# Lecture 1: Introduction & Data

## KI-Workshop
## (HFT Stuttgart, 8-9 Nov 2023)

Michael Mommert
University of St. Gallen (soon-to-be HFT Stuttgart)

What this course is about...

Who am I?

Course modalities
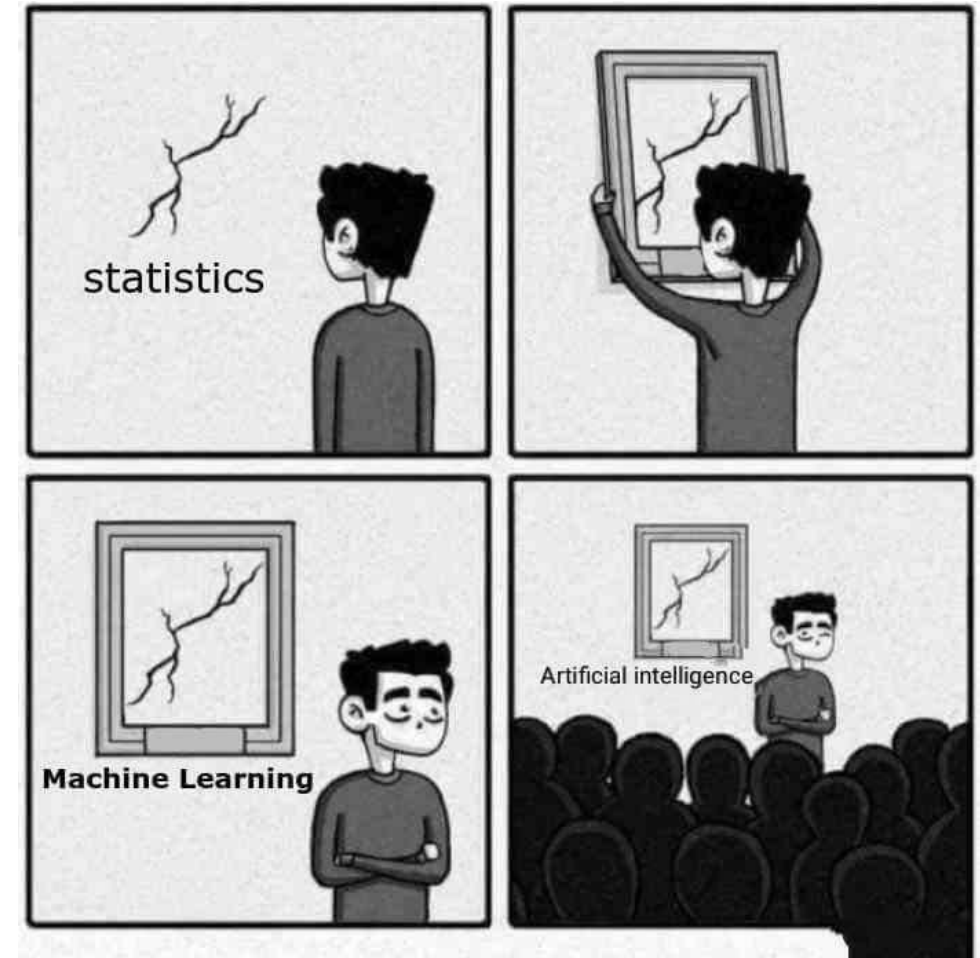
Course syllabus

Types of data

Features and feature engineering
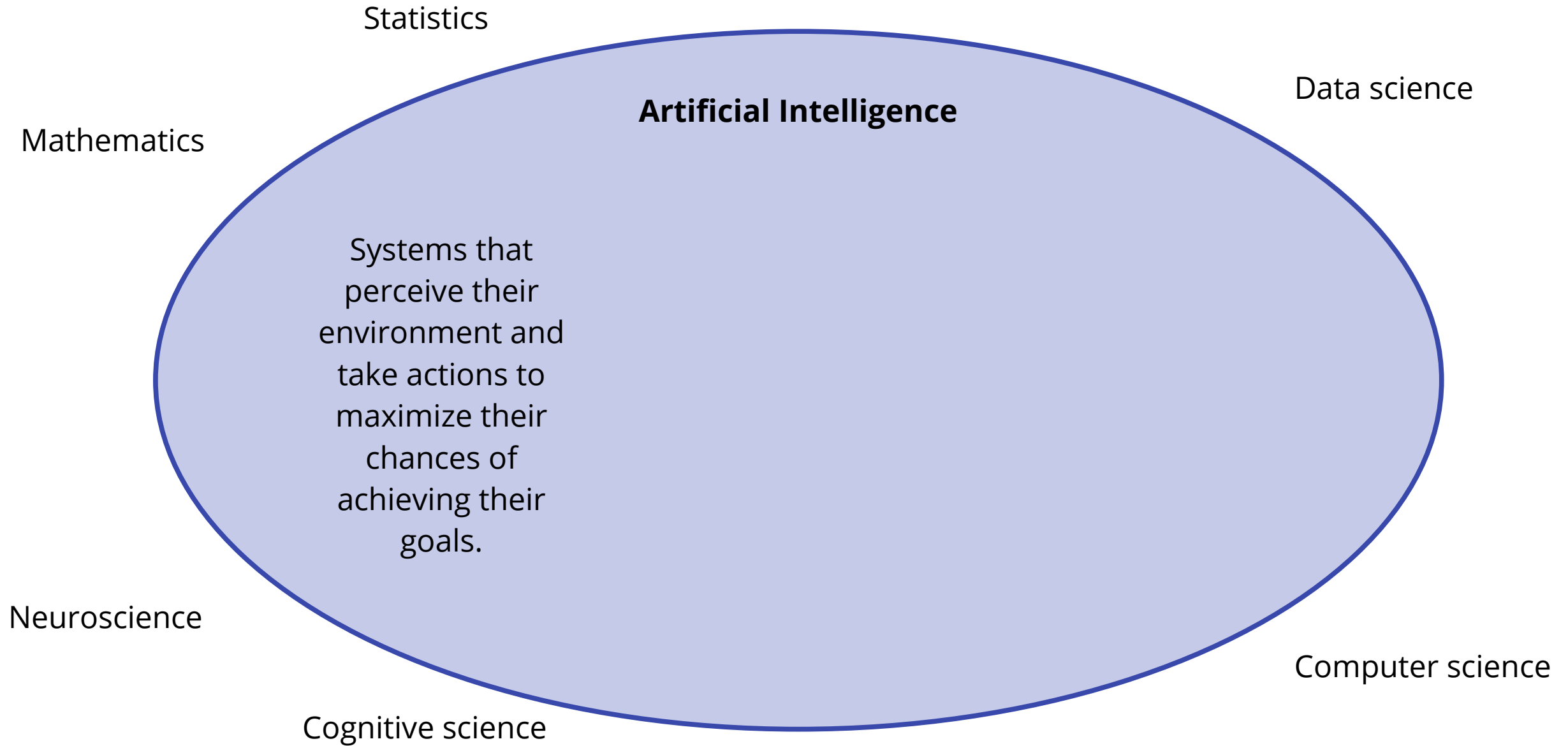
Data scaling

# What this course is about...
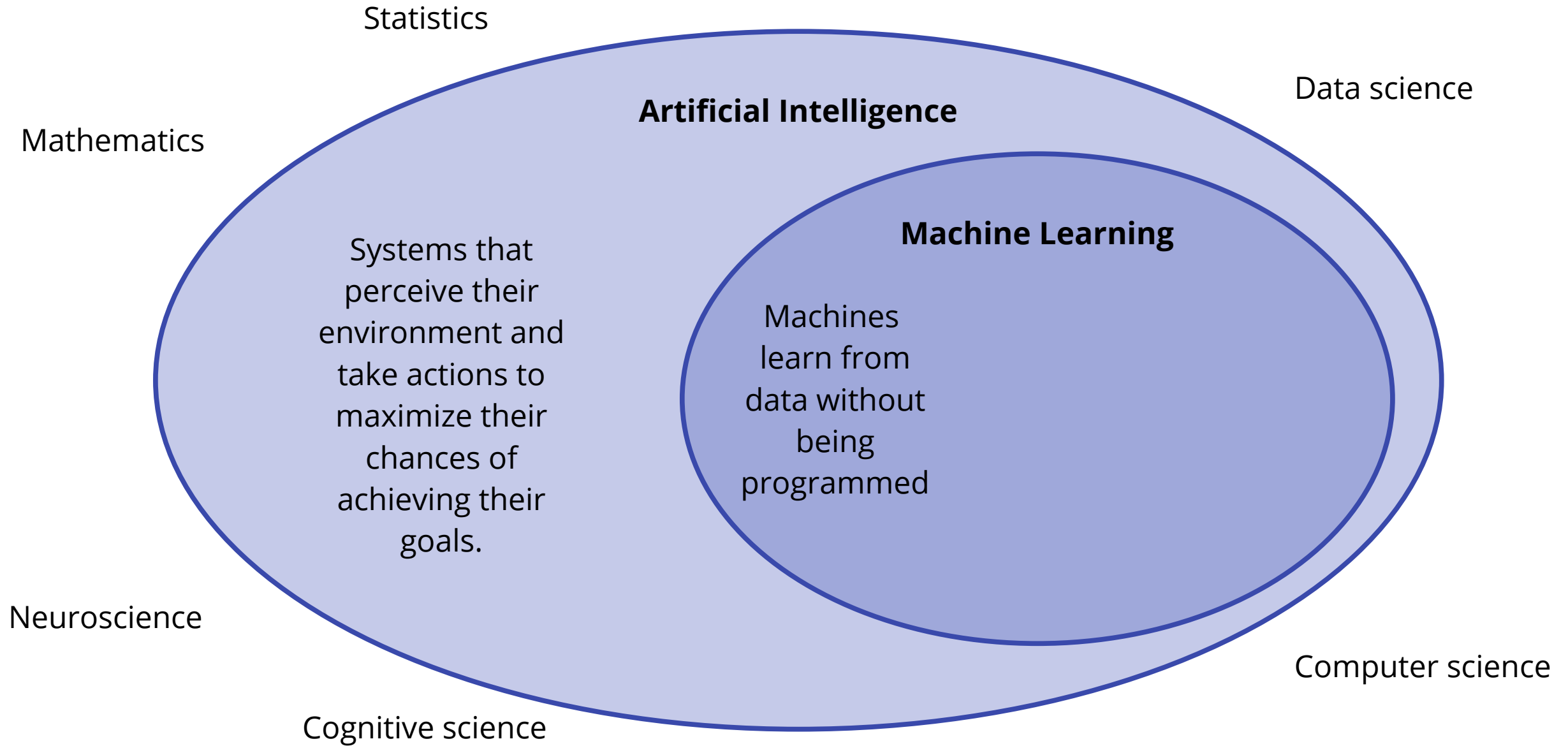


sandserif

# Mapping terminology

**Artificial Intelligence**

Systems that perceive their environment and take actions to maximize their chances of achieving their goals.
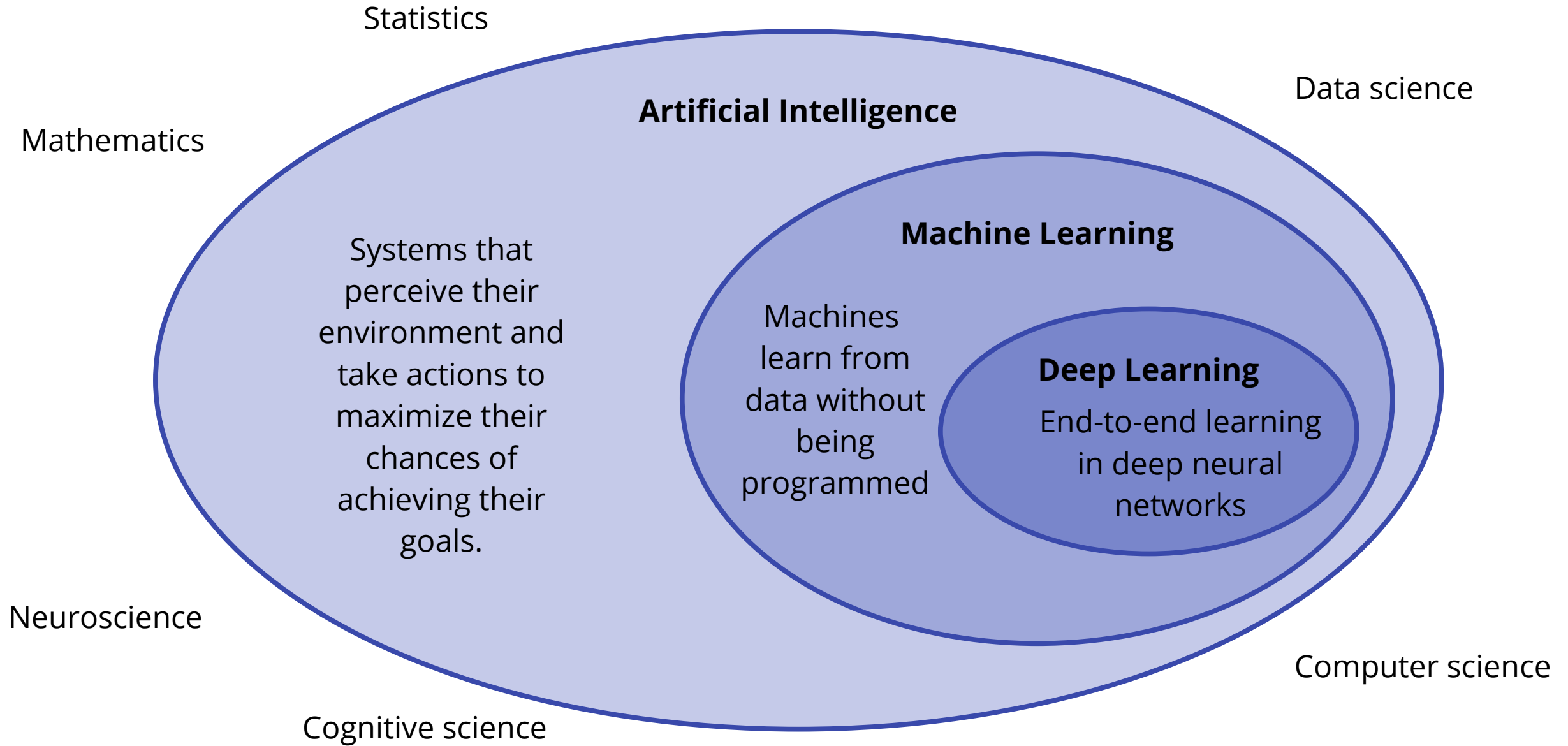
Statistics

Data science

Mathematics

**Artificial Intelligence**

Systems that perceive their environment and take actions to maximize their chances of achieving their goals.

Neuroscience

Computer science

Cognitive science

Statistics

Data science

Mathematics

**Artificial Intelligence**

**Machine Learning**

Systems that perceive their environment and take actions to maximize their chances of achieving their goals.

Machines learn from data without being programmed

Neuroscience

Computer science

Cognitive science

Statistics

Data science

Mathematics

**Artificial Intelligence**

**Machine Learning**

Systems that perceive their environment and take actions to maximize their chances of achieving their goals.
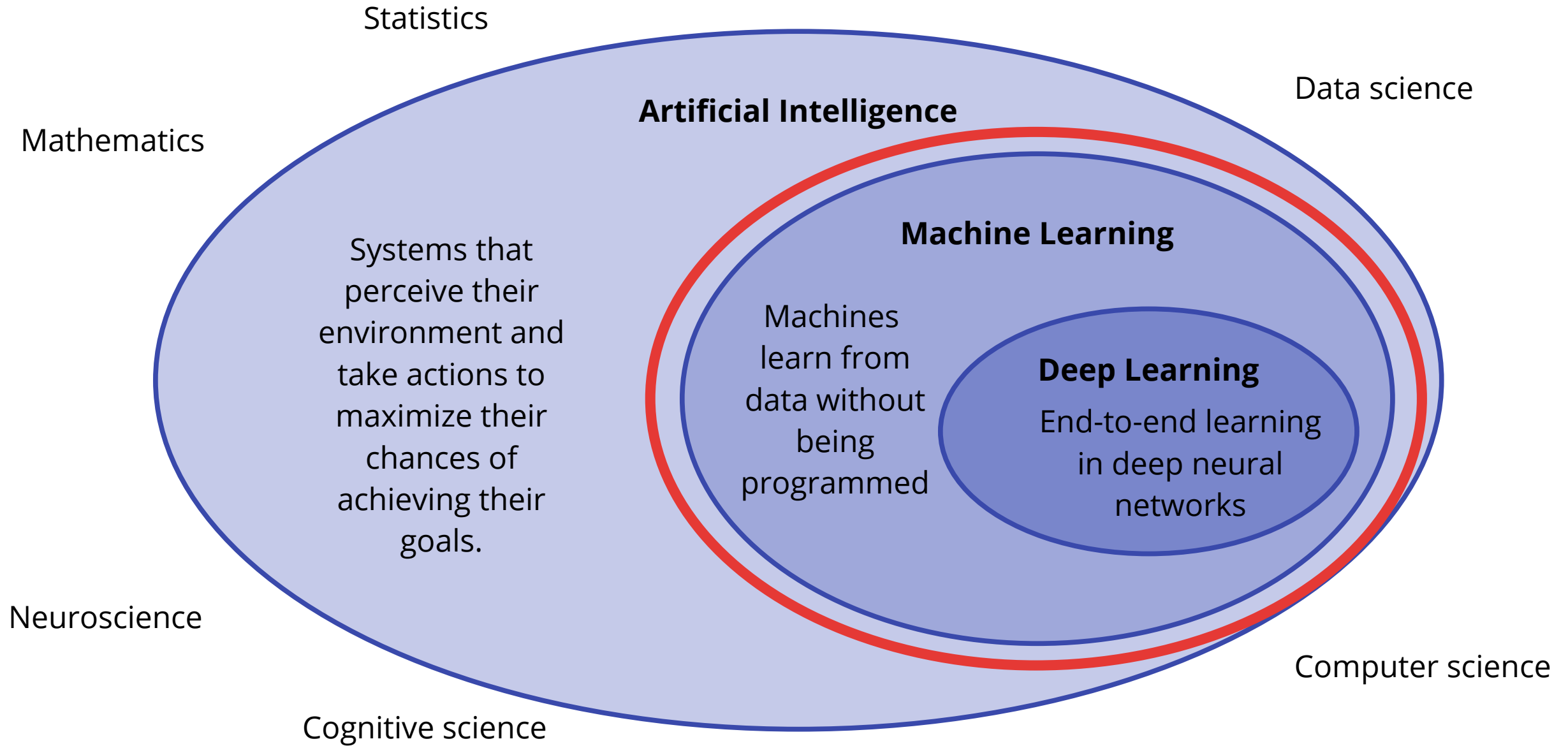
Machines learn from data without being programmed

**Deep Learning**
End-to-end learning in deep neural networks

Neuroscience

Computer science

Cognitive science

# What is Machine Learning (ML)?

> *"The field of study that gives computers the ability to learn without being explicitly programmed."*
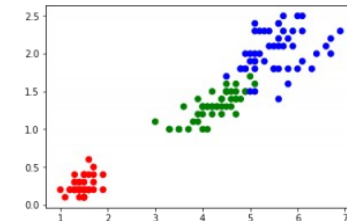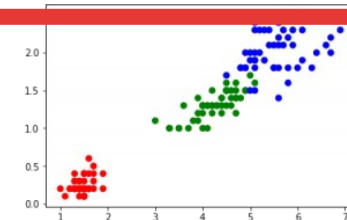> - Arthur Samuel (1959)

Different approaches:



Iris Versicolor

- **Supervised learning**
  Find a function that relates input data to output data by learning a specific task.



- **Unsupervised learning**
  Find structure within a data set.



- **Reinforcement learning**
  Learn a task in a dynamic and responsive environment.

# What is Machine Learning (ML)?

> *"The field of study that gives computers the ability to learn without being explicitly programmed."*
> - Arthur Samuel (1959)

Different approaches:


Iris Versicolor

- **Supervised learning**
  Find a function that relates input data to output data by learning a specific task.

- **Unsupervised learning**
  Find structure within a data set.



- **Reinforcement learning**
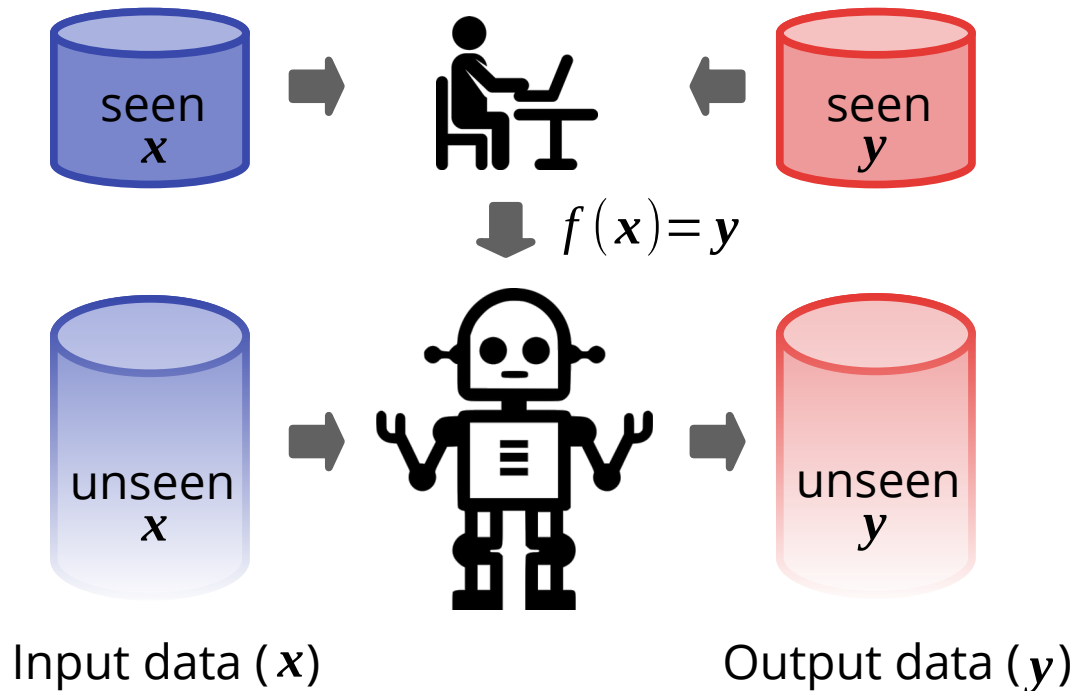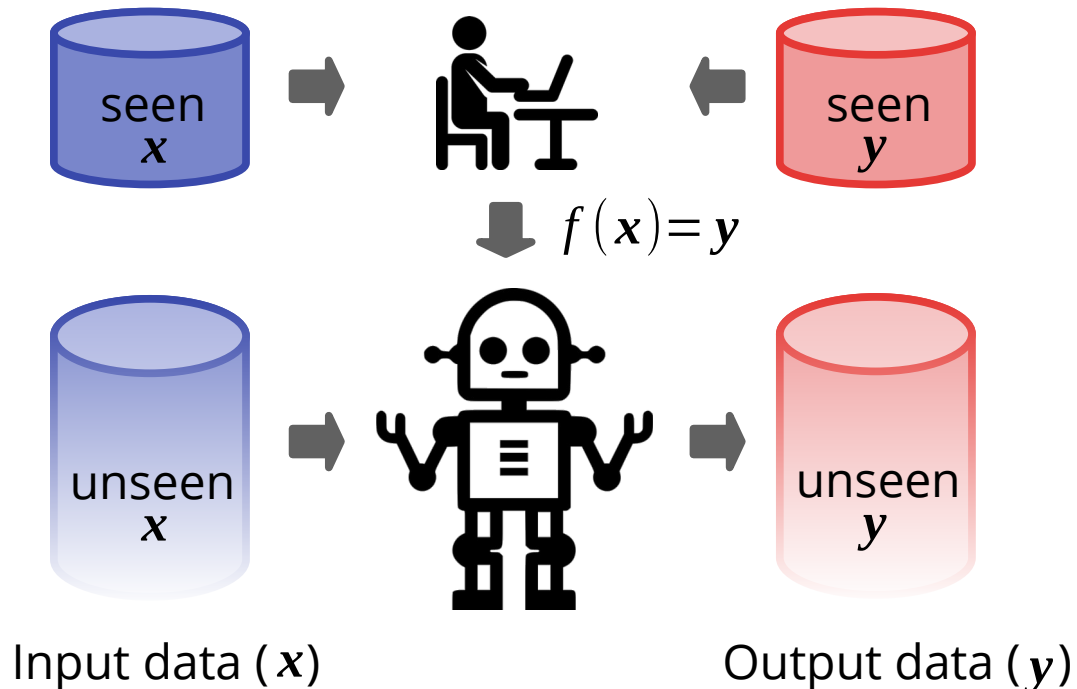  Learn a task in a dynamic and responsive environment.

# Supervised ML

**General goal for supervised problems**:
Find a function ("task") that relates input data ($x$) to output data ($y$) such that: $f(x) = y$

**General goal for supervised problems**:
Find a function ("task") that relates input data ($x$) to output data ($y$) such that: $f(x) = y$

**Traditional (Rule-based) Approach:**



Input data ($x$)          Output data ($y$)

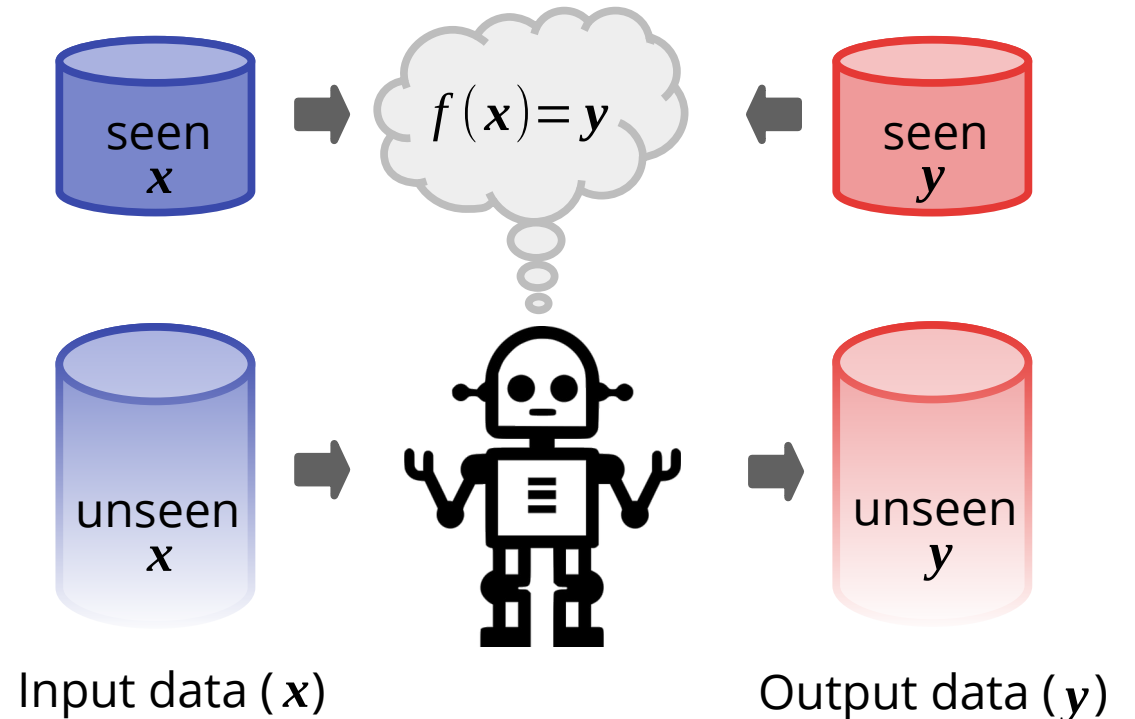**General goal for supervised problems**:

Find a function ("task") that relates input data ($x$) to output data ($y$) such that: $f(x) = y$

**Traditional (Rule-based) Approach:**



$f(x) = y$

Input data ($x$)          Output data ($y$)

**Machine-Learning Approach:**



$f(x) = y$

Input data ($x$)          Output data ($y$)

# Who am I?

# About myself

Physics

2009
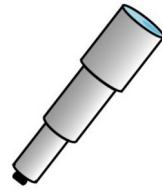
Physics

Dr. rer. nat.
(Earth Sciences)

2009                    2013

ISAS/JAXA

Gerald Rhemann

Deutsches Zentrum
DLR für Luft- und Raumfahrt
German Aerospace Center

NAU
NORTHERN
ARIZONA
UNIVERSITY

UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

Freie Universität Berlin

LOWELL
OBSERVATORY
125 YEARS | 1894 - 2019

Physics

Dr. rer. nat.
(Earth Sciences)

2009

2013

ISAS/JAXA

Gerald Rhemann

UNIVERSITÄT HEIDELBERG
ZUKUNFT SEIT 1386

Deutsches Zentrum DLR für Luft- und Raumfahrt
German Aerospace Center

Freie Universität Berlin

NAU NORTHERN ARIZONA UNIVERSITY

LOWELL OBSERVATORY
125 YEARS | 1894 - 2019

| Physics | Dr. rer. nat. (Earth Sciences) | Postdoc @HSG-AIML |
|---|---|---|

2009          2013          2020

**About myself**

Physics

Dr. rer. nat.
(Earth Sciences)

Postdoc
@HSG-AIML

2009

2013

2020

# About myself



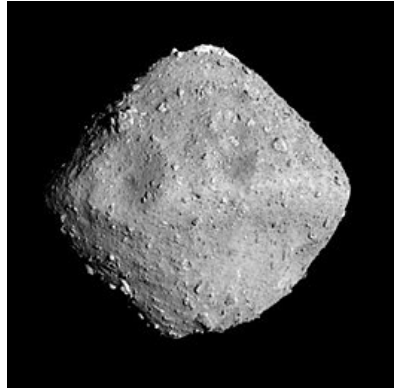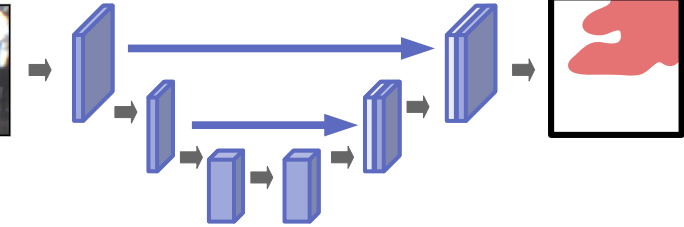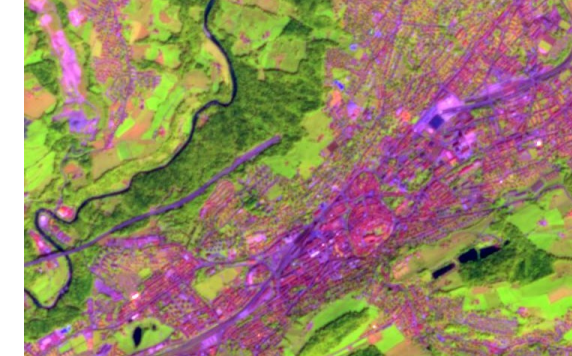| Physics | Dr. rer. nat. (Earth Sciences) | Postdoc @HSG-AIML | Asst. Prof. Computer Vision |
|---------|-------------------------------|-------------------|----------------------------|
| 2009 | 2013 | 2020 | 2022 |

# About myself



ISAS/JAXA

Gerald Rhemann

| Physics | Dr. rer. nat. (Earth Sciences) | Postdoc @HSG-AIML | Asst. Prof. Computer Vision | Prof. AI in Remote Sensing |
|---|---|---|---|---|

| 2009 | 2013 | 2020 | 2022 | 2024 |

Commercial Vehicle Traffic Monitoring
(*Blattner et al. 2021*)

# What I work on...



Commercial Vehicle Traffic Monitoring
(*Blattner et al. 2021*)

Characterization of Plumes and Estimation of
Power Generation from Remote Sensing Data
(*Mommert et al. 2020, Hanna et al. 2023*)



True Color

Segmentation Map

IoU=0.71  IoU=0.13  IoU=0.71

R: ground-truth, G: prediction

Commercial Vehicle Traffic Monitoring
(*Blattner et al. 2021*)

Characterization of Plumes and Estimation of Power Generation from Remote Sensing Data
(*Mommert et al. 2020, Hanna et al. 2023*)

True Color

Segmentation Map

IoU=0.71    IoU=0.13    IoU=0.71

R: ground-truth, G: prediction

Fossil Hard Coal

Hydro Water Reservoir

Power Plant
Classification from
Remote Imaging with
Deep Learning
(*Mommert et al. 2021*)

# What I work on...

Commercial Vehicle Traffic Monitoring
(*Blattner et al. 2021*)

Characterization of Plumes and Estimation of Power Generation from Remote Sensing Data
(*Mommert et al. 2020, Hanna et al. 2023*)



True Color

Segmentation Map

IoU=0.71   IoU=0.13   IoU=0.71

R: ground-truth, G: prediction

Power Plant Classification from Remote Imaging with Deep Learning
(*Mommert et al. 2021*)

Fossil Hard Coal

Hydro Water Reservoir

Contrastive Self-supervised data fusion for Satellite Imagery
(Scheibenreif et al. 2022)

Sentinel-1               Sentinel-2

Location 1   S1 Encoder   attract   S2 Encoder   Location 1

Location 2   S1 Encoder   contrast   S2 Encoder   Location 2

Location 3   S1 Encoder   attract   S2 Encoder   Location 3

Latent space

# Course modalities

- **Goal** of this course:
  *To understand and be able to implement and utilize supervised traditional Machine Learning and Deep Learning models.*

- **Setup**: Combination of lectures and voluntary hands-on lab courses

- **Lecture mode**: This course is supposed to be bi-directional: let me know if anything is unclear, ask questions anytime!

- We will use **Google Colab** for running our Lab Notebooks (they offer free GPUs!). If you don't have a Google account, please let me know as soon as possible!

# Literature resources

- Stuart Russell, Peter Norvig: **Artificial Intelligence: A Modern Approach** (2020 and earlier versions, MIT Press)
*Part V ("Learning") is especially relevant to this course and provides good introductions*

ebook@HSG

- Andreas Müller & Sarah Guido: **Introduction to Machine Learning with Python** (2017, O'Reilly)
*Easy-to-understand introduction to Python for ML, uses scikit-learn*

- Ian Goodfellow, Yoshua Bengio, Aaron Courville: **Deep Learning** (2016, MIT Press)
*All you need to know about Deep Learning*

free online

# Course syllabus

# Content

| Slot | Wednesday | Thursday |
|------|-----------|----------|
| 09:00 - 10:30 | **Intro & Data** | **Neural Networks** |
| 10:30 – 10:45 | Break | Break |
| 10:45 – 12:15 | **Supervised ML: Concepts** | **Convolutional Neural Networks & Computer Vision** |
| 12:15 – 13:45 | Lunch break | Lunch break |
| 13:45 – 15:15 | **Supervised ML: Methods** | **Lab: Neural Networks** |
| 15:15 – 15:30 | Break | Break |
| 15:30 – 17:00 | **Lab: Supervised ML** | **Advanced Deep Learning** |

**Data**

- Data storage used to be a bottleneck – not anymore!

### Historical Cost of Computer Memory and Storage



Legend:
- Flip-Flops
- Core
- ICs on boards
- SIMMs
- DIMMs
- Big Drives
- Floppy Drives
- Small Drives
- Flash Memory
- SSD

RAM

SSD

HDDs

John C. McCallum

- Data storage used to be a bottleneck – not anymore!

# Data storage

- Data storage used to be a bottleneck – not anymore!

- Vast amounts of data can now be stored easily

# Data storage

- Data storage used to be a bottleneck – not anymore!

- Vast amounts of data can now be stored easily

Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2025 (in zettabytes)

$$zeta = 10^{21} = 1000000000000000000000$$

Data volume in zettabytes

| Year | Value |
|------|-------|
| 2010 | 2 |
| 2011 | 5 |
| 2012 | 6.5 |
| 2013 | 9 |
| 2014 | 12.5 |
| 2015 | 15.5 |
| 2016 | 18 |
| 2017 | 26 |
| 2018* | 33 |
| 2019* | 41 |
| 2020* | 64.2 |
| 2021* | 79 |
| 2022* | 97 |
| 2023* | 120 |
| 2024* | 147 |
| 2025* | 181 |

Sources
IDC; Seagate; Statista estimates
© Statista 2021

Additional Information:
Worldwide; 2010 to 2020

Statista

# Data storage

- Data storage used to be a bottleneck – not anymore!

- Vast amounts of data can now be stored easily

- Is all this data technically accessible for analysis?
  (of course not, since most of it is privately owned, but…)

**Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2025 (in zettabytes)**

zeta = $10^{21}$ = 1000000000000000000000

Data volume in zetabytes

| Year | Value |
|------|-------|
| 2010 | 2 |
| 2011 | 5 |
| 2012 | 6.5 |
| 2013 | 9 |
| 2014 | 12.5 |
| 2015 | 15.5 |
| 2016 | 18 |
| 2017 | 26 |
| 2018* | 33 |
| 2019* | 41 |
| 2020* | 64.2 |
| 2021* | 79 |
| 2022* | 97 |
| 2023* | 120 |
| 2024* | 147 |
| 2025* | 181 |

Sources
IDC; Seagate; Statista estimates
© Statista 2021

Additional Information:
Worldwide; 2010 to 2020

Statista

# Structured vs unstructured data

University of St.Gallen

**Structured data**

Preprocessed and formatted data that
is easily queryable.

## Structured data

Preprocessed and formatted data that
is easily queryable.

Quantitative data

# Structured vs unstructured data

## Structured data

Preprocessed and formatted data that is easily queryable.

Quantitative data



## Unstructured data

Unprocessed and unformatted data is not easily queryable.

# Structured vs unstructured data

## Structured data

Preprocessed and formatted data that
is easily queryable.

Quantitative data

## Unstructured data

Unprocessed and unformatted data is
not easily queryable.

Qualitative
data

Image data

Video data

Textual data

Data stream

Audio data

# Structured vs unstructured data

## Structured data

Preprocessed and formatted data that
is easily queryable.

Quantitative data

## Unstructured data

Unprocessed and unformatted data is
not easily queryable.

Qualitative
data

Image data

Video data

Data complexity

Textual data
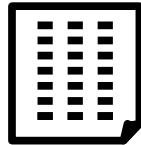
Data stream

Audio data

# Structured vs unstructured data

## Structured data

Preprocessed and formatted data that
is easily queryable.

Quantitative data



Most data analysis techniques require data to be
available in a structured form for easier processing.

Structured data can always be represented in a
database **schema** (e.g., a table in 2 dimensions).

## Unstructured data

Unprocessed and unformatted data is
not easily queryable.

Qualitative data



Image data



Video data



Data complexity

Textual data



Data stream

Audio data

# Structured vs unstructured data

## Structured data

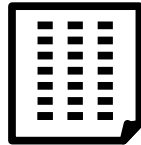Preprocessed and formatted data that is easily queryable.

Quantitative data

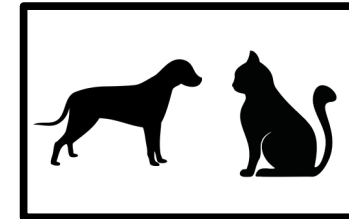Most data analysis techniques require data to be available in a structured form for easier processing.

Structured data can always be represented in a database **schema** (e.g., a table in 2 dimensions).

## Unstructured data

Unprocessed and unformatted data is not easily queryable.

Qualitative data

Image data

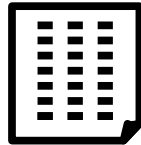Video data

Data complexity

Textual data

Data stream

Audio data

# Quantitative and qualitative data

**Quantitative data**
(can be measured; distances can be defined)

**Qualitative (categorical) data**
(cannot be measured; distances not defined)

# Quantitative and qualitative data

| Quantitative data (can be measured; distances can be defined) | Qualitative (categorical) data (cannot be measured; distances not defined) |
| --- | --- |

**Continuous data**

Real-valued numbers; potentially within a given range

*Examples*:
- Temperatures
- A person's height
- Prices

# Quantitative and qualitative data

| Quantitative data (can be measured; distances can be defined) | | Qualitative (categorical) data (cannot be measured; distances not defined) |
|---|---|---|
| **Continuous data**<br><br>Real-valued numbers; potentially within a given range<br><br><br>*Examples*:<br>• Temperatures<br>• A person's height<br>• Prices | **Discrete data**<br><br>Discrete numbers; whole numbers or real numbers, potentially within a given range<br><br>*Examples*:<br>• Number of people in a room<br>• Inventory counts | |

# Quantitative and qualitative data

| Quantitative data (can be measured; distances can be defined) | | Qualitative (categorical) data (cannot be measured; distances not defined) |
|---|---|---|
| **Continuous data** | **Discrete data** | **Nominal data** |
| Real-valued numbers; potentially within a given range | Discrete numbers; whole numbers or real numbers, potentially within a given range | Labels for different categories without ordering |
| *Examples*: <br> • Temperatures <br> • A person's height <br> • Prices | *Examples*: <br> • Number of people in a room <br> • Inventory counts | *Examples*: <br> • Color of hair <br> • Names of persons <br> • Types of fruit |

# Quantitative and qualitative data

| Quantitative data<br>(can be measured; distances can be defined) | | Qualitative (categorical) data<br>(cannot be measured; distances not defined) | |
|---|---|---|---|
| **Continuous data**<br><br>Real-valued numbers; potentially within a given range<br><br><br><br>*Examples*:<br>• Temperatures<br>• A person's height<br>• Prices | **Discrete data**<br><br>Discrete numbers; whole numbers or real numbers, potentially within a given range<br><br>*Examples*:<br>• Number of people in a room<br>• Inventory counts | **Nominal data**<br><br>Labels for different categories without ordering<br><br><br><br>*Examples*:<br>• Color of hair<br>• Names of persons<br>• Types of fruit | **Ordinal data**<br><br>Labels for different categories following an inherent ranking scheme.<br><br><br>*Examples*:<br>• Rank in a competition<br>• Grades<br>• Day of the week |

# Turning unstructured data into structured data

**Structured data**

Preprocessed and formatted data that
is easily queryable.

Quantitative

**Unstructured data**

Unprocessed and unformatted data is
not easily queryable.

Qualitative
data

Image data

Video data

Textual data

Data stream

Audio data

# Turning unstructured data into structured data

**Structured data**

**Unstructured data**

Preprocessed and formatted data that is easily queryable.

Unprocessed and unformatted data is not easily queryable.

Quantitative

Qualitative data

Image data

Video data

Before ML methods can be applied to unstructured data, we have to process those and extract useful features from them.

This process is called **feature engineering**.

Textual data

Data stream

Audio data

**Features and Feature Engineering**

# What are features?

Features are quantitative and independent variables based on which our ML models learn.

Features are quantitative and independent variables based on which our ML models learn.



Raw data
(qualitative)

Features are quantitative and independent variables based on which our ML models learn.



$x$

Raw data
(qualitative)

Feature engineering

Features
(quantitative)

Features are quantitative and independent variables based on which our ML models learn.

Features are quantitative and independent variables based on which our ML models learn.

# Feature engineering

Extract or create features that may provide a ML model with rich information on its task based on **domain knowledge**. Feature engineering can be applied to raw data, resulting in quantitative data that can be directly fed into the ML model (features).

Extract or create features that may provide a ML model with rich information on its task based on **domain knowledge**. Feature engineering can be applied to raw data, resulting in quantitative data that can be directly fed into the ML model (features).

# Feature engineering – quantitative data

Create meaningful features through mathematical transformations.

*Examples*:

# Feature engineering – quantitative data

Create meaningful features through mathematical transformations.

*Examples*:

**Arithmetic**

*Situation*: You have two variables, $x_1$ and $x_2$ , but you are more interested in their difference, $\delta$ .

*Transformation*:

$$\delta = x_1 - x_2$$

# Feature engineering – quantitative data

Create meaningful features through mathematical transformations.

*Examples*:

| Arithmetic | Aggregation of Features |
|---|---|
| *Situation*: You have two variables, $x_1$ and $x_2$, but you are more interested in their difference, $\delta$. | *Situation*: You have results from different business units, $x_i$, but your ML model should not consider the results separately, but as an aggregated overall result, $x$. |
| *Transformation*: $$\delta = x_1 - x_2$$ | *Transformation*: $$x = \sum_i x_i$$ |

# Feature engineering – quantitative data

Create meaningful features through mathematical transformations.

*Examples*:

| **Arithmetic** | **Aggregation of Features** | **Geometric Transformations** |
|---|---|---|
| *Situation*: You have two variables, $x_1$ and $x_2$, but you are more interested in their difference, $\delta$. | *Situation*: You have results from different business units, $x_i$, but your ML model should not consider the results separately, but as an aggregated overall result, $x$. | *Situation*: To identify common wind speed patterns, you have measurements of two orthogonal wind speed components, $u$ and $v$. Since only the magnitude of the resulting wind vector, $w$, matters, you can utilize its magnitude, $\lvert w \rvert$. |
| *Transformation*: $$\delta = x_1 - x_2$$ | *Transformation*: $$x = \sum_i x_i$$ | *Transformation*: $$\lvert w \rvert = \sqrt{u^2 + v^2}$$  |

# Feature engineering – qualitative data

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

# Feature engineering – qualitative data

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**:  ordinal (ranked) data → discrete quantitative data
  The intuition is that the ranking/order of the classes is conserved in a discrete numerical schema and a "distance" can be defined.

  *Examples*:

  - Competition ranks: [1st, 2nd, 3rd, 4th, 5th] → [1, 2, 3, 4, 5]

  - Cloudiness scale: [clear, mostly clear, partly cloudy, mostly cloudy] → [0, 1, 2, 3]

  - Quality scale: [very good, good, satisfying, sufficient, insufficient] → [0, 1, 2, 3, 4]

  - Days of the week: [Mon, Tue, Wed, Thu, Fri, Sat, Sun] → [1, 2, 3, 4, 5, 6, 7]

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**:  ordinal (ranked) data → discrete quantitative data
  The intuition is that the ranking/order of the classes is conserved in a discrete numerical schema and a "distance" can be defined.

*Examples*:

  - Competition ranks: [1st, 2nd, 3rd, 4th, 5th] → [1, 2, 3, 4, 5]

  - Cloudiness scale: [clear, mostly clear, partly cloudy, mostly cloudy] → [0, 1, 2, 3]

  - Quality scale: [very good, good, satisfying, sufficient, insufficient] → [0, 1, 2, 3, 4]

  - Days of the week: [Mon, Tue, Wed, Thu, Fri, Sat, Sun] → [1, 2, 3, 4, 5, 6, 7]

      ← be careful: day of week is cyclical!

# Feature engineering – qualitative data

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**:  ordinal (ranked) data → discrete quantitative data
  The intuition is that the ranking/order of the classes is conserved in a discrete numerical schema and a "distance" can be defined.

  *Examples*:

  - Competition ranks: [1st, 2nd, 3rd, 4th, 5th] → [1, 2, 3, 4, 5]

  - Cloudiness scale: [clear, mostly clear, partly cloudy, mostly cloudy] → [0, 1, 2, 3]

  - Quality scale: [very good, good, satisfying, sufficient, insufficient] → [0, 1, 2, 3, 4]

  - Days of the week: [Mon, Tue, Wed, Thu, Fri, Sat, Sun] → [1, 2, 3, 4, 5, 6, 7]       ← be careful: day of week is cyclical!

(Caveat: Label encoding can also be used if a large number of classes is present)

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**: *ordinal (ranked) data → discrete quantitative data*

# Feature engineering – qualitative data

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**:  *ordinal (ranked) data → discrete quantitative data*

- **One-hot encoding**: *nominal (unranked) data → binary coding of labels*
  For each possible class in a feature, a binary feature is introduced; for each sample, all one-hot features are zero, only those that match have a value of one.
  Examples:

    - House properties: [balcony, cellar, fireplace, jacuzzi]      →
      samples:      house 1: "balcony"                              →
                    house 2: "fireplace"                            →
                    house 3: "balcony and jacuzzi"                  →
                    house 4: "cellar, fireplace and jacuzzi"        →

# Feature engineering – qualitative data

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**:  *ordinal (ranked) data → discrete quantitative data*

- **One-hot encoding**: *nominal (unranked) data → binary coding of labels*
  For each possible class in a feature, a binary feature is introduced; for each sample, all one-hot features are zero, only those that match have a value of one.
  Examples:

  - House properties: [balcony, cellar, fireplace, jacuzzi]   →
    samples:     house 1: "balcony"                           →
                 house 2: "fireplace"                          →
                 house 3: "balcony and jacuzzi"                →
                 house 4: "cellar, fireplace and jacuzzi"      →

| balcony | cellar | fireplace | jacuzzi |
|---------|--------|-----------|---------|
| 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 |

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**: *ordinal (ranked) data → discrete quantitative data*

- **One-hot encoding**: *nominal (unranked) data → binary coding of labels*
  For each possible class in a feature, a binary feature is introduced; for each sample, all one-hot features are zero, only those that match have a value of one.
  Examples:

- House properties: [balcony, cellar, fireplace, jacuzzi]    →
  samples:    house 1: "balcony"    →
  house 2: "fireplace"    →
  Multi-class ⎰ house 3: "balcony and jacuzzi"    →
  feature ⎱ house 4: "cellar, fireplace and jacuzzi"    →

| balcony | cellar | fireplace | jacuzzi |
|---------|--------|-----------|---------|
| 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 |

Qualitative (categorical) data cannot be fed into ML models directly, they have to be turned into quantitative data first. There are two common methods available, depending on the data type:

- **Label encoding**: *ordinal (ranked) data → discrete quantitative data*

- **One-hot encoding**: *nominal (unranked) data → binary coding of labels*
  For each possible class in a feature, a binary feature is introduced; for each sample, all one-hot features are zero, only those that match have a value of one.
  Examples:

|  | balcony | cellar | fireplace | jacuzzi |
|---|---|---|---|---|
| House properties: [balcony, cellar, fireplace, jacuzzi] → | | | | |
| samples: house 1: "balcony" → | 1 | 0 | 0 | 0 |
| house 2: "fireplace" → | 0 | 0 | 1 | 0 |
| Multi-class { house 3: "balcony and jacuzzi" → | 1 | 0 | 0 | 1 |
| feature { house 4: "cellar, fireplace and jacuzzi" → | 0 | 1 | 1 | 1 |

(Caveat: if too many classes present, use label encoding instead; see *curse of dimensionality*)

# Final data set nomenclature

Feature engineering results in a compilation of features that we can use to train our ML models.

*Example*:

| Weight | Height | Wings | Legs | Cuteness |
|--------|--------|-------|------|----------|
| 0.1 | 0.1 | true | 2 | 1 |
| 3.5 | 0.3 | false | 4 | 1 |
| 12.0 | 0.7 | false | 4 | 1 |
| 500 | 1.8 | false | 4 | 2 |
| 800 | 3.0 | true | 4 | 3 |
| ... | ... | ... | ... | ... |

| Pet | Type |
|-----|------|
| true | bird |
| true | cat |
| true | dog |
| false | rhinoceros |
| false | chimera |
| ... | ... |

# Final data set nomenclature

Feature engineering results in a compilation of features that we can use to train our ML models.

*Example*:

**Features**/Attributes (input variables, $x$)    $f(x) = y$    **Targets**/Labels (output variables, $y$)
**Ground-Truth**

| Weight | Height | Wings | Legs | Cuteness |
|--------|--------|-------|------|----------|
| 0.1 | 0.1 | true | 2 | 1 |
| 3.5 | 0.3 | false | 4 | 1 |
| 12.0 | 0.7 | false | 4 | 1 |
| 500 | 1.8 | false | 4 | 2 |
| 800 | 3.0 | true | 4 | 3 |
| ... | ... | ... | ... | ... |

| Pet | Type |
|------|------|
| true | bird |
| true | cat |
| true | dog |
| false | rhinoceros |
| false | chimera |
| ... | ... |

# Final data set nomenclature

Feature engineering results in a compilation of features that we can use to train our ML models.

*Example*:

**Features**/Attributes (input variables, $x$)  $f(x) = y$  **Targets**/Labels (output variables, $y$)

**Ground-Truth**

| | Weight | Height | Wings | Legs | Cuteness | | Pet | Type |
|---|---|---|---|---|---|---|---|---|
| | 0.1 | 0.1 | true | 2 | 1 | | true | bird |
| | 3.5 | 0.3 | false | 4 | 1 | | true | cat |
| | 12.0 | 0.7 | false | 4 | 1 | | true | dog |
| | 500 | 1.8 | false | 4 | 2 | | false | rhinoceros |
| | 800 | 3.0 | true | 4 | 3 | | false | chimera |
| | ... | ... | ... | ... | ... | | ... | ... |

**Samples**/Instances

# Final data set nomenclature

Feature engineering results in a compilation of features that we can use to train our ML models.

*Example*:

**Features**/Attributes (input variables, $x$)   $f(x) = y$   **Targets**/Labels (output variables, $y$)
**Ground-Truth**

| | Weight | Height | Wings | Legs | Cuteness | | Pet | Type |
|---|---|---|---|---|---|---|---|---|
| **Samples**/Instances | 0.1 | 0.1 | true | 2 | 1 | | true | bird |
| | 3.5 | 0.3 | false | 4 | 1 | | true | cat |
| | 12.0 | 0.7 | false | 4 | 1 | | true | dog |
| | 500 | 1.8 | false | 4 | 2 | | false | rhinoceros |
| | 800 | 3.0 | true | 4 | 3 | | false | chimera |
| | ... | ... | ... | ... | ... | | ... | ... |

**classes** of label "Type"

Feature engineering results in a compilation of features that we can use to train our ML models.

*Example*:

**Features**/Attributes (input variables, $x$)     $f(x) = y$     **Targets**/Labels (output variables, $y$)

**Ground-Truth**

| | Weight | Height | Wings | Legs | Cuteness | | Pet | Type | |
|---|---|---|---|---|---|---|---|---|---|
| **Samples**/Instances | 0.1 | 0.1 | true | 2 | 1 | | true | bird | |
| | 3.5 | 0.3 | false | 4 | 1 | | true | cat | |
| | 12.0 | 0.7 | false | 4 | 1 | | true | dog | **classes** of label "Type" |
| | 500 | 1.8 | false | 4 | 2 | | false | rhinoceros | |
| | 800 | 3.0 | true | 4 | 3 | | false | chimera | |
| | ... | ... | ... | ... | ... | | ... | ... | |

Data Types:

continuous        binary        ordinal        categorical (multi-class)

continuous        discrete        binary

# Data scaling

Data scaling means to linearly transform your data in order to normalize them.

Data scaling means to linearly transform your data in order to normalize them.

**Why scale data?**

Data scaling means to linearly transform your data in order to normalize them.
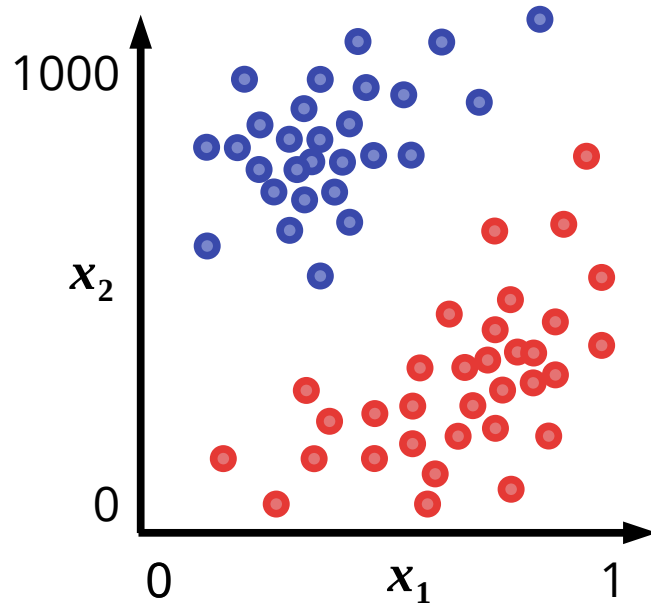
**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.
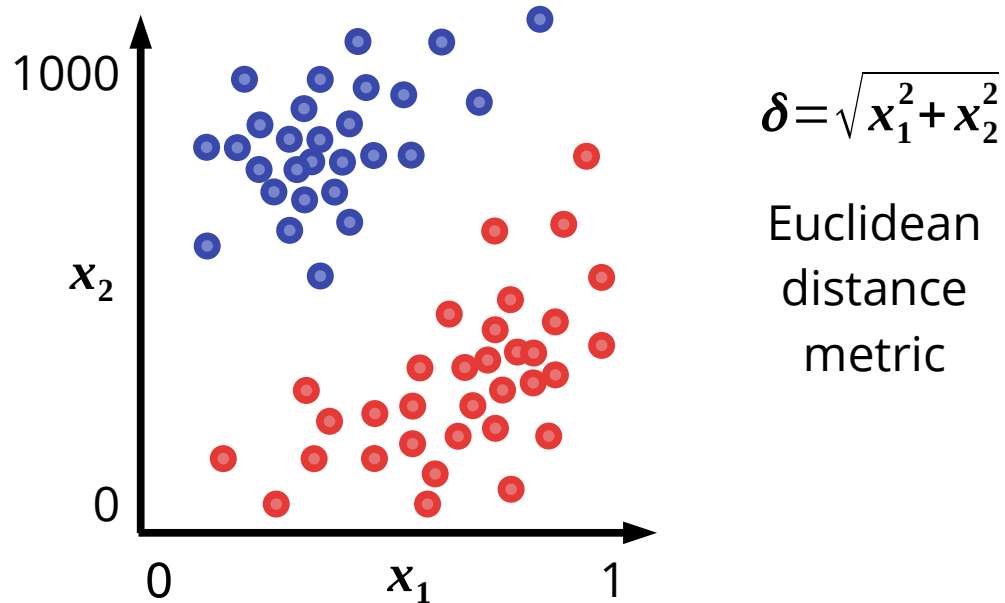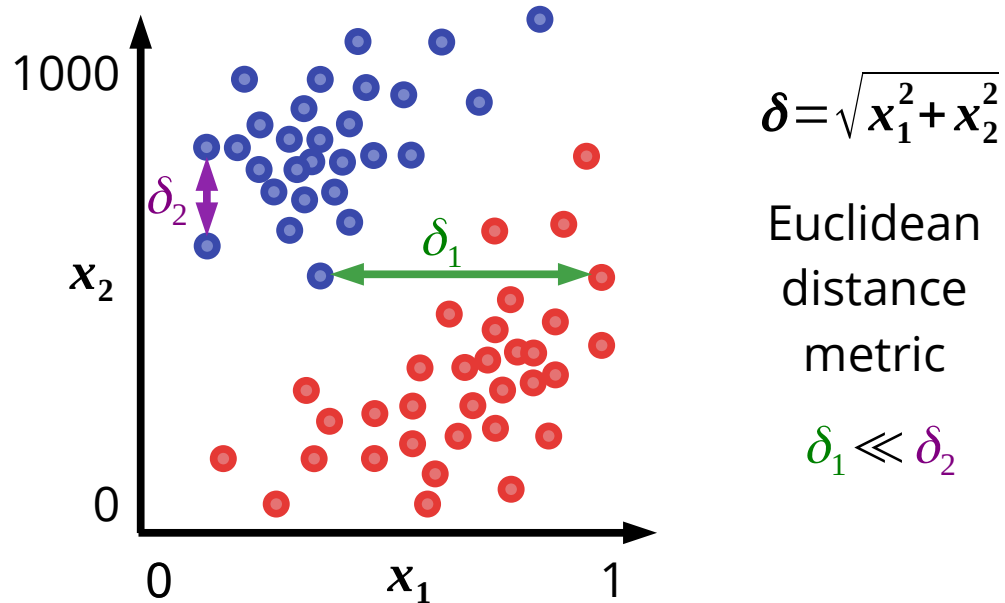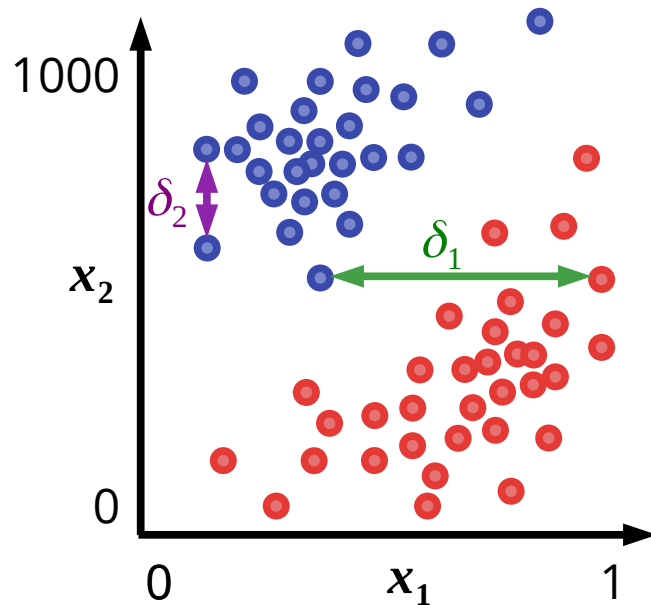
Data scaling means to linearly transform your data in order to normalize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.

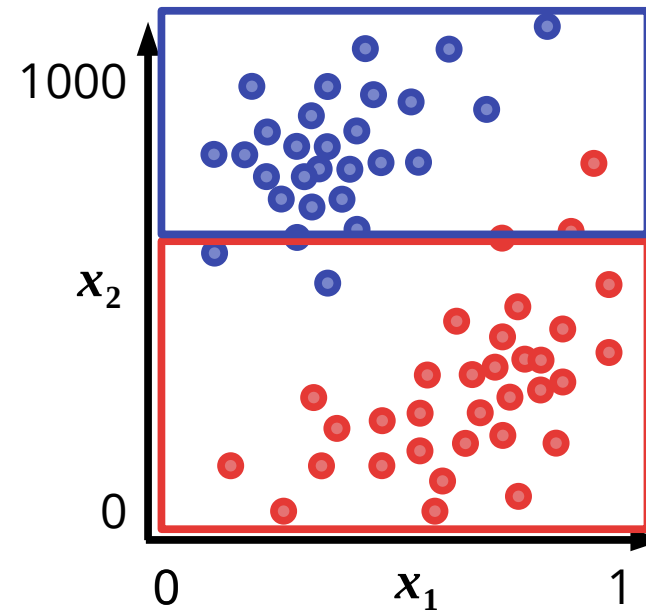Data scaling means to linearly transform your data in order to normalize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.



$$\delta = \sqrt{x_1^2 + x_2^2}$$

Euclidean
distance
metric

Data scaling means to linearly transform your data in order to normalize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.

$$\delta = \sqrt{x_1^2 + x_2^2}$$

Euclidean distance metric
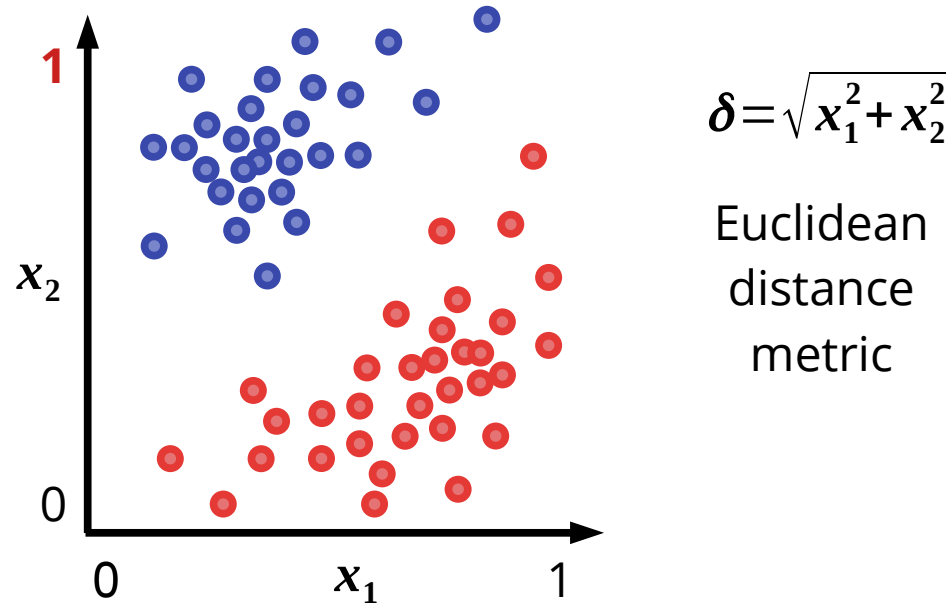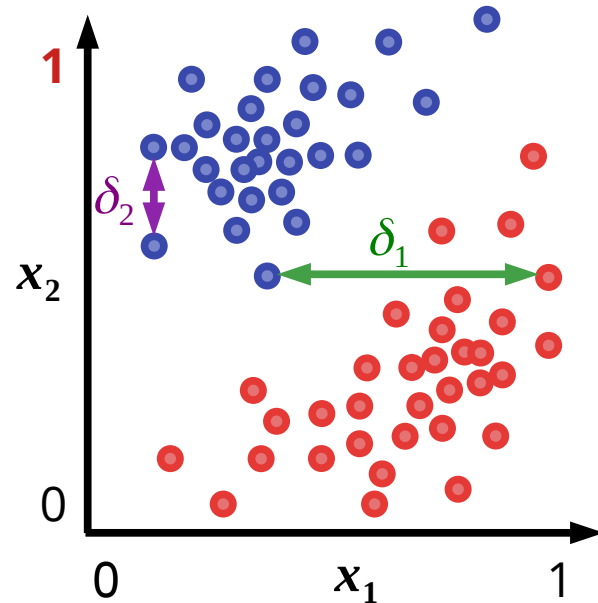
$$\delta_1 \ll \delta_2$$

Data scaling means to linearly transform your data in order to normalize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.
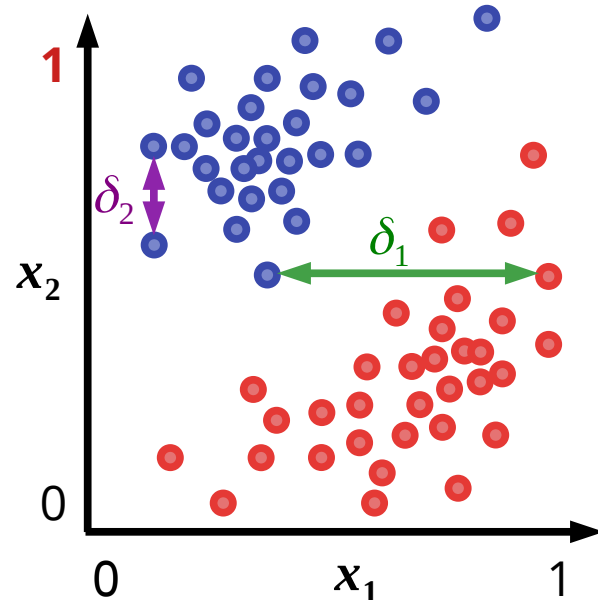


$$\delta = \sqrt{x_1^2 + x_2^2}$$

Euclidean distance metric

$$\delta_1 \ll \delta_2$$

Decision regions of a hypothetical distance-based classifier.

Results are ok-ish, but could be much better...

Data scaling means to linearly transform your data in order to standardize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.

$$\delta = \sqrt{x_1^2 + x_2^2}$$

Euclidean
distance
metric

Data scaling means to linearly transform your data in order to standardize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.

$$\delta = \sqrt{x_1^2 + x_2^2}$$

Euclidean distance metric

$$\delta_1 > \delta_2$$

Data scaling means to linearly transform your data in order to standardize them.
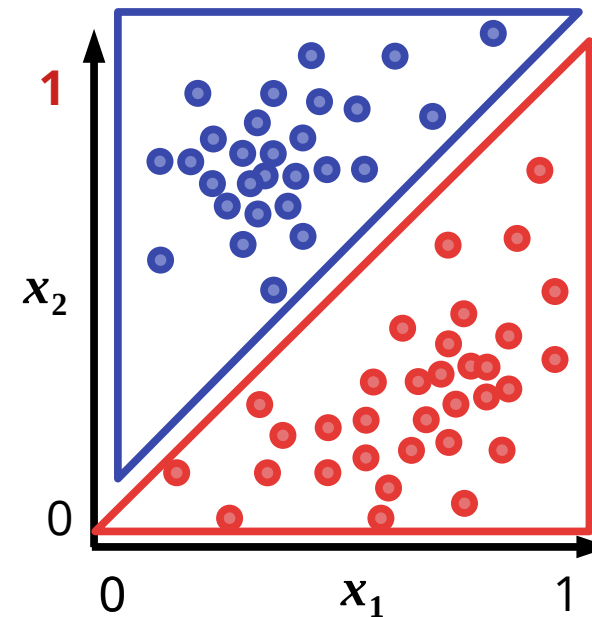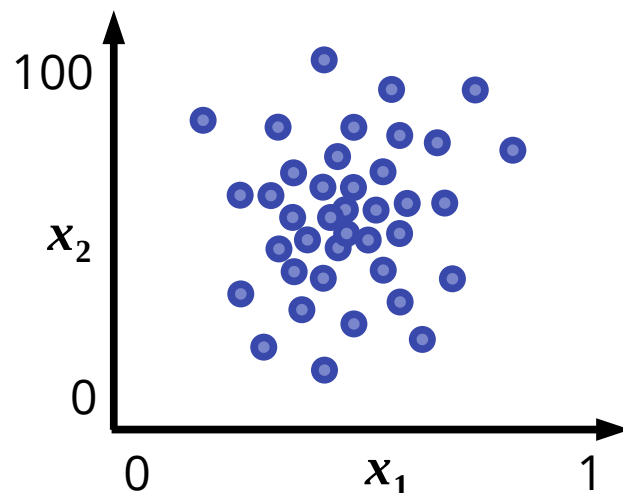
## Why scale data?

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.



$$\delta = \sqrt{x_1^2 + x_2^2}$$

Euclidean distance metric

$$\delta_1 > \delta_2$$

Decision regions of a hypothetical distance-based classifier.

This is much better!

Data should be scaled!

# Data scaling

Data scaling means to linearly transform your data in order to standardize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.
- Some ML models intrinsically presume that data are distributed following a Gaussian fashion with similar variances along all features; high variance along one feature leads to bias.
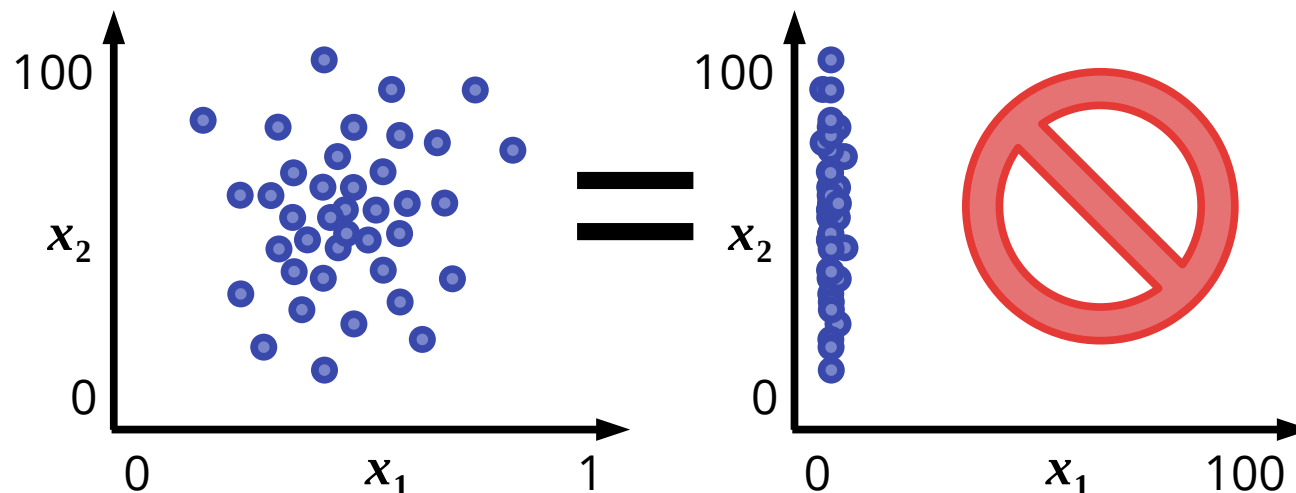
# Data scaling

Data scaling means to linearly transform your data in order to standardize them.
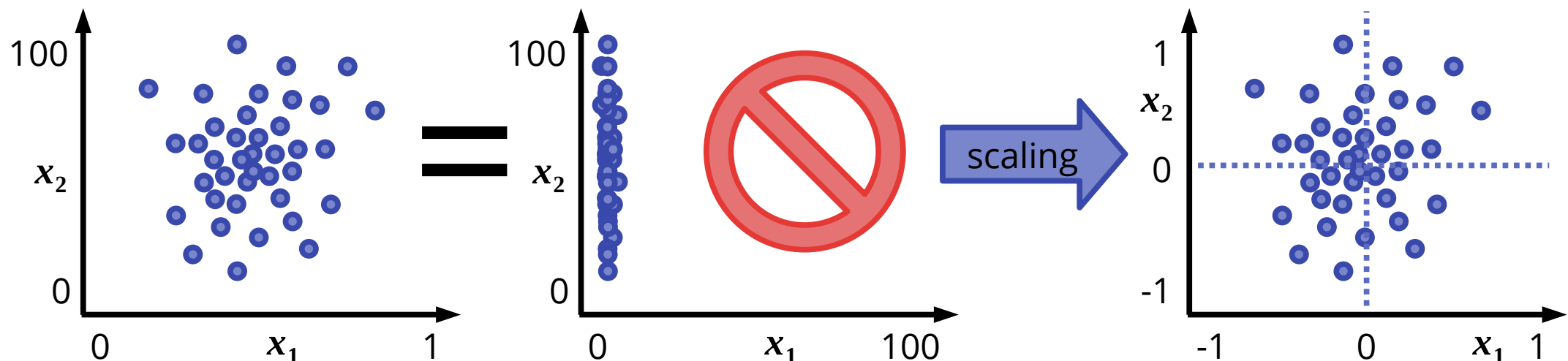
**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.
- Some ML models intrinsically presume that data are distributed following a Gaussian fashion with similar variances along all features; high variance along one feature leads to bias.

# Data scaling

Data scaling means to linearly transform your data in order to standardize them.
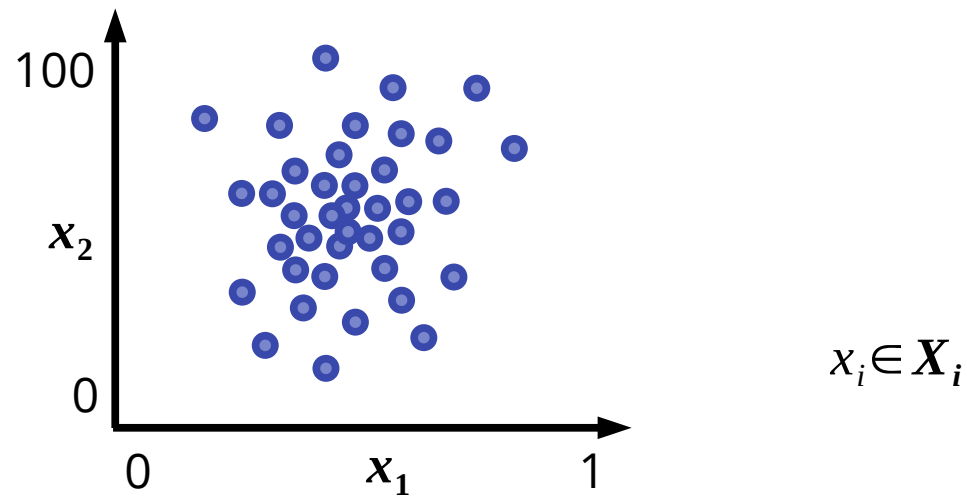
**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.
- Some ML models intrinsically presume that data are distributed following a Gaussian fashion with similar variances along all features; high variance along one feature leads to bias.

Data scaling means to linearly transform your data in order to standardize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.
- Some ML models intrinsically presume that data are distributed following a Gaussian fashion with similar variances along all features; high variance along one feature leads to bias.

# Data scaling

Data scaling means to linearly transform your data in order to standardize them.

**Why scale data?**

- Many ML models are based on a notion of "distance" between samples; improperly scaled data may jeopardize the learning capability of such models.
- Some ML models intrinsically presume that data are distributed following a Gaussian fashion with similar variances along all features; high variance along one feature leads to bias.

**How to scale data?**

- Normalize feature variances  (to give similar weights to the different features)
- Normalize feature mean values (assumed by a number of ML models)

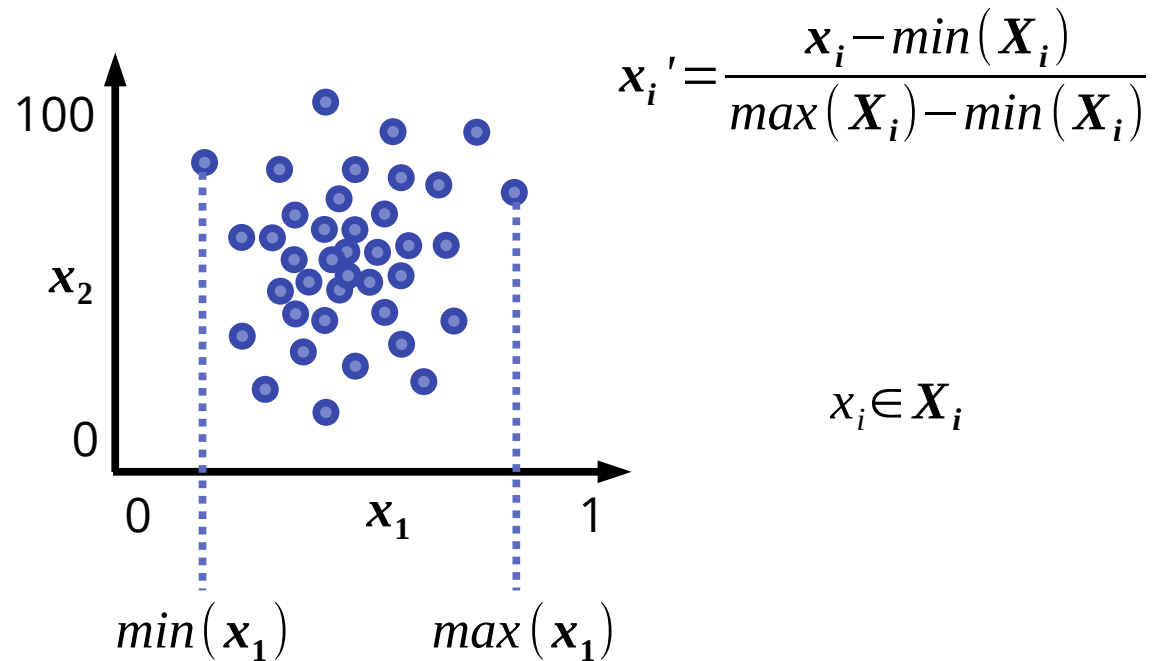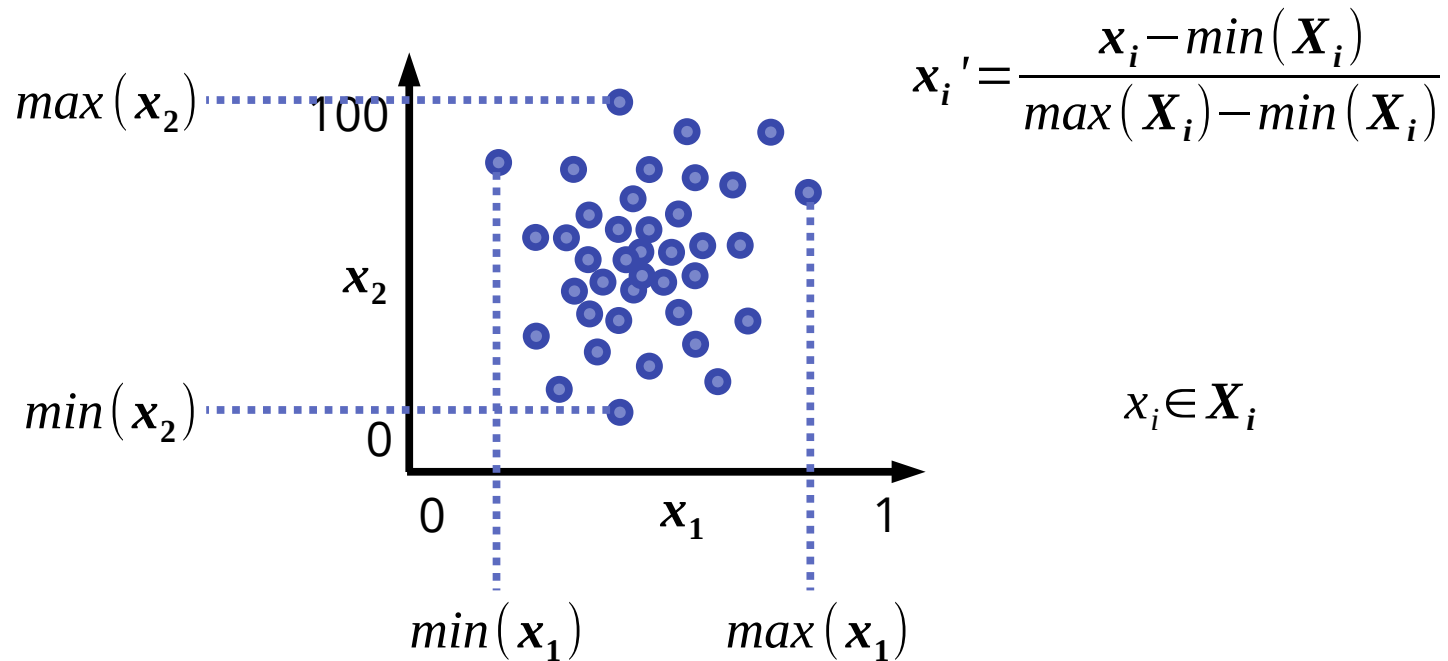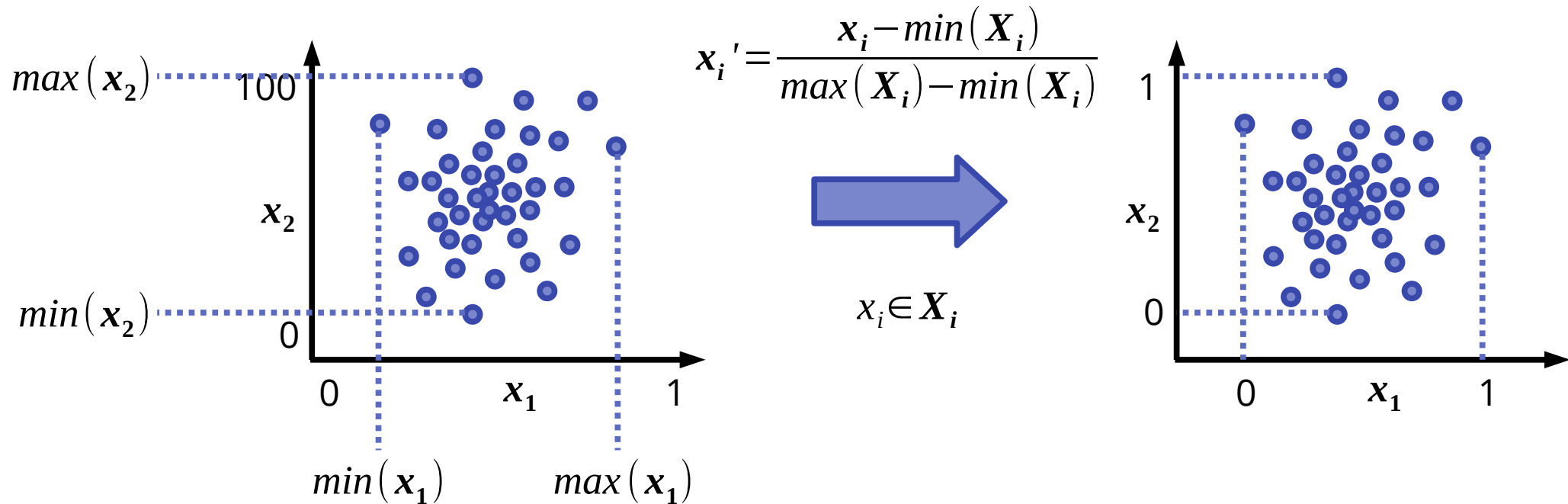Scale every feature onto a range from 0 to 1 based on the minimum and maximum of the underlying distribution.



$$x_i \in \boldsymbol{X_i}$$

Scale every feature onto a range from 0 to 1 based on the minimum and maximum of the underlying distribution.



$$x_i{}' = \frac{x_i - min(\boldsymbol{X_i})}{max(\boldsymbol{X_i}) - min(\boldsymbol{X_i})}$$

$$x_i \in \boldsymbol{X_i}$$

Scale every feature onto a range from 0 to 1 based on the minimum and maximum of the underlying distribution.

$$x_i' = \frac{x_i - min(X_i)}{max(X_i) - min(X_i)}$$

$$x_i \in X_i$$

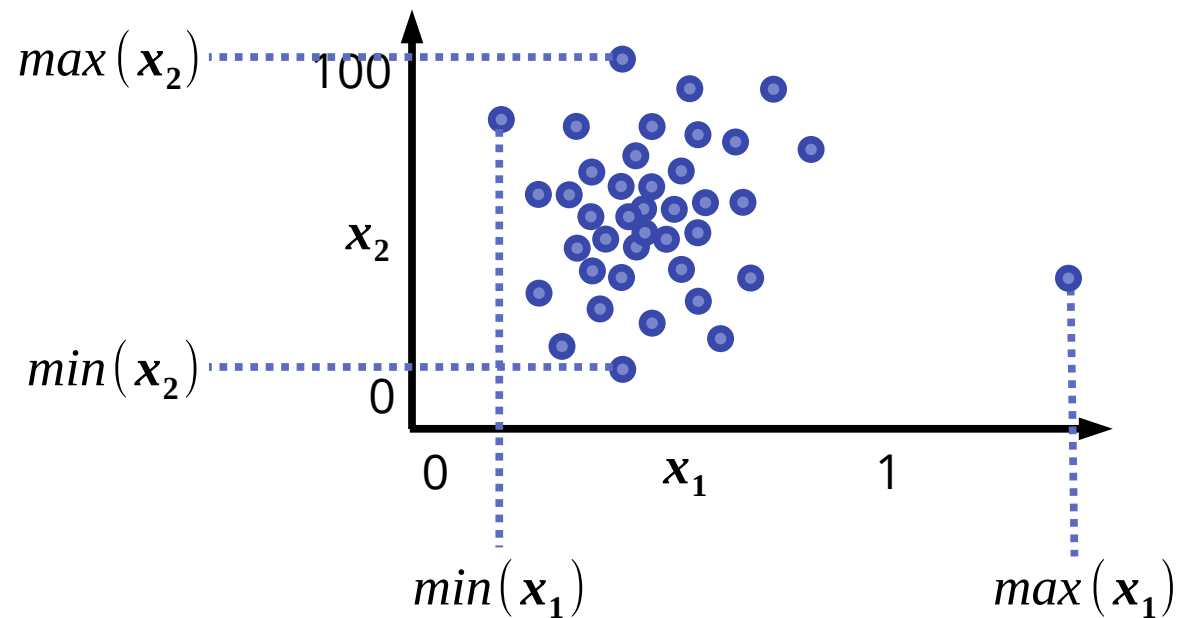Scale every feature onto a range from 0 to 1 based on the minimum and maximum of the underlying distribution.



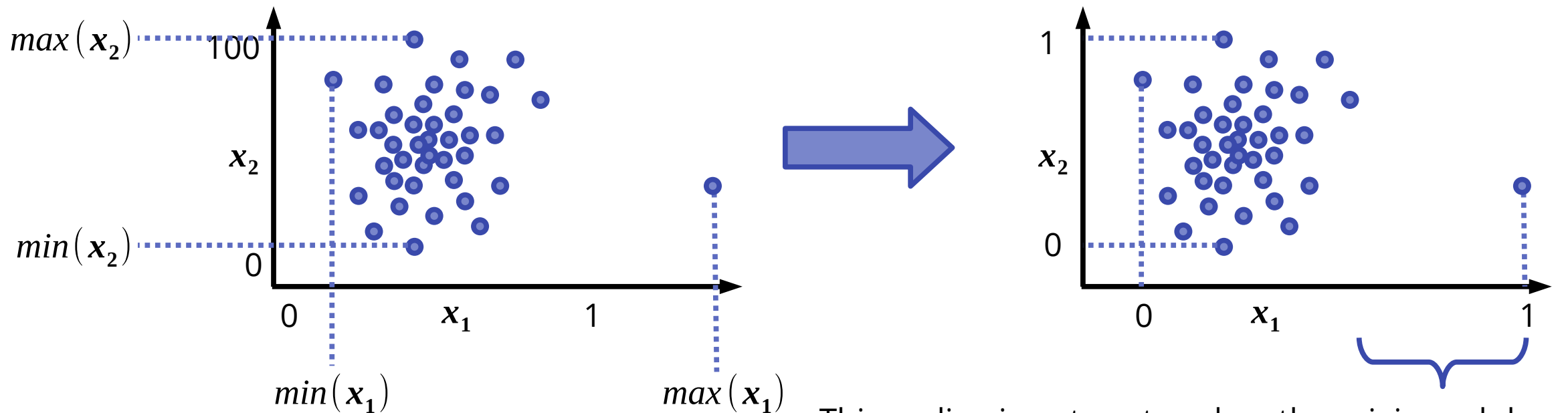$$x_i' = \frac{x_i - min(X_i)}{max(X_i) - min(X_i)}$$

$$x_i \in X_i$$

Scale every feature onto a range from 0 to 1 based on the minimum and maximum of the underlying distribution.



$$x_i' = \frac{x_i - min(X_i)}{max(X_i) - min(X_i)}$$

$$x_i \in X_i$$

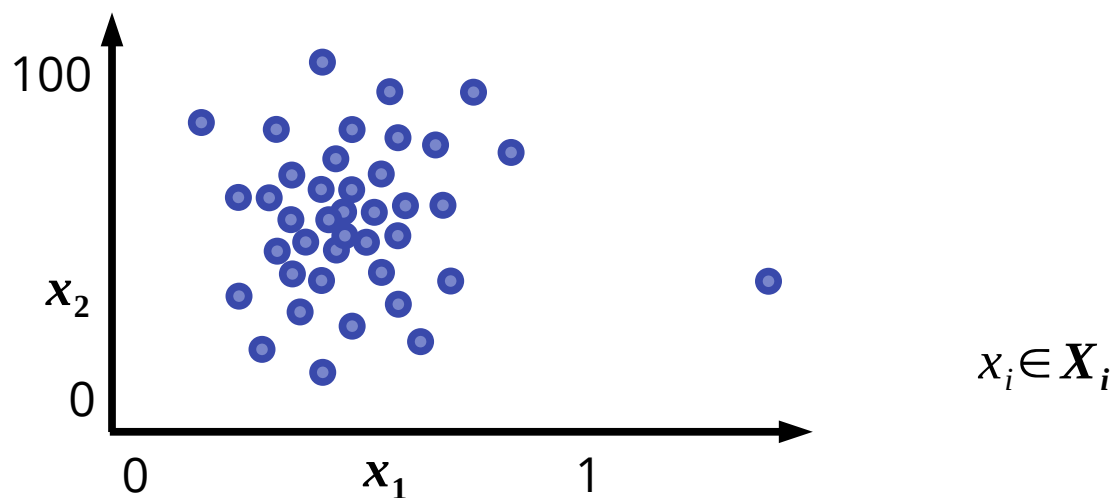Scale every feature onto a range from 0 to 1 based on the minimum and maximum sof the underlying distribution.

*Disadvantage*: the MinMax scaler is prone to outliers and does not center the distribution in the origin.

Scale every feature onto a range from 0 to 1 based on the minimum and maximum $s$of the underlying distribution.

*Disadvantage*: the MinMax scaler is prone to outliers and does not center the distribution in the origin.



This scaling is not centered on the origin and does not describe the data distribution well.
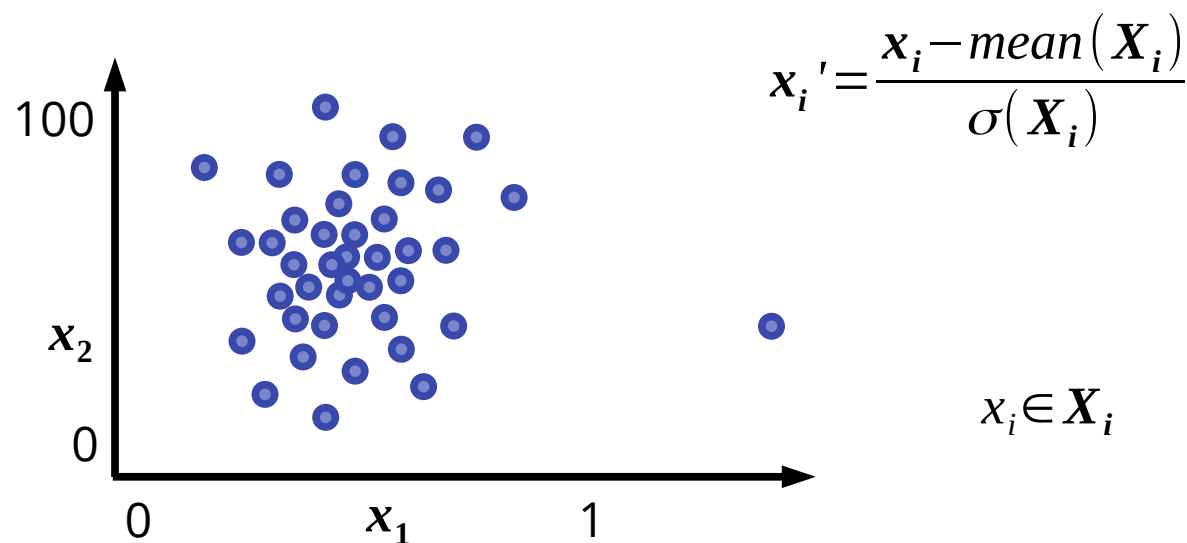
Scale every feature onto a range from -1 to 1 based on the mean and standard deviation of the underlying distribution.



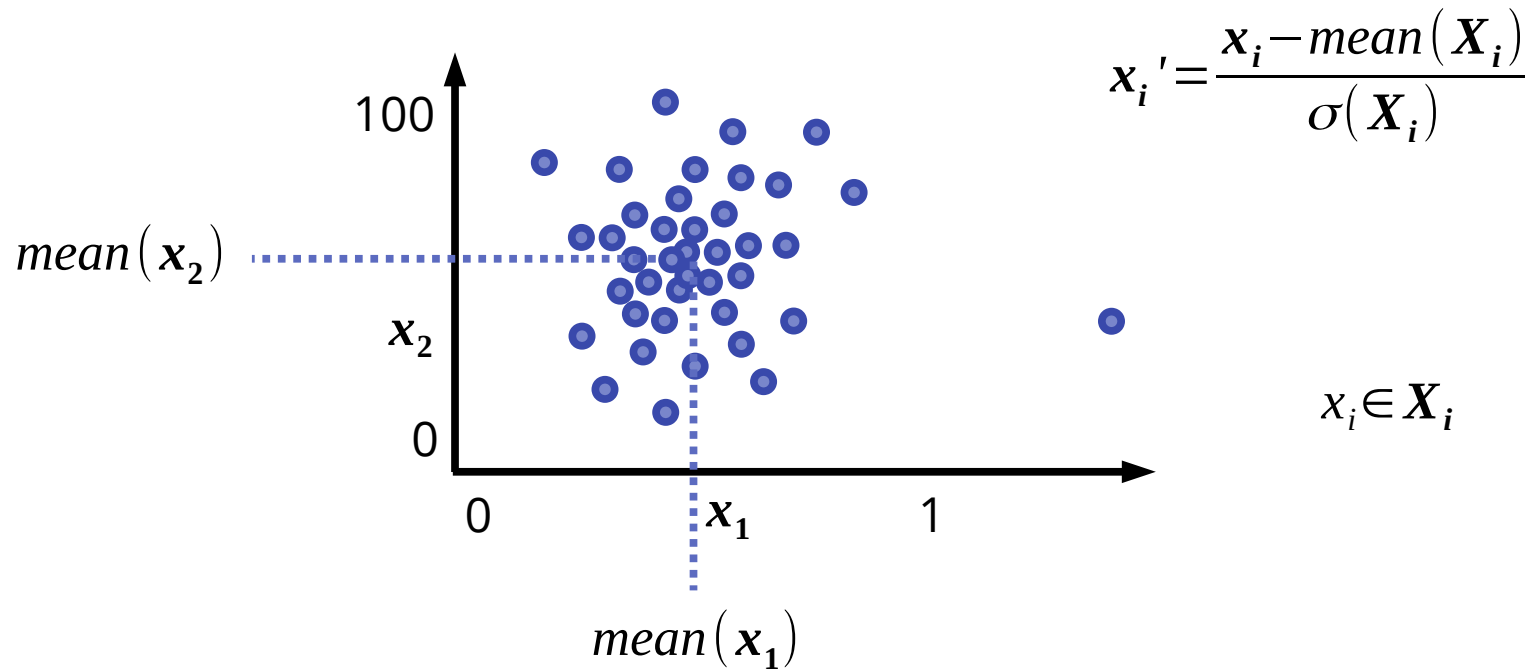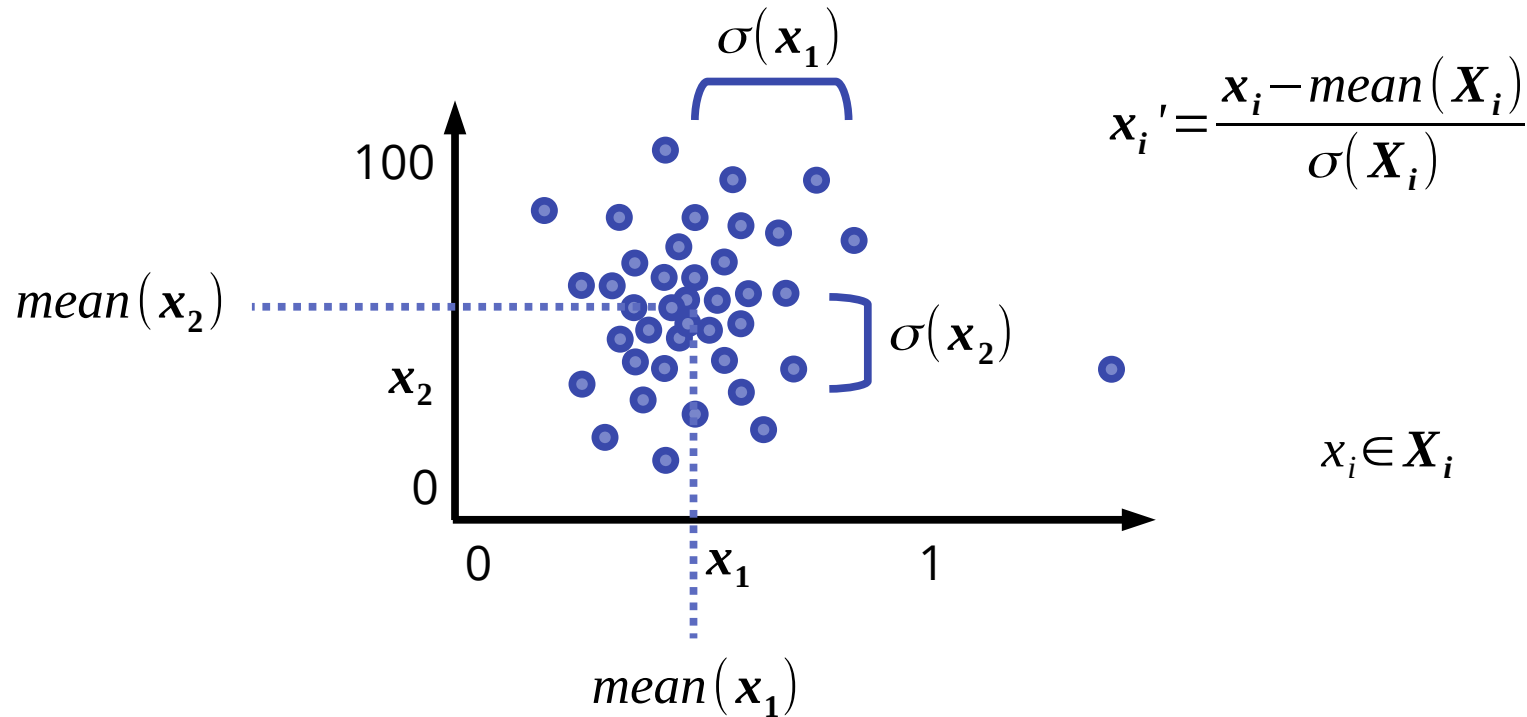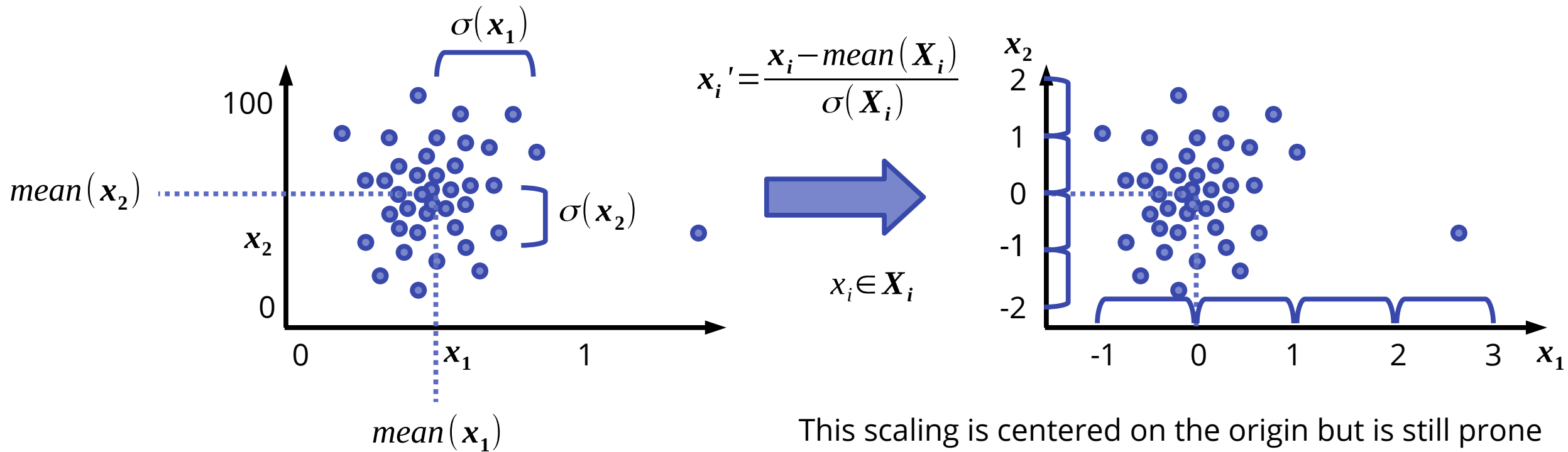$$x_i \in \boldsymbol{X_i}$$

# Data scaling – Standard scaler

Scale every feature onto a range from -1 to 1 based on the mean and standard deviation of the underlying distribution.

$$x_i' = \frac{x_i - mean(X_i)}{\sigma(X_i)}$$

$$x_i \in X_i$$

Scale every feature onto a range from -1 to 1 based on the mean and standard deviation of the underlying distribution.

$$x_i' = \frac{x_i - mean(X_i)}{\sigma(X_i)}$$

$$x_i \in X_i$$

Scale every feature onto a range from -1 to 1 based on the mean and standard deviation of the underlying distribution.



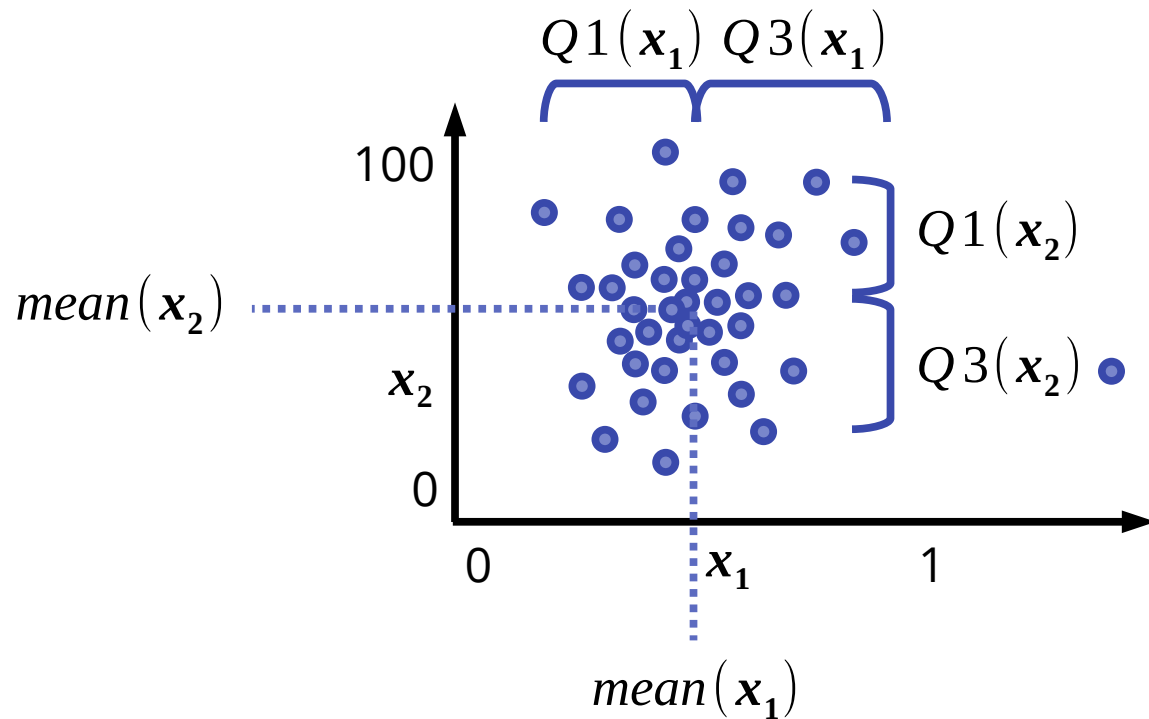$$x_i' = \frac{x_i - mean(X_i)}{\sigma(X_i)}$$

$$x_i \in X_i$$

# Data scaling – Standard scaler

Scale every feature onto a range from -1 to 1 based on the mean and standard deviation of the underlying distribution.



$$x_i' = \frac{x_i - mean(X_i)}{\sigma(X_i)}$$

$$x_i \in X_i$$

This scaling is centered on the origin but is still prone to outliers to some extent.

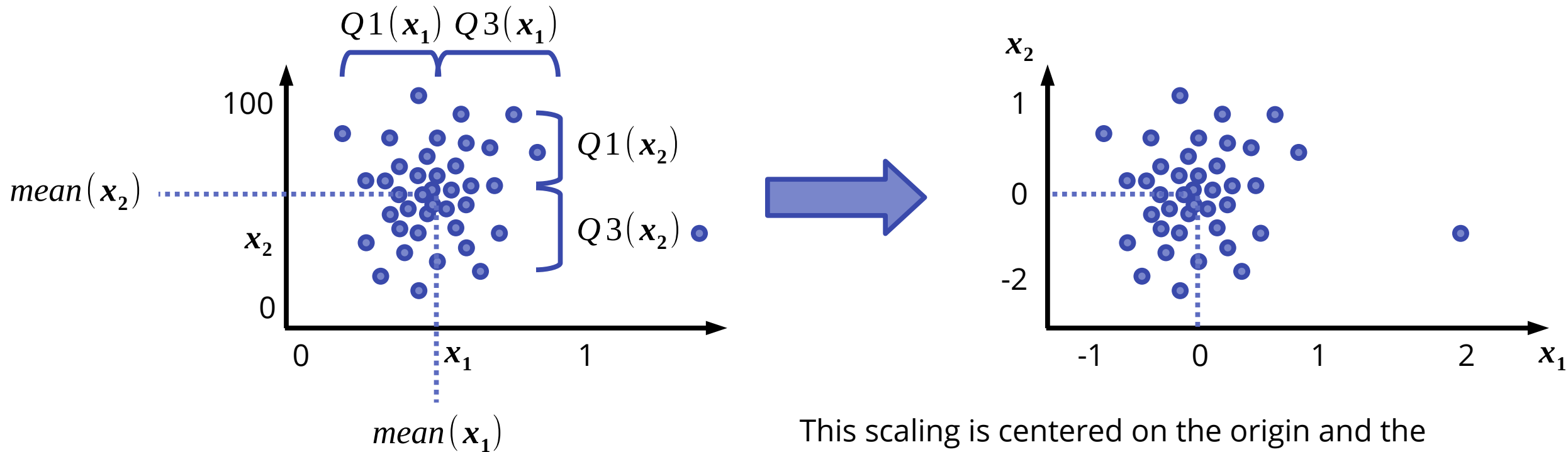Scale every feature onto a range from -1 to 1 based on the mean and the quantiles of the underlying distribution.

Scale every feature onto a range from -1 to 1 based on the mean and the quantiles of the underlying distribution.



This scaling is centered on the origin and the resulting distribution is less affected by outliers

# That's all folks!