained digital model used for e.g. prediction purposes (e.g., FK(q) $\approx X_E$) in other applications (e.g., robotized indust. automation)

```
% Approximate the gradient contribution using a difference
grad(i) = (costFunctionPertubedRight - costFunctionPertubedLeft)...
/ (2 * espsilonValue);
```

Model capture using maths, statistics, and data, etc

(3)

**Model output during training**

Loss function $f$

- $\epsilon$: small and positive number
- $e_1$: [1 0 0 0 ….]
- $e_2$: [0 1 0 0 ….]

(1)

**Pertinent data** (e.g., $q$) as **input examples**

(5)

Computation of the gradient $\nabla f$ of the loss function $f$ w.r.t the weights of the neural network

Model training (i.e., optimization of model parameters to capture the mapping between input-output examples and generalize)

$$\frac{\partial f(w)}{\partial w_1} = \frac{f(x + \epsilon e_1) - f(x - \epsilon e_1)}{2\epsilon}$$

$$\frac{\partial f(w)}{\partial w_2} = \frac{f(x + \epsilon e_2) - f(x - \epsilon e_2)}{2\epsilon}$$

(2)

$$w_{k+1} = w_k - r\nabla f w_k$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \frac{\partial f}{\partial w_3} \\ \vdots \\ \frac{\partial f}{\partial w_n} \end{bmatrix}$$

$$x = f(w_1, \ldots, w_n)_u \rightarrow \nabla f = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_n} \end{bmatrix}$$

**Physical**

**Digital**

**Weight update**

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT UNIVERSITY OF APPLIED SCIENCES

## Machine Learning – Adapting the learning rate via moment estimation

- $m_{k+1} = \beta_1 m_k + (1 - \beta_1)\nabla f$

- $v_{k+1} = \beta_2 v_k + (1 - \beta_2)(\nabla f)^2$

$\beta_1$ and $\beta_2$ induce an exponential decay (from $\beta_1 = 0.9, \beta_2 = 0.999$)

- $\widehat{m}_{k+1} = \dfrac{m_{k+1}}{1 - \beta_1^{k+1}}$

adaptation

- $\widehat{v}_{k+1} = \dfrac{v_{k+1}}{1 - \beta_2^{k+1}}$

- $w_{k+1} = w_k - \alpha \dfrac{1}{\sqrt{\widehat{v}_k} + \epsilon} \widehat{m}_k$

- $\alpha = 0.001$
- $\epsilon = 10^{-8}$

```
% Compute the contribution to the gradient of the loss for each weight
for i = 1:length(trainedWeights)
    % Perturb the current weight using espsilonValue
    trainedWeightsPerturbedRight = trainedWeights;
    trainedWeightsPerturbedRight(i) = trainedWeightsPerturbedRight(i) + espsilonValue;
```

**Data label**

④

```
% Approximate the gradient contribution using a difference
grad(i) = (costFunctionPertubedRight - costFunctionPertubedLeft)...
/ (2 * espsilonValue);
```

③

**Model output during training**

⑤

- $\epsilon$: small and positive number
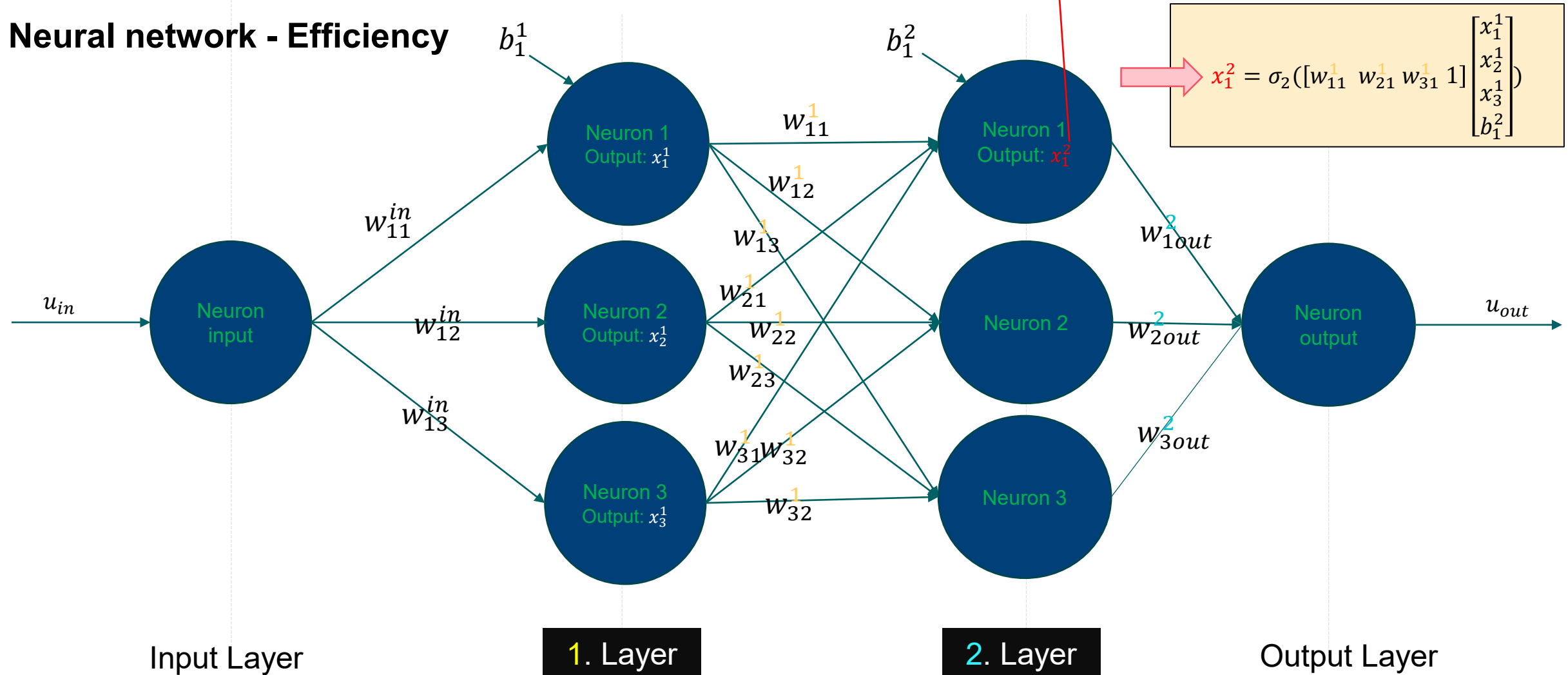- $e_1$: [1 0 0 0 0 ….]
- $e_2$: [0 1 0 0 0 ….]

$w_{k+1} = w_k - r\nabla f(w)_{w_k}$

**Instead of …**
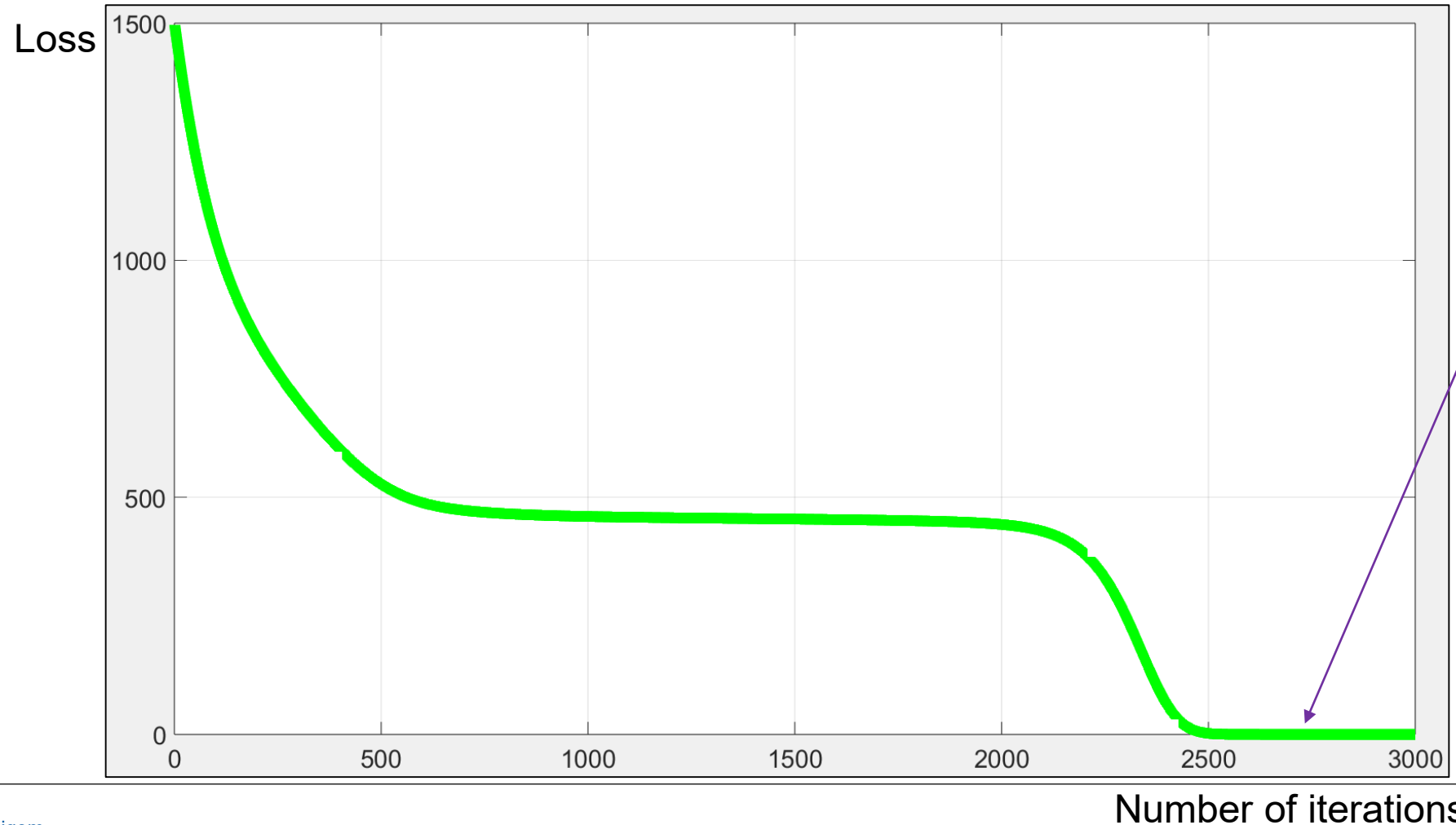
FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

## Optimization of the prediction accuracy (under adaptive moment-based learning rate) - Result



Loss

Number of iterations

Eric Kaigom
kaigom@fb2.fra-uas.de

**Project (presentation: 29.01.25, 8:15 h, groups of min. 2 /max. 3 students, Duration max. 20 Min)**

- Improve our model (see the file withStudentsNew.m) to learn

  - **X1 = q1 + q2 - 1**
  - **X2 = q1 – q2 + 1**

  **and**

  - **X1 = q1 + q2.*q2 - 1**
  - **X2 = q1 – q2 + 1**

- When and why is an enhancement necessary?
- How did you enhance the model?
- What did you observe? Explain your observation with your own words

Eric Kaigom
kaigom@fb2.fra-uas.de

FRANKFURT
UNIVERSITY
OF APPLIED SCIENCES

# Thank you
# for your attention!