# Introduction to
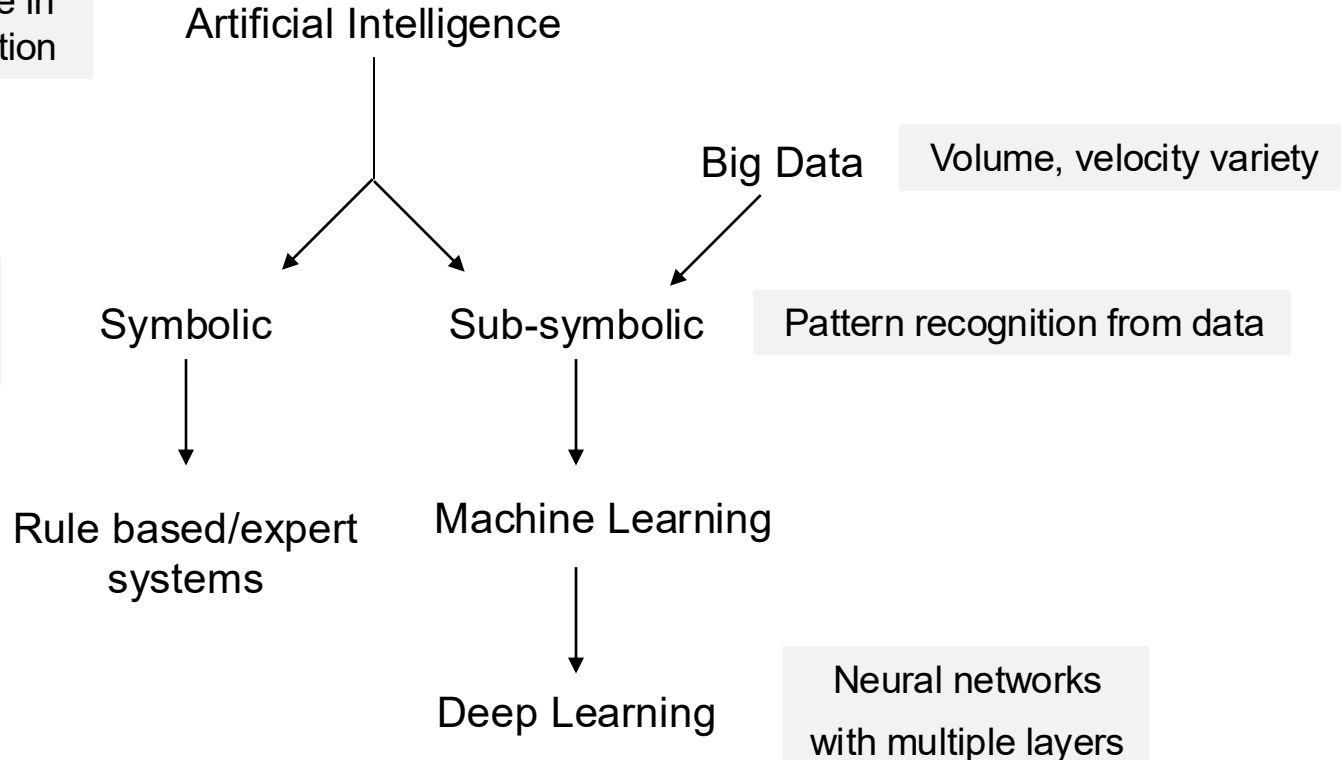# Machine Learning

# AI and Big Data

Simulate human intelligence in reasoning, learning, preception
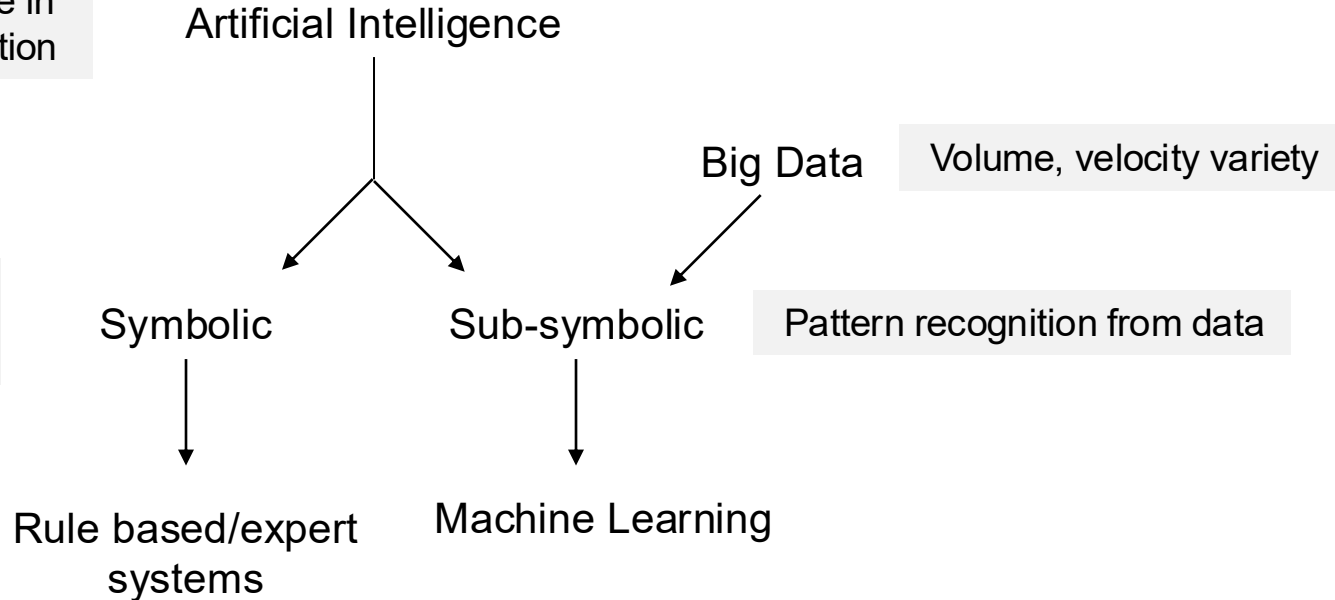
Artificial Intelligence

Big Data — Volume, velocity variety

Information is represented in human-readable form

Symbolic

Sub-symbolic — Pattern recognition from data

Rule based/expert systems

Machine Learning

Deep Learning — Neural networks with multiple layers

# AI and Big Data

Simulate human intelligence in reasoning, learning, preception

Artificial Intelligence

Big Data

Volume, velocity variety

Information is represented in human-readable form

Symbolic

Sub-symbolic

Pattern recognition from data
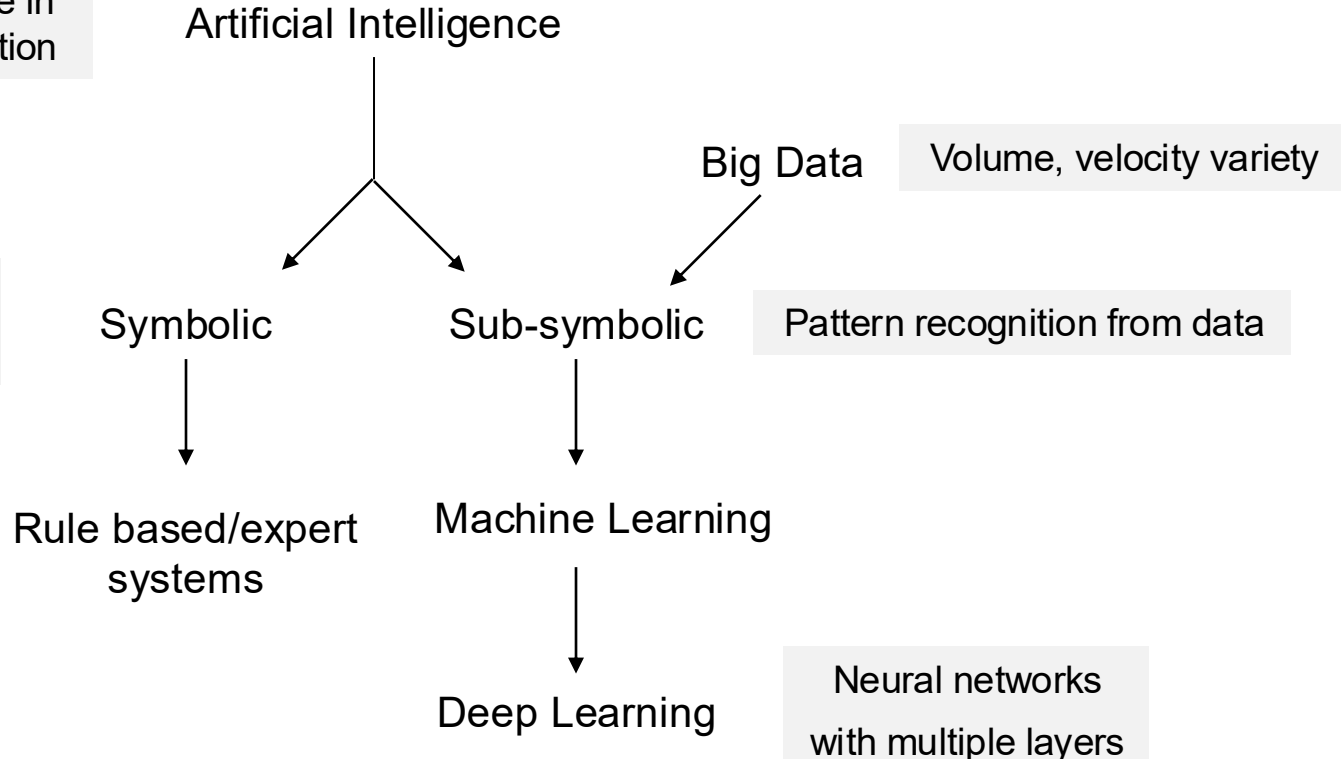
Rule based/expert systems

Machine Learning

„Machine learning is a field of study that gives computers the ability to learn without being explicitly programmed." (Arthur Samuel, 1959)

# AI and Big Data

Simulate human intelligence in reasoning, learning, preception

Artificial Intelligence

Big Data

Volume, velocity variety

Information is represented in human-readable form

Symbolic

Sub-symbolic

Pattern recognition from data

Rule based/expert systems

Machine Learning

Deep Learning

Neural networks with multiple layers

# What is machine learning used for?

The computational methods in Machine learning
are used to discover patterns in the data and/or derive a corresponding generating process to

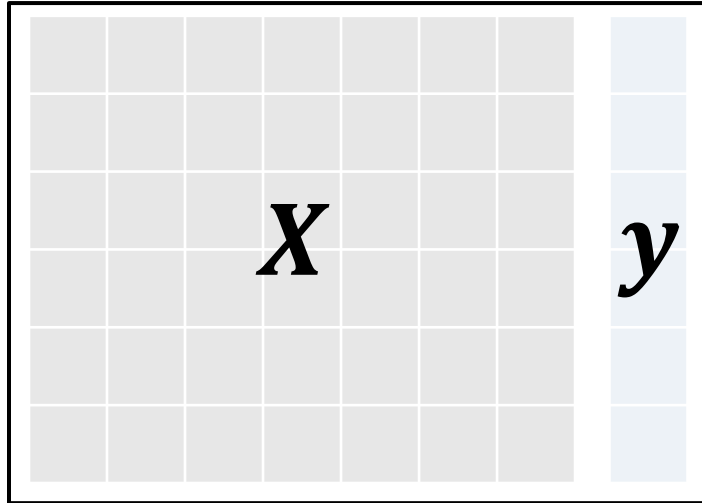1) gain insights

2) predict events

In order to
- provide a quantitative basis for decisions (actionable insights)
  e.g. determine target segment for marketing campagne

- influence the underlying process of the data
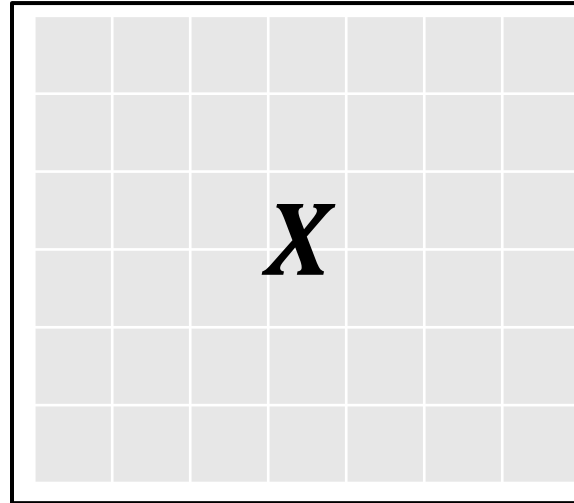  e.g. adapt the user features of an app

# Machine learning paradigms

# Supervised vs. unsupervised learning

### Supervised learning

### Unsupervised learning

$M$: Number of training samples

$N$: Number of features

Dimension $X$: $M \times N$

Dimensions $y$: $M$

The training data consists of input samples $\mathbf{x}_{m,:}$ and their associated output values $y_m$

The training data does not contain any output values

# Supervised Learning

**Goal:** Derive a model that is able to accurately predict output values from new input values

**Pre-requisite**: Training data - labeled samples (input features + output values)

**Approach**: find a function $f$, which systematically produces the output values $y_m$ associated with the input values $\mathbf{x}_{m,:}$ from the training data:

$$f\left(\mathbf{x}_{m,:}; \boldsymbol{\theta}\right) \rightarrow y_m$$

**Process**: Algorithm adapts parameters $\boldsymbol{\theta}$ of function $f$ to predict the correct outputs for the known training samples.

$\rightarrow$ Use $f$ to make predictions on new data (unseen during training)

# Model and Learning

A **model** is a mathematical, statistical, or logical representation that describes the relation-ship between variables and can be used to make predictions or understand patterns in data.

**Learning:** Machine Learning employs adaptive models, which are configured and parametrised automatically based on the training data.
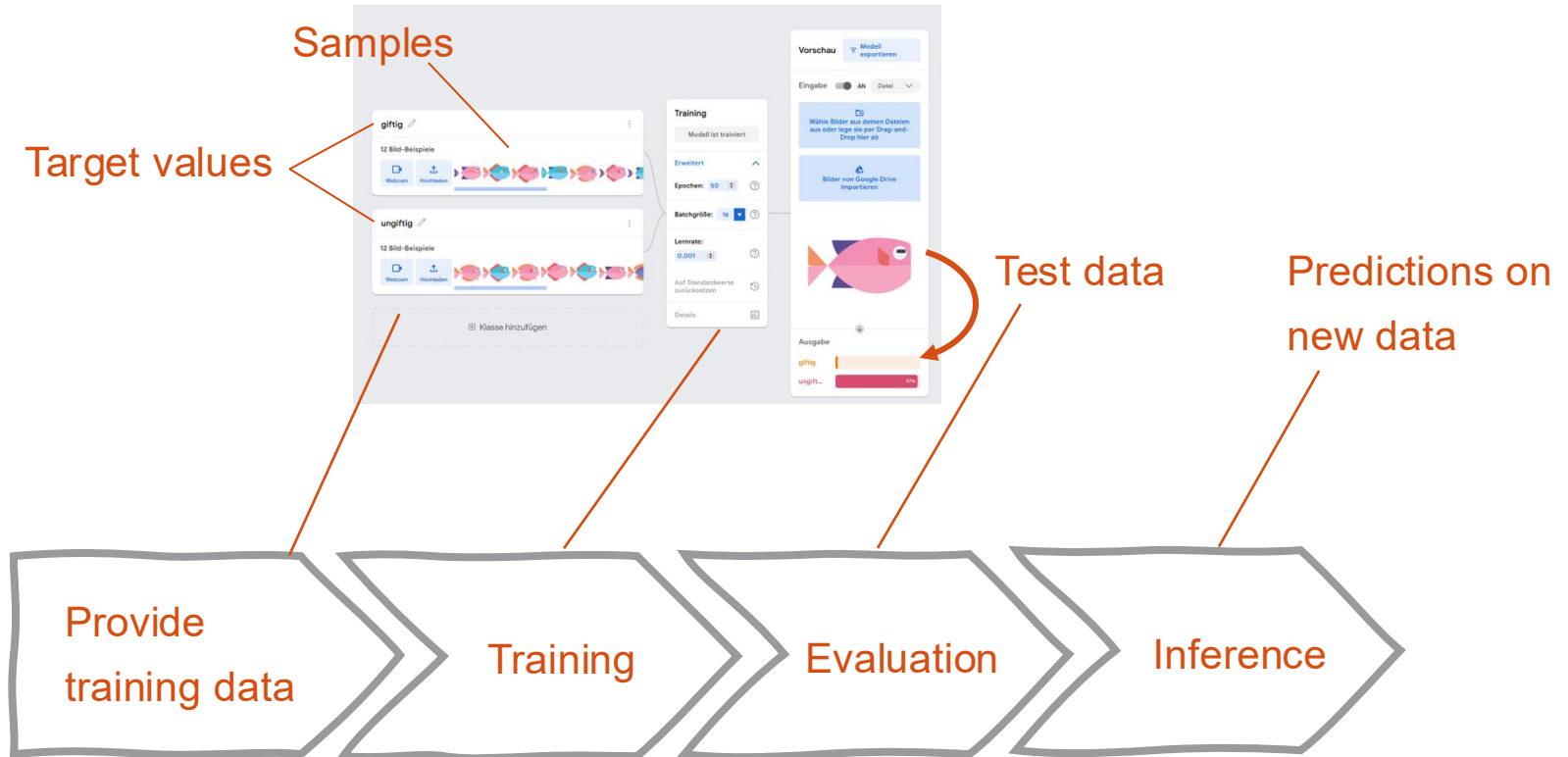
# AI in Action



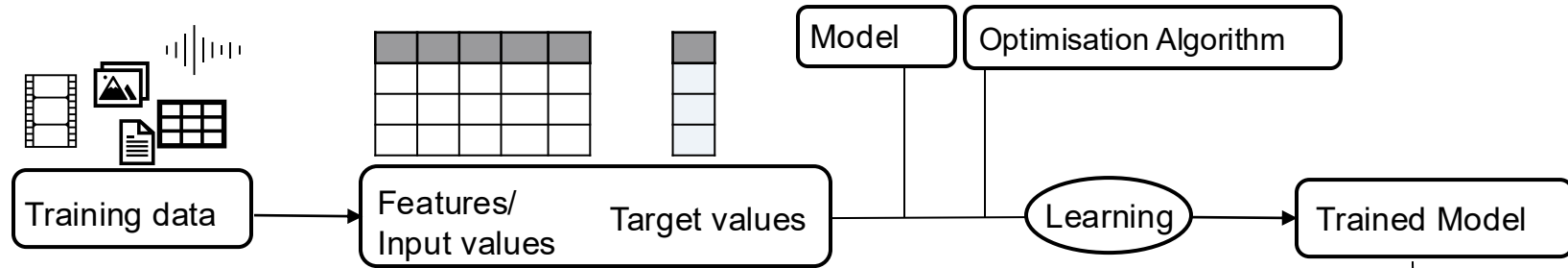**https://youtu.be/FnigvS_uI1w?feature=shared**

# Teachable Machine

The data: https://tinyurl.com/mvvhj2n5

The machine: https://teachablemachine.withgoogle.com

# Structure of teaching the Teachable Machine



Samples

Target values

Test data

Predictions on new data

Provide training data

Training

Evaluation

Inference

# Structure of a supervised learning problem

**Training**

Training data → Features/ Input values   Target values → Model | Optimisation Algorithm → Learning → Trained Model

**Evaluation**

Test data → Features/ Input values   Target values → Predictions on test data → Evaluate

**Inference**

New data → Features/ Input values → Predictions on new data

# Classification vs. Regression

In supervised learning we try to find a function $f$, which systematically produces the output values $y_m$ associated with the input values $\mathbf{x}_{m,:}$:

$$f(\mathbf{x}_{m,:}) \to y_m$$

**Classification**

**Target variable $y$: categorical**

$$y_m \in \{C_1, C_2, \dots, C_K\}$$



**Regression**

**Target variable $y$: numerical - continuous**
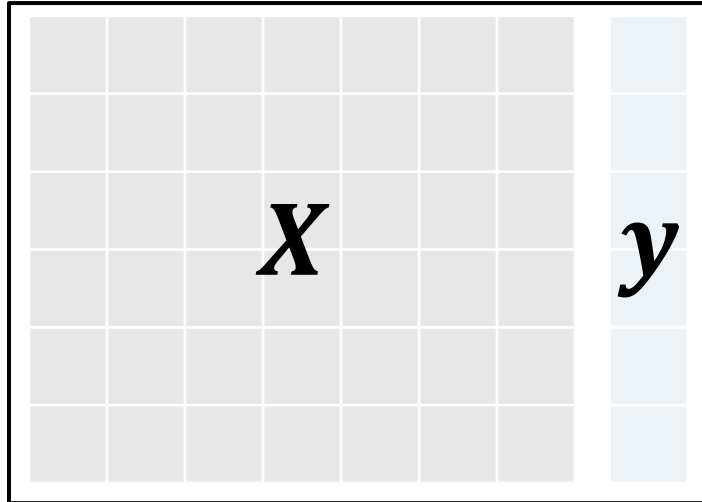
$$y_m \in \mathbb{R}$$

# Terminology

Input data: X
Output data: y
Sample: one row in X (and y)
Covariates = predictors = independent variables = features =attributes: columns of X
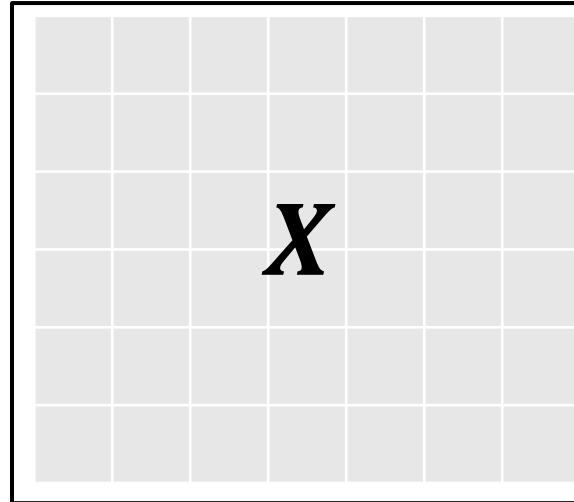Dependent variable, target variable, outputs, labels: y

# Supervised vs. unsupervised learning

Supervised learning

Unsupervised learning

$X$     $y$

$X$

The training data consists of input samples $\mathbf{x}_{m,:}$ and their associated output values $y_m$

The training data does not contain any output values

$M$: Number of training samples

$N$: Number of features

Dimension $X$: $M \times N$

Dimensions $y$: $M$

# Unsupervised learning

In unsupervised learning the goal is to model the underlying distribution without labels in the training data.

Tasks:

- Dimensionality reduction
- Clustering
- Anomaly detection



Challenges:

- Problem is often less clearly defined as in supervised learning
- Evaluation is difficult without labelled test data

# Dimensionality reduction

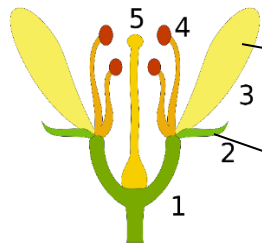Goal: Transforming the data into an optimal lower dimensional representation
- for visualisation (normally 2D, sometimes 3D) of the data
- to generate more informative features for supervised learning

Some methods for dimensionality reduction:

- **Principal Component Analysis - PCA**

- *t*-disributed Stochastic Neighbour Embedding (**t-SNE**)

# Dimensionality reduction on the Iris dataset

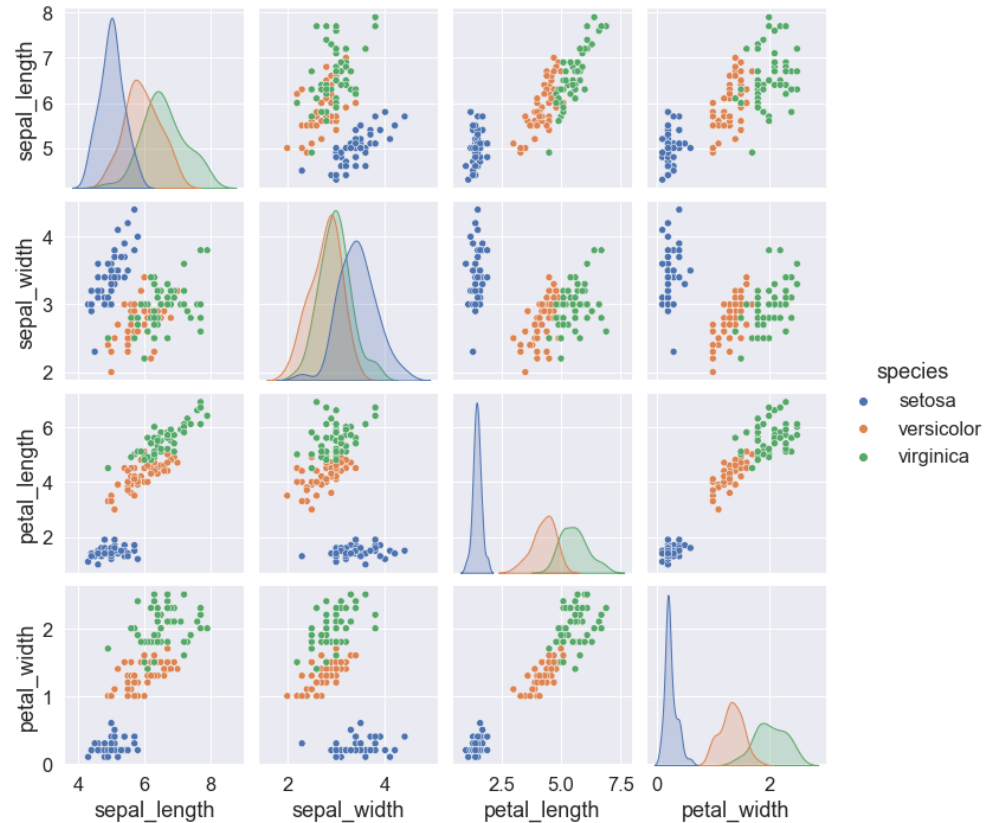The dataset consists of measurements on 150 Iris flowers from 3 species

with **4 Features**: ⟶ Visualisation in 4 dimensions difficult

- petal width (Kronblatt Breite)
- petal length (Kronblatt Länge)
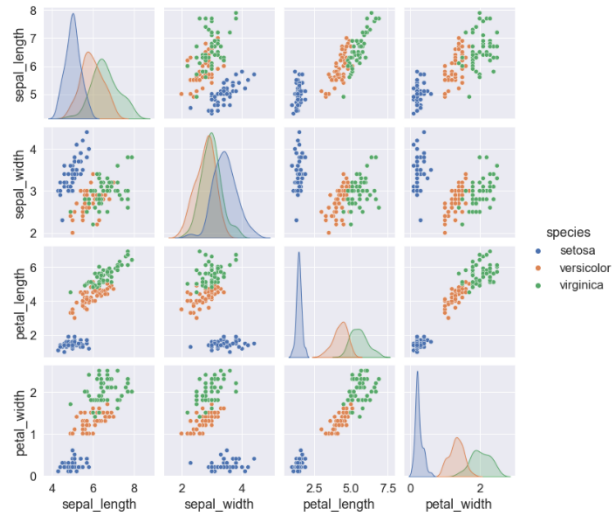- sepal width (Kelchblatt Breite)
- sepal length (Kelchblatt Länge

(https://de.wikipedia.org/wiki/Kronblatt)

| sepal_length | sepal_width | petal_length | petal_width | class |
|---|---|---|---|---|
| 5 | 3.3 | 1.4 | 0.2 | Iris-setosa |
| ... | | | | |
| 5.7 | 2.8 | 4.1 | 1.3 | Iris-versicolor |
| ... | | | | |
| 6.3 | 3.3 | 6 | 2.5 | Iris-virginica |
| ... | | | | |

# Iris dataset: Visualisations of the four features

# Dimensionality reduction on the Iris dataset

4 Dimensions:

Visualisation difficult



Dimensionality reduction

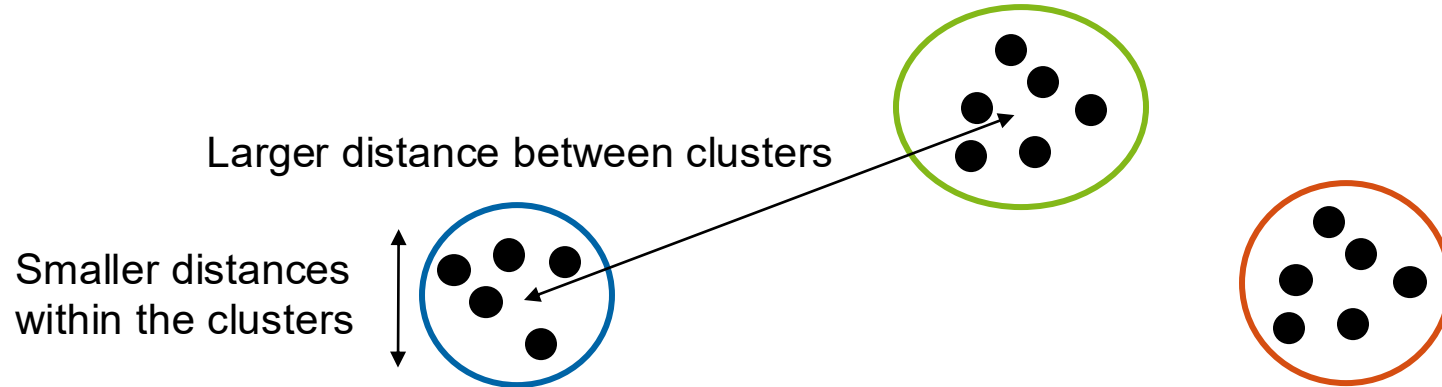2 new dimensions, that contain the «most information»:

# Clustering

Goal**: Identify subgroups** of datapoints that are more similar to each other than to the elements in other subgroups.
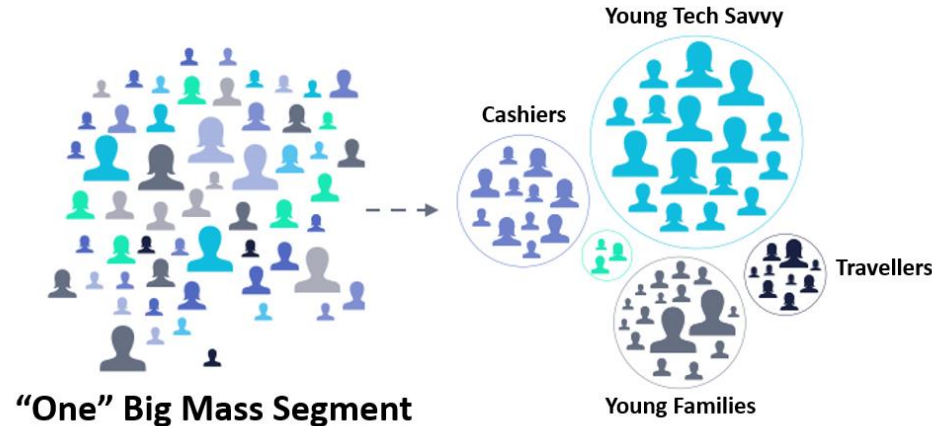
# Clustering

Goal**: Identify subgroups** of datapoints that are more similar to each other than to the elements in other subgroups.
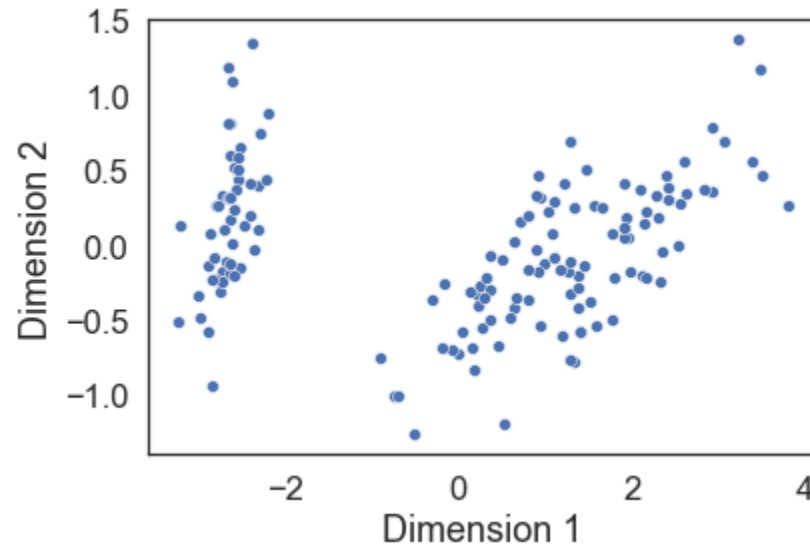
Larger distance between clusters

Smaller distances within the clusters

→ Needs metric to quantify similarities.

# Unsupervised machine learning

**Example: Clustering** is the task of ***grouping a set of objects*** in such a way that objects in the same group (called a cluster) are more similar (in some sense) to each other than to those in other groups (clusters)



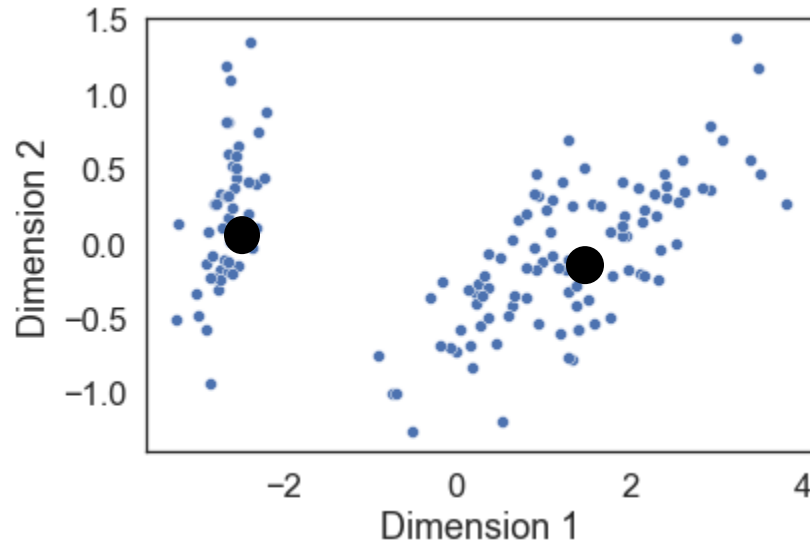Source: **https://www.smartera3s.com/products/customer-segmentation/**

# Example in 2D

# *K*-Means: A simple clustering method

Hyperparameter *k*: number of clusters to determine

*Means*: The centroids of the clusters

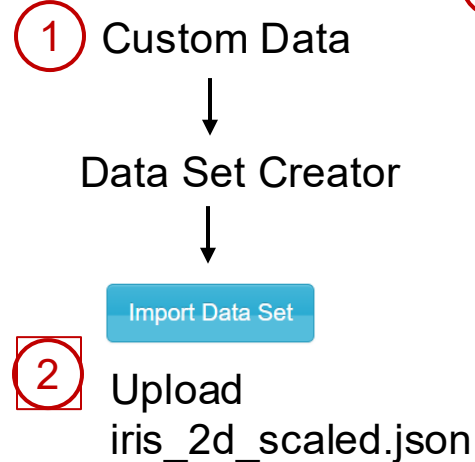Assumption: **spherical distribution** within clusters



Iterative algorithm – until stopping criterium is statisfied:
1.  Random **Initialisation** of the *Means*
2.  **Each datapoint is assigned** to closest *Mean*.
3.  **Recalculate the *Means*** from the newly assigned datapoints.
4.  **Repeat steps 2 and 3** until *Means* do not change anymore.

# *K*-Means on 2D Example-Dataset

## https://educlust.dbvis.de/



(1) Custom Data

↓

Data Set Creator

↓

Import Data Set

(2) Upload iris_2d_scaled.json

Download: https://tinyurl.com/5f9mkxmb

(3) Choose method: k-means

(4) Parameters: Choose *k*

(5) Set speed

(6) Start animation

# Structure of an unsupervised learning problem

**Training**

Training data

Features/
Input values

Model

Optimisation Algorithm

Learning

Trained Model

Evaluate

**Inference**

New data

Features/
Input values

Assignment
of new data

# Reinforcement-Learning



**"Play games without knowing the rules"**