# Introduction to Machine Learning

## Robert Haase, Maximilian Joas

CENTER FOR SCALABLE DATA ANALYTICS AND ARTIFICIAL INTELLIGENCE

Reusing materials from Johannes Soltwedel, Till Korten, Johannes Soltwedel, Laura Žigutytė (TU Dresden), Ryan Savill (MPI-CBG Dresden), Matthias Täschner (ScaDS.AI/Uni Leipzig) and the Scikit-learn community.
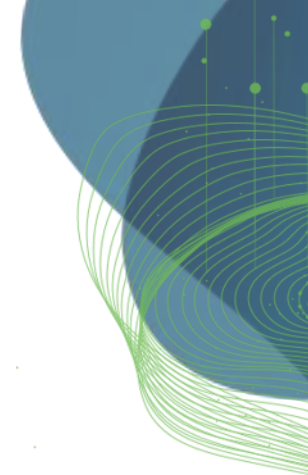
Funded by

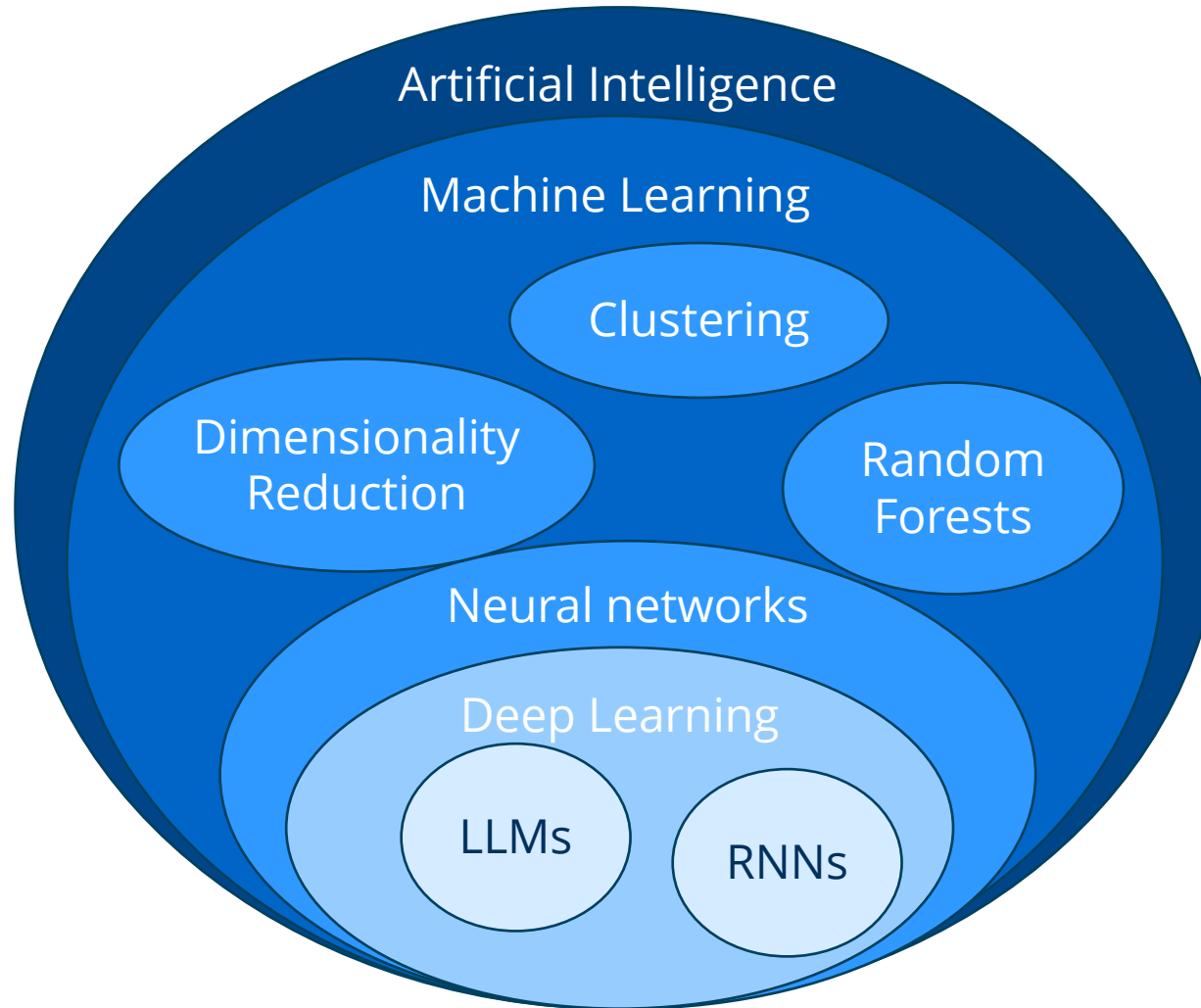Bundesministerium für Bildung und Forschung

SACHSEN

Diese Maßnahme wird gefördert durch die Bundesregierung aufgrund eines Beschlusses des Deutschen Bundestages. Diese Maßnahme wird mitfinanziert durch Steuermittel auf der Grundlage des von den Abgeordneten des Sächsischen Landtags beschlossenen Haushaltes.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

These slides can be reused under the terms of the CC-BY 4.0 license, unless mentioned otherwise.

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Artificial intelligence

# Artificial intelligence

Narrow AI
- Application specific
- Trained on labelled data
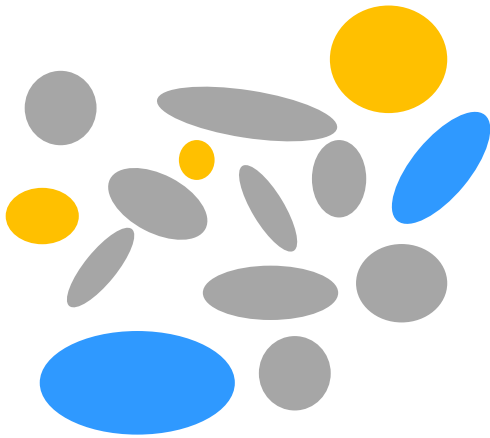- Reflexive tasks
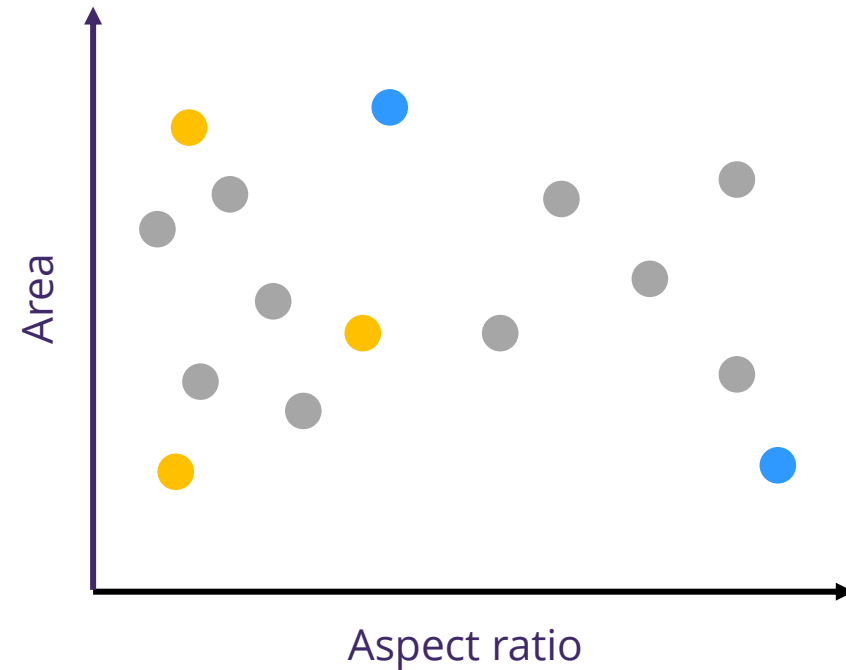- Cannot extrapolate

Great for data analysis tasks

General AI
- Human capabilities
- Access to knowledge of humanity, beyond individuals
- Can create *new* solutions by working creatively

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Labelled data

- E.g. for shape differentiation of objects

- Partially labelled data ◁ **Bias?**

Elongated
Round
Unlabelled

Area

Aspect ratio

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

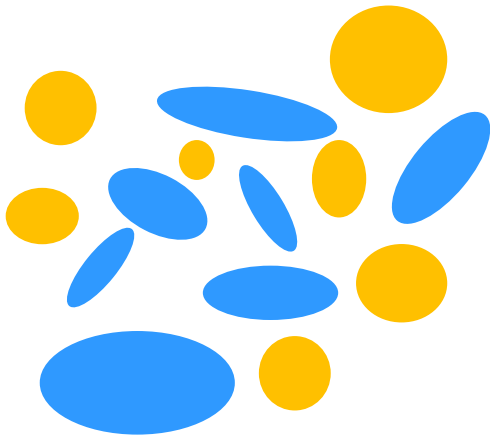ScaDS.AI
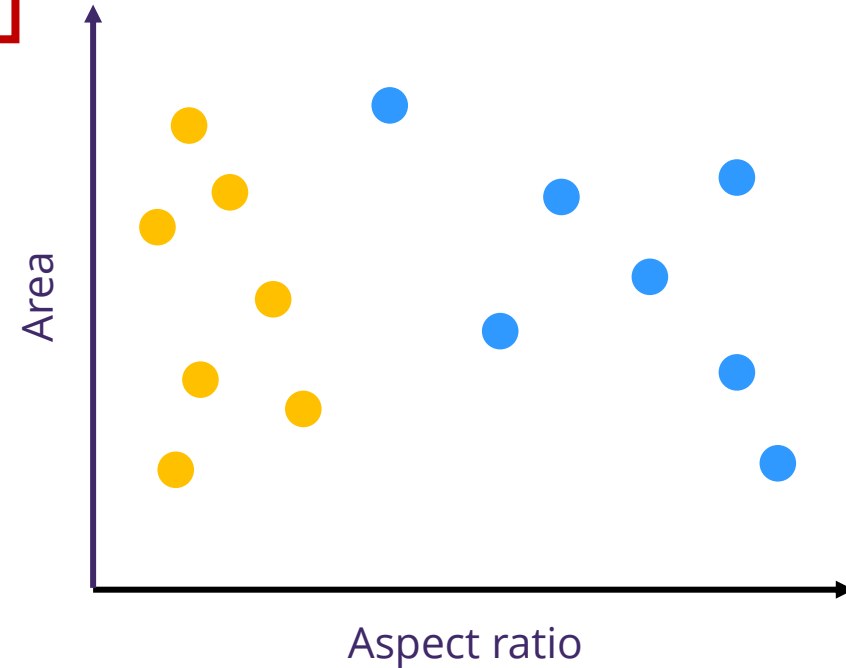DRESDEN LEIPZIG

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Labelled data

- E.g. for shape differentiation of objects
- Fully labelled data
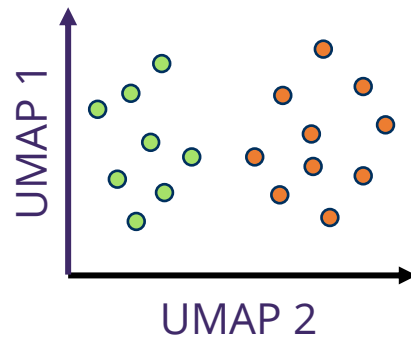
Typically expensive



Elongated
Round
Unlabelled

Area

Aspect ratio

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Artificial intelligence

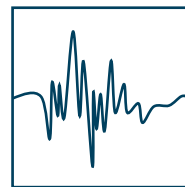| Explorative | → | Analytic | → | Generative |
|:-:|:-:|:-:|:-:|:-:|

## Unsupervised ML

- Dimensionality reduction
- Clustering
- Detecting patterns in unlabeled data
- Hypothesis generation

## Supervised ML

- Learning tasks otherwise only humans could do
- Train a model based on labeled data, predict a classification

## Generative AI

- Produces new data provided a context, often with human language prompts
- Hyped since 2022, with yet unclear limitations

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# <u>Unsupervised</u> Machine Learning

## Robert Haase

Reusing materials from Johannes Soltwedel, Till Korten, Johannes Müller, Laura Žiguty (TU Dresden), Ryan Savill (MPI-CBG), Matthias Täschner (ScaDS.AI/Uni Leipzig) and the Scikit-learn community.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Hypothesis-driven quantitative science

Hypothesis: The amplitude of a given signal is an indicator for upcoming earthquakes.

Null-Hypothesis: There is no relationship between the amplitude and future earthquakes.

Data download

> Shall we use a different dataset / sensor?

Data preprocessing

> Shall we use a different denoising algorithm?

> Shall we modify our measurement + hypothesis?

Amplitude measurement

Statistics

Reject / accept null-hypothesis

> Shall we use a different statistical test?

> Be careful going down this rabbit hole, you may be leaving good scientific practice behind.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 8

8

# Data-driven quantitative science

~~Hypothesis: The amplitude of a given signal is an indicator for upcoming earthquakes.~~

Question: Which measurement is a good predictor for upcoming earthquakes?

Which sensor / data is the most reliable?
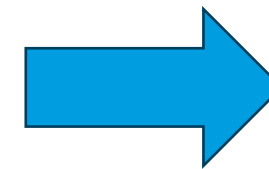
Data download (multiple sources, sensors, …)

Why?

Data preprocessing using Method A, B, C

Amplitude, frequency, wavelength, … measurement

Which parameter shows any relationship with upcoming earthquakes?

Statistics

Hypothesis generation

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Feature selection

- Which measurement / parameter / feature is related to the effect I'm investigating?

- Example goals:



- Amplitude
- Energy
- Duration
- …

- Silence
- Tourists jumping on a sensor
- Earthquake approaching

Signal classification

original     top_hat(10)

gaussian_sobel(1)     random

Pixel classification

- Area
- Perimeter
- Aspect ratio
- …

- Round
- Elongated

Object classification

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Feature selection

Question: Which features shall I analyse?

Challenges:

- Physical properties versus measurable features

- Correlation versus causation

- Too many features

If you have no idea -> unsupervised machine learning

- Dimensionality reduction

- Clustering

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Dimensionality reduction: Principal Component Analysis (PCA)

Linear transformation of high-dimensional data to concentrate information in a lower dimensional *embedding*

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slides adapted from Johannes Soltwedel, PoL TU-Dresden (licensed CC-BY 4.0)

# Embeddings

- N-dimensional latent space

- Axes typically have no meaningful/physical name (PCA1, UMAP1, …) and no physical unit

- Allow representing complex measurements, things, relationships in numeric space.

- Example:
  - You measure amplitude, frequency, wavelength, etc.,
  - derive a 2D-embedding from it,
  - to visualize the data or
  - to better process data

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 13

# Non-Euclidian spaces

## Not all features might be distances



Use travel time between W and S as metric for distance

→ Travelling from **W**ehlen to **S**trand by bike is probably faster if you make a detour through **R**athen

Slide 14

14

# Uniform Manifold Approximation Projection (UMAP)

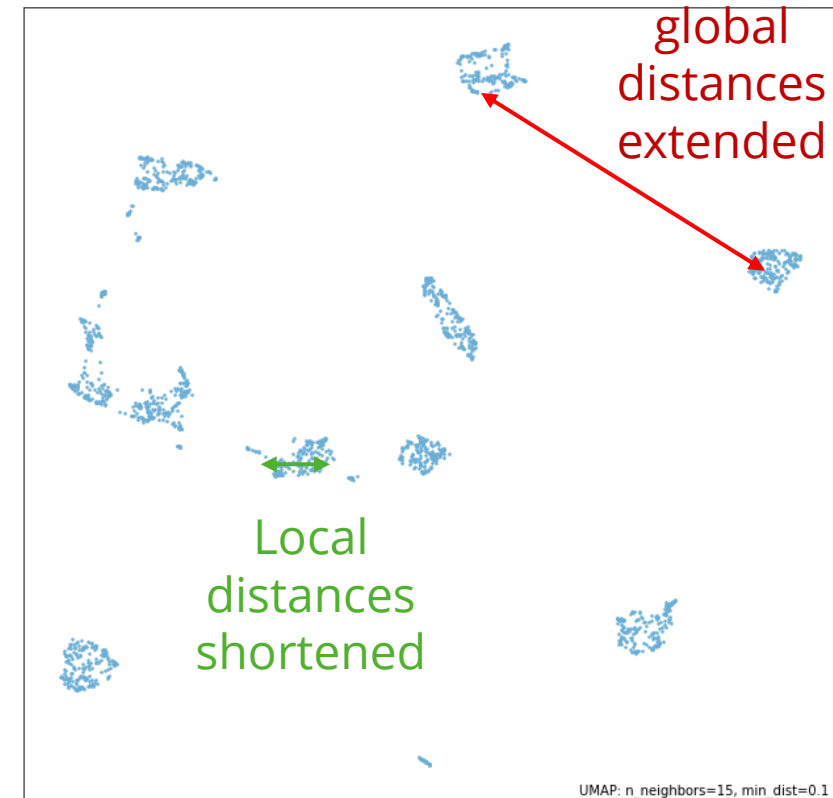Structural, hierarchical, <span style="color:red">non-linear</span> transformation

Modifies density of data points.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025
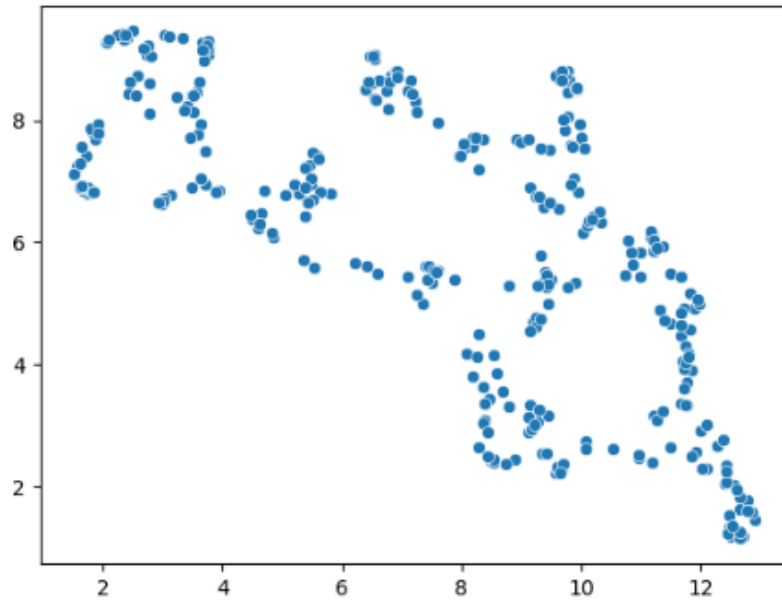
# Uniform Manifold Approximation Projection (UMAP)

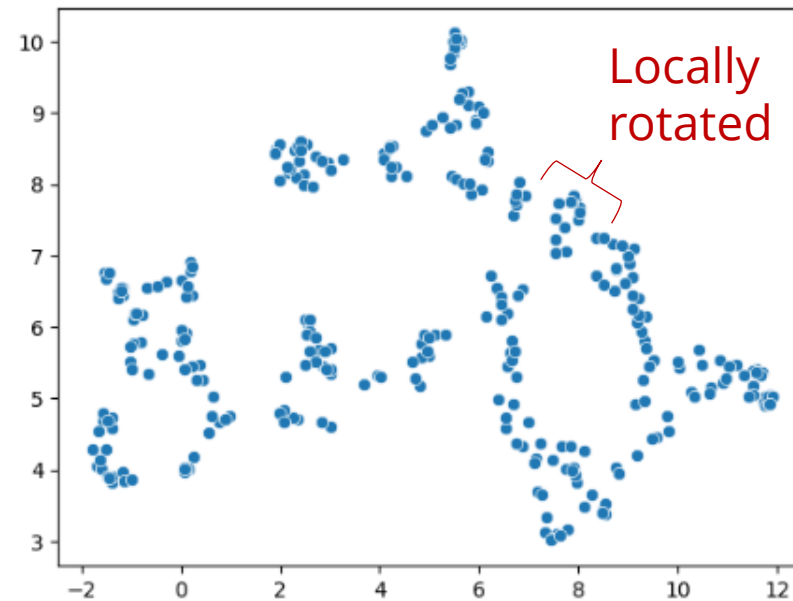Non-deterministic algorithm: You execute it twice, you get different results.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

https://haesleinhuepf.github.io/BioImageAnalysisNotebooks/47_clustering/umap.html?highlight=umap#a-note-on-repeatability

# Dimensionality reduction

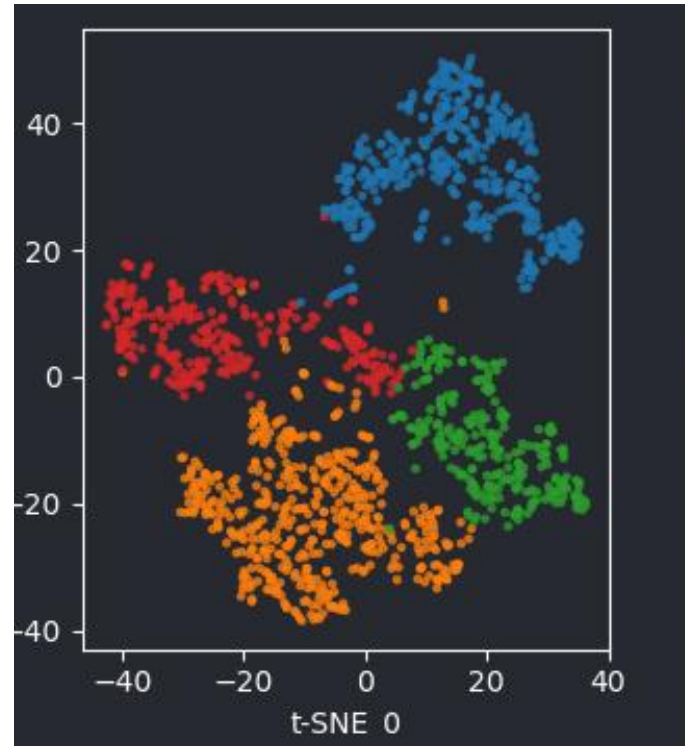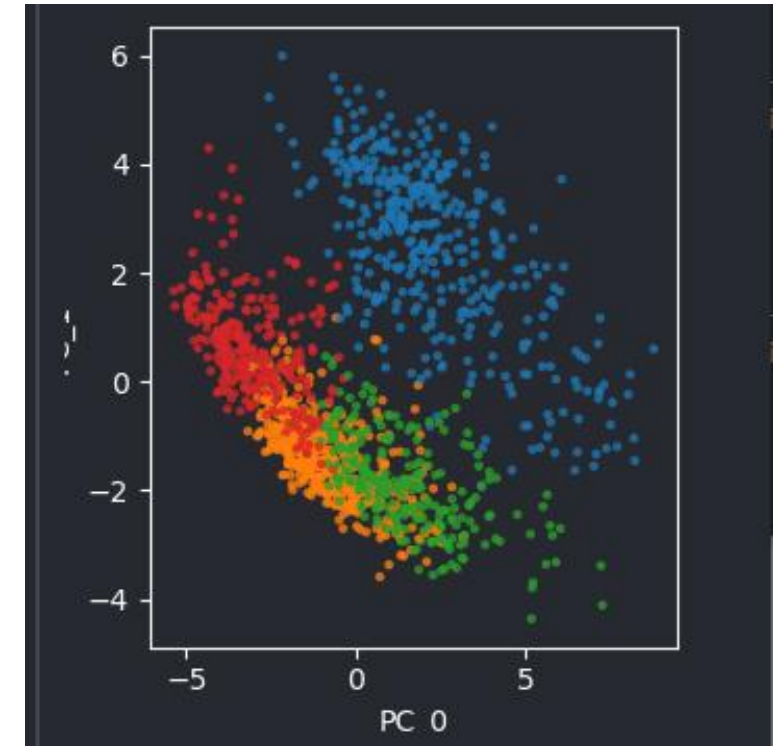Uniform manifold approximation and projection (UMAP)

t-distributed stochastic neighbor embedding (t-SNE)

Principal component analysis (PCA)

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 17

17

# Clustering

Unsupervised machine learning may include grouping objects without given ground truth

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Clustering

Unsupervised machine learning may include grouping objects without given ground truth



Round  Elongated

Names given by human observer *after* grouping / clustering

UMAP 1

frequency

UMAP 2

frequency

Robert Haase,
Max Joas
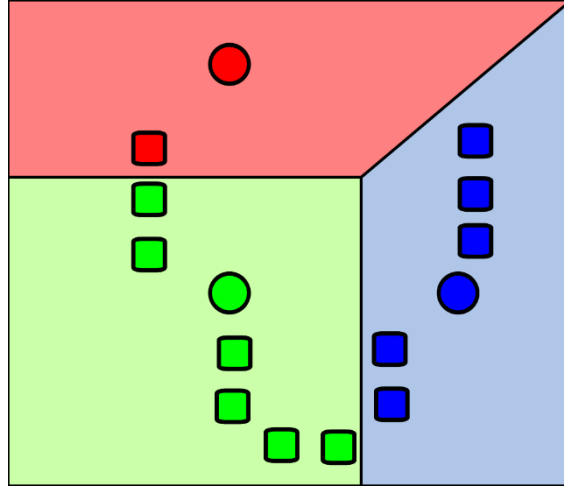Intro to ML & DL
AI4Seismology
May 5th 2025

# K-Means Clustering

Clustering algorithm, where you *only* need to specify the number of clusters.

Step1: Random initialization of cluster centers
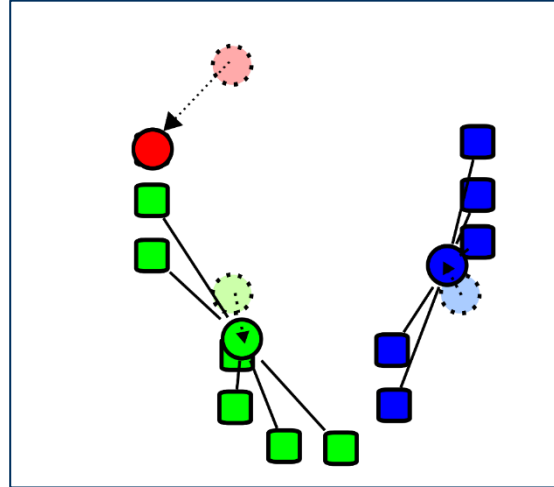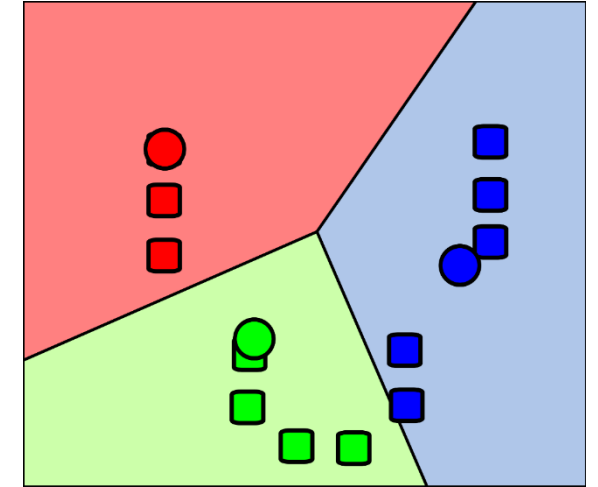
Step2: Tessellation of space into cluster regions

Step3: Replace cluster center with centroids

Step4: Repeat 2&3 until convergence
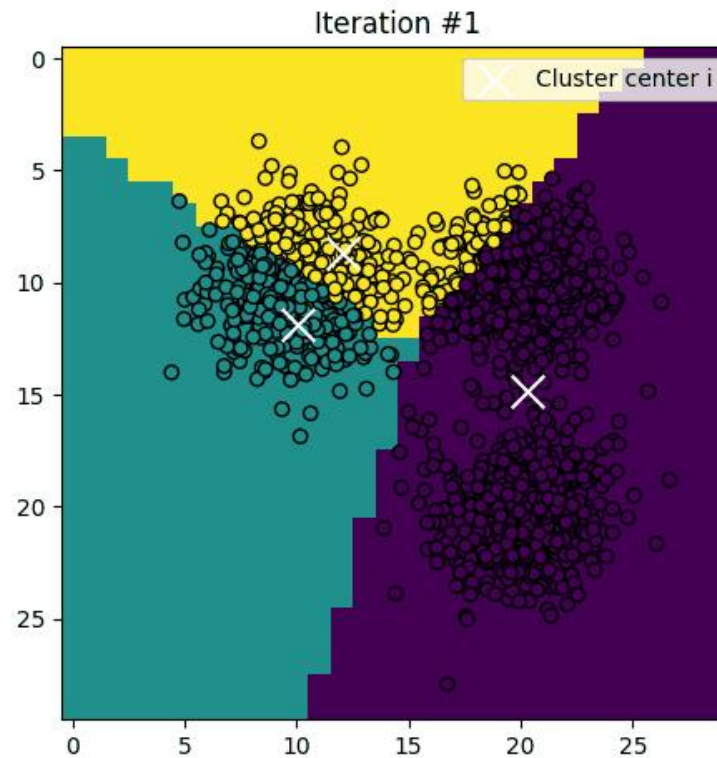
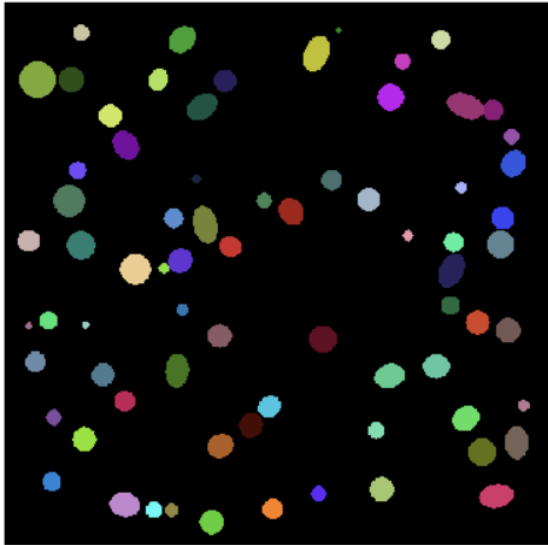Slide adapted from Johannes Soltwedel, TU Dresden
Licensed CC-BY 4.0

# K-Means Clustering

Clustering algorithm, where you *only* need to specify the number of clusters.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

Goal: Understand shape measurements

Data: Shape measurements from *randomly* shaped blobs.



| | label | area | perimeter | minor_axis_length | major_axis_length | circularity | solidity | aspect_ratio | elongation |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 97.0 | 32.970563 | 11.092860 | 11.092860 | 1.121318 | 0.788288 | 1.000000 | 0.000000 |
| 1 | 2 | 285.0 | 60.284271 | 19.052651 | 19.052651 | 0.985477 | 0.785116 | 1.000000 | 0.000000 |
| 2 | 3 | 473.0 | 79.597980 | 21.823280 | 27.594586 | 0.938138 | 0.785448 | 1.264456 | 0.209146 |
| 3 | 4 | 321.0 | 63.112698 | 19.033334 | 21.456036 | 1.012701 | 0.786033 | 1.127287 | 0.112915 |
| 4 | 5 | 407.0 | 72.769553 | 22.155138 | 23.384406 | 0.965839 | 0.785586 | 1.055485 | 0.052568 |

...

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
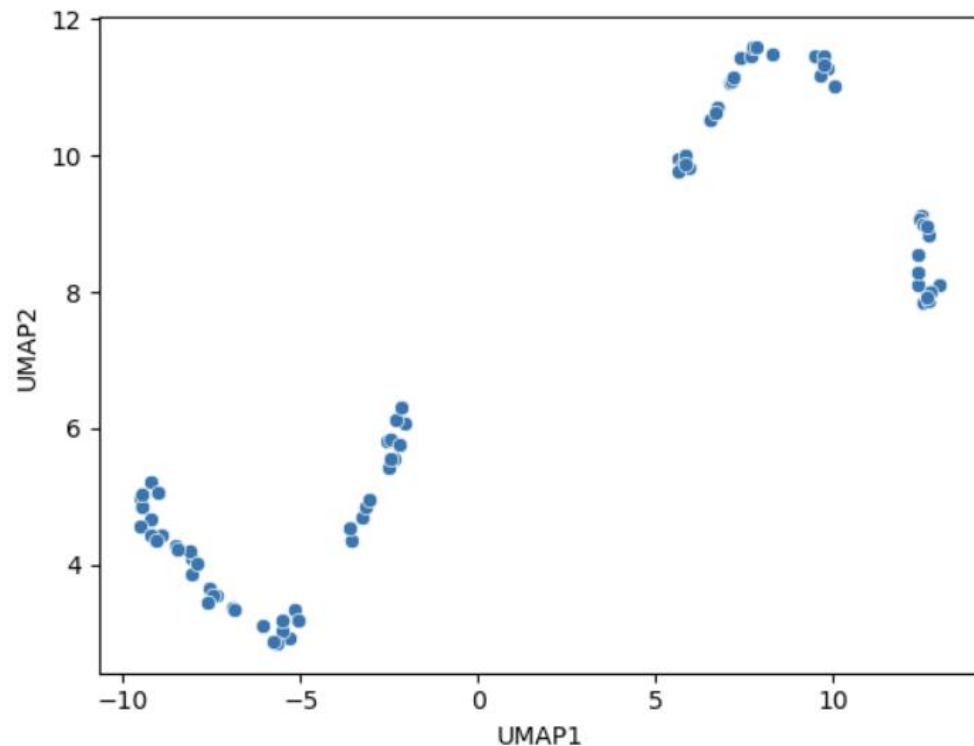May 5th 2025

# Walk-through: Data Exploration

Step 1: Dimensionality reduction (UMAP)

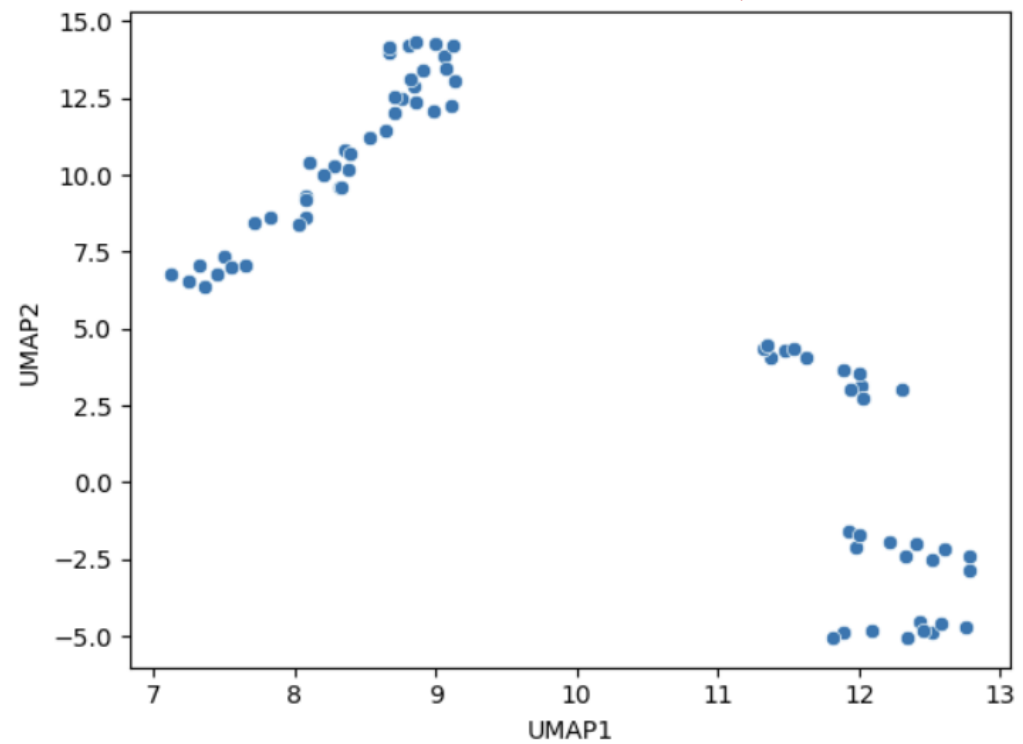Observation: There appear to be *2 distinct groups*

Pinning the random seed is no solution to this general problem.

Beware: UMAPs are non-deterministic. Different runs lead to different results.

Run 1:

Run 2:

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

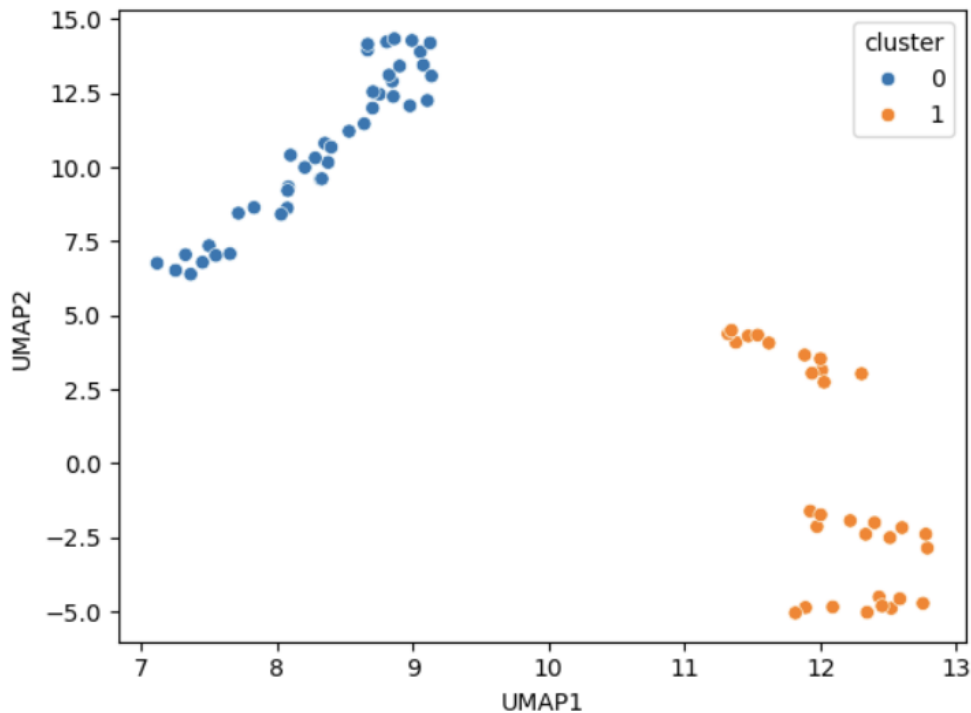# Walk-through: Data Exploration
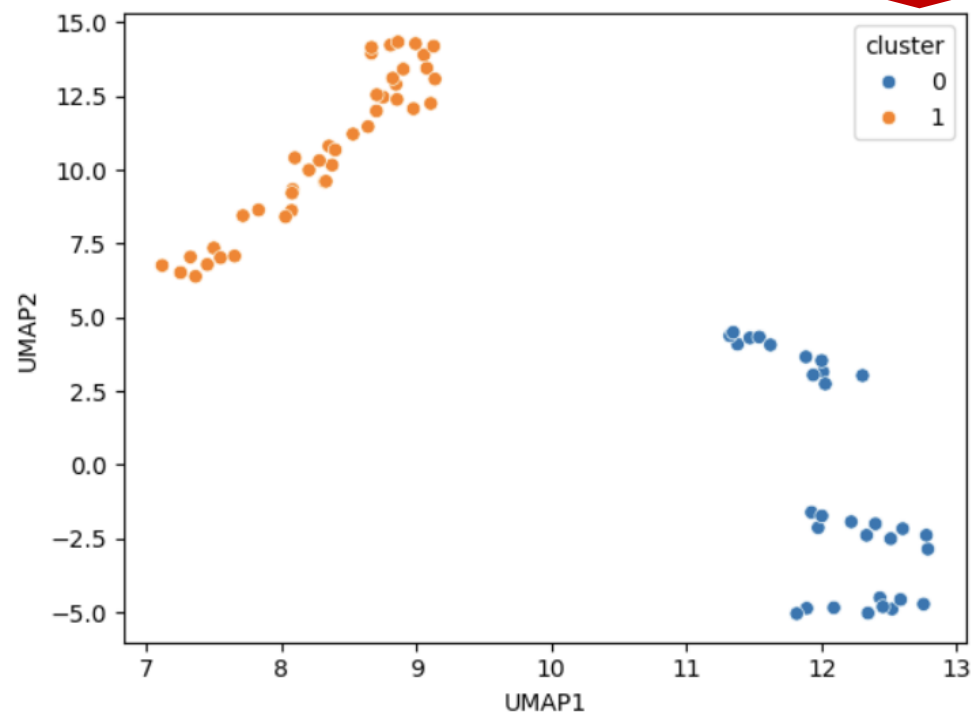
Step 2: Clustering data into 2 clusters

Using K-Means clustering

Pinning the random seed is no solution to this general problem.

Beware: Clustering-algorithms are non-deterministic. Different runs lead to different results.

Run 1:



Run 2:

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

Step 3: Feature selection

Based on correlation
with distance to cluster-centers

Hypothesis:
"Circularity and minor_axis_length
allow to predict round vs.
elongated classification."

*Hypothesis generation*

*Side note: beware of feature correlation.*

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Walk-through: Data Exploration

Step 4: Train a classifier (supervised ML)

Goal: Eliminate non-determinism



Clustering result (non-deterministic)

Classification result (deterministic, repeatable)

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Supervised Machine Learning

## Robert Haase

Funded by

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

CENTER FOR SCALABLE DATA ANALYTICS AND ARTIFICIAL INTELLIGENCE

# Supervised Machine learning

Automatic construction of predictive models from given data

Annotated raw data, often generated by humans

Pixels,    Objects,    Images, Audio, Sensor data, Text, Measurements, …

Classification (categorical)

Cat ~~Dog~~

Earth-quake ~~Wind~~

Regression (continuous numerical)

n = 11

green_magenta_ratio=0.3

$P_{Cat}$= 0.5
$P_{Microscope}$= 0.4

Height = 80 cm

Raw data

Training

Ground truth

Prediction

Model

Classification / regression

Quality

Accuracy, Precision, Recall, …

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Goal

Guess classification (color) from position of a sample in parameter space.



Input data     Decision Tree .95     Random Forest .93     Neural Net .90

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Adapted from https://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html
© 2007 - 2019, scikit-learn developers (BSD License).

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Machine learning for image segmentation

*Supervised* machine learning: We give the computer some ground truth to learn from

The computer derives a *model* or a *classifier* which can judge if a pixel should be foreground (white) or background (black)

Example: Binary classifier

Training



Model / classifier

?

Raw image

Binary image

Robert Haase
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

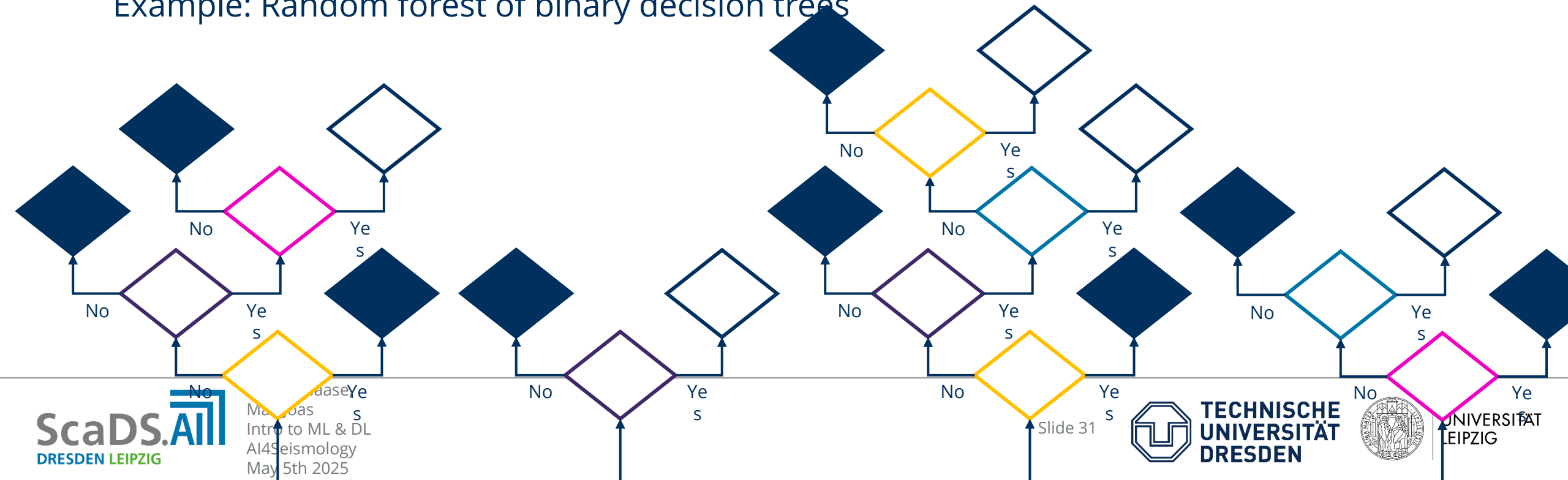# Random forest based image segmentation

Decision trees are classifiers, they decide if a pixel should be white or black

Random decision trees are randomly initialized, afterwards evaluated and selected
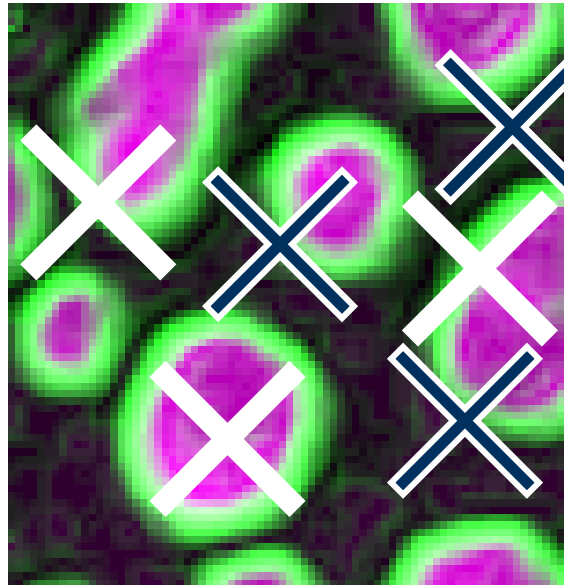
Random forests consist of many random decision trees

Example: Random forest of binary decision trees

Matthias
Intro to ML & DL
AI4Seismology
May 5th 2025

# Deriving random decision trees

For efficient processing, we randomly *sample* our data set
- Individual pixels, their intensity and their classification



Threshold

$X_2$

$X_1$

Note: You cannot use a single threshold to make the decision correctly

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 32

# Deriving random decision trees

Decision trees combine several thresholds on several parameters

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Deriving random decision trees

Depending on sampling, the decision trees are different

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Random Forest Pixel Classifiers

Combination of individual tree decisions by voting or max / mean

Prediction

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Random Forest Pixel Classifiers

Typical numbers for pixel classifiers in microscopy

Available features:

- Gaussian blur image
- DoG image
- LoG image
- Hessian
- ....

Depth: 4

Number of trees: > 100

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Model validation

In order to assess model quality, we split the ground truth into two set
- Training set (50%-90% of the available data)
- Test set (10%-50% of the available data)

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 38

# Model validation

Based on the theory of sets



Legend:
- A — Prediction
- B — Reference / ground truth
- ROI — Region of interest
- TP — True-positive
- FN — False-negative
- FP — False-positive
- TN — True-negative

$$\text{Accuracy} = \frac{\text{correct classifications}}{\text{all classifications}}$$

This means: $= \dfrac{TP + TN}{FP + FN + TP + TN}$

$$\text{Precision} = \frac{\text{Relevant retrieved instances}}{\text{All } \textbf{retrieved} \text{ instances}}$$

This *may* mean: $= \dfrac{TP}{FP + TP}$

# Model validation: Accuracy versus precision



Accurate and precise

Accurate and but not precise

Not accurate and but precise

Neither accurate nor precise

Lesson learned:
A single quality metric cannot describe the whole situation

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Model validation: Accuracy versus Jaccard Index

## Side-effect of number of true negatives

Nuclei

Reference

Segmentation result



$$A = \frac{TP + TN}{FN + FP + TP + TN}$$

$$J = \frac{TP}{FN + FP + TP}$$

Accuracy: 0.97
Jaccard Index: 0.73

Accuracy decreases because there are less correct black pixels (TN)

Accuracy: 0.95
Jaccard Index: 0.73

# Explainable AI

Depending on the target group [for the explanation], the influence of data is more important than how AI algorithms work.

- Many computer scientists want to explain and understand AI methods.

- Geoscientists use AI as a method to explain geological processes.

- Example: "What parameters distinguish round objects from elongated ones?"

https://haesleinhuepf.github.io/xai/30_shap/object_classification.html

# Pitfall: Correlation

Correlated features may harm interpretability



Features may appear less valuable.

# Deep Learning

## Robert Haase

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Machine learning for image analysis

In classical machine learning, we typically select features for training our classifier

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Deep learning for image analysis

In deep learning, this selection becomes part of the black box



Neural networks

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 55

# Neural networks

- How biologists see neurons

- How computer scientists see neurons

"perceptron"

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Neural Networks

- Early form: "Multilayer Perceptron"
- fully connected class of feedforward artificial neural network

If there are *many* hidden layers, we speak of a *deep* neural network

Input layer          n hidden layer(s)          Output layer



Feed forward network

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

https://en.wikipedia.org/wiki/Multilayer_perceptron

# Convolutional neural networks

- Layer types

**Fully connected layer**

**Convolutional layer**

**Pooling layer**
**("Max pool", "Average pool")**

Field
of View
(FoV)



| 3 | 15 | 1 | 13 |
|---|----|---|----|
| 9 | 7  | 0 | 10 |
| 11| 5  | 5 | 3  |
| 1 | 8  | 9 | 6  |

Max pooling →

| 15 | 13 |
|----|----|
| 11 | 9  |

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Activation functions

- Introduction of *non-linearity* and *activation functions* enabled what we call *deep-learning* today.



$$y = f(w_1 x_1 + w_2 x_2 + w_3 x_3 + b)$$

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Convolutional neural networks

- Assuming we had no activation functions in the network layers can be reduced by eliminating brackets!

input

$x_1$  $x_2$  $x_3$

$w_1$  $w_2$  $w_4$

$w_3$

$w_6$

$w_5$

output  $y$

$$y = w_5(w_1 x_1 + w_2 x_2) + w_6(w_3 x_2 + w_4 x_3)$$

$$y = w_5 w_1 x_1 + w_5 w_2 x_2 + w_6 w_3 x_2 + w_6 w_4 x_3$$

$$y = w_5 w_1 x_1 +$$

$$v_1 = w_5 w_1$$
$$v_2 = w_5 w_2 + w_6 w_3$$
$$v_3 = w_6 w_4$$

$$y = v_1 x_1 + v_2 x_2 + v_3 x_3$$

$x_1$  $x_2$  $x_3$

$v_1$  $v_2$  $v_3$

$y$

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 60

# Activation functions

- Introduction of *non-linearity* and *activation functions* enabled what we call *deep-learning* today.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Source: https://en.wikipedia.org/wiki/Activation_function
Licensed CC-BY-SA 4.0 by Laughsinthestocks

# Learning: Back propagation

- Step 0: Initialize the network randomly (weights, bias)

- Step 1: Forward pass the input through the network, get an initial prediction

- Step 2: Compare the output with the ground truth, compute the error (loss function)
  - The loss function can be freely defined.
  - Example: mean squared error

- Step 3: Update weights



- Silence
- Tourists jumping on a sensor
- Earthquake approaching

Prediction: 0.3   0.4   0.4

Ground truth: 0   0   1

Loss   0.18

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

ScaDS.AI
DRESDEN LEIPZIG

# Learning: Back propagation

- Updating weights:
  - Set output to the error (per-parameter gradient)
  - Backward-pass: add/subtract gradients from weights, to push the network towards giving the right answer.
- Execute the same procedure for next sample
- Execute the same for multiple *epochs*



- Silence
- Tourists jumping on a sensor
- Earthquake approaching

Set: 0.3  0.3  -0.6

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Training NNs: Batch size & epochs

Problem:

* Assume you have $10^{10}$ samples and attempt to train for 1000 epochs

-> $10^{13}$ backprop steps required.

Data

epoch 1

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Training NNs: Batch size & epochs

Problem:

- Assume you have $10^{10}$ samples and attempt to train for 1000 epochs

-> $10^{13}$ backprop. steps required.

Solution:

- Draw n=1000 random samples from the training data to train for one epoch.

- Next epoch: different n samples.

-> $10^6$ backprop. steps required.



Data

Batch size

epoch 1

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Training NNs: Batch size & epochs

Problem:

- Assume you have $10^{10}$ samples and attempt to train for 1000 epochs
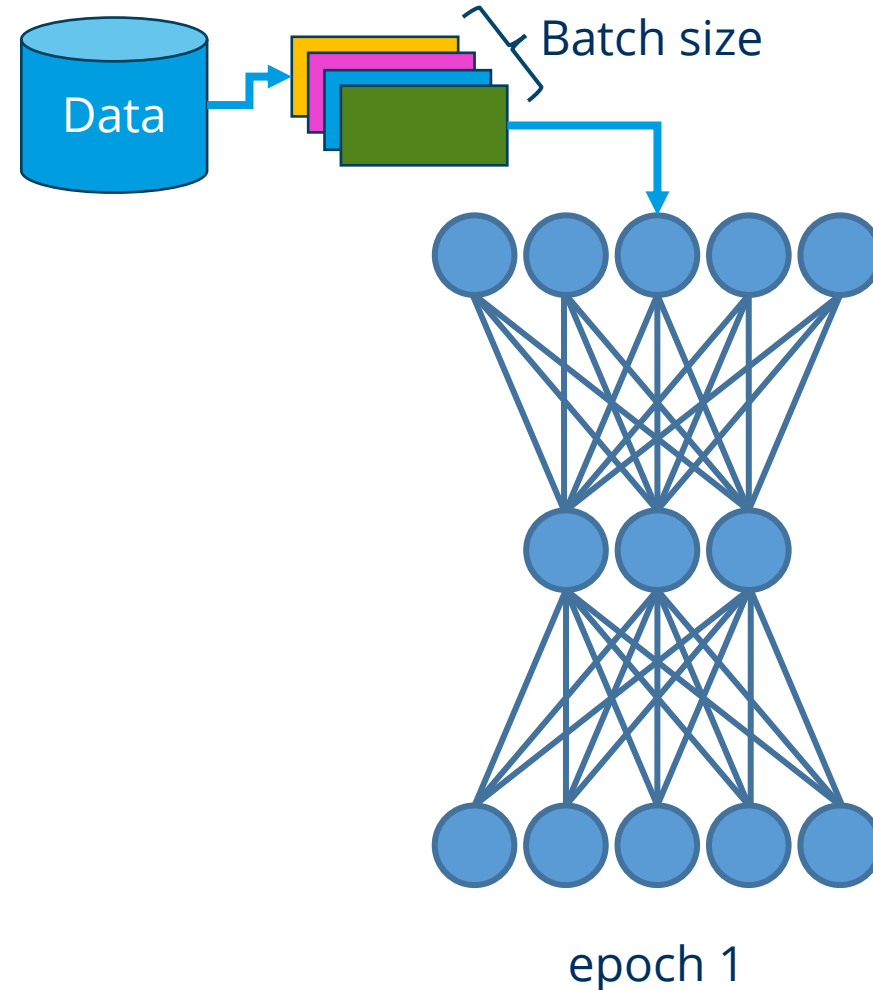
-> $10^{13}$ backprop steps required.

Solution:

- Draw n=1000 random samples from the training data to train for one epoch.

- Next epoch: different n samples.

-> $10^6$ backprop steps required.

Batch size

Data

epoch 2

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# Training NNs: Drop-out

- Drop-out: deactivating individual neurons during training

- Helps with over-fitting, because the network cannot rely on individual neurons by chance being well trained, while others remain randomly initialized

- Example: drop-out-rate: 30%



epoch 1

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Training NNs: Drop-out

- Drop-out: deactivating individual neurons during training

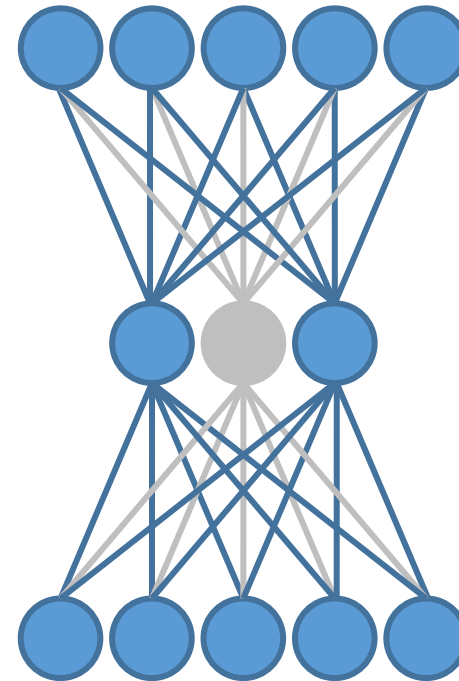- Helps with over-fitting, because the network cannot rely on individual neurons by chance being well trained, while others remain randomly initialized

- Example: drop-out-rate: 30%



epoch 2

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Active Learning

## Maximilian Joas

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Train- Validation- Test-split

Training dataset (~80% of the data)

Used for training directly.

Validation dataset (~10% of the data)

Used to tune parameters, select features, and make other architecture decisions (also called **Dev set).**

Test dataset (~10% of the data)

Final evaluation after training is finished (once).



Underfitting

Overfitting

Loss (lower is better)

Training duration (epochs)

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

https://cs230.stanford.edu/blog/split/

Slide 71

# Loss Curve Analysis

Underfitting

Overfitting

Loss (lower is better)

Training duration (epochs)

**Questions answered**:

- Is my model converging?
- Is the learning rate appropriate?
- Am I training for the right number of epochs?
- When should I apply early stopping?

**Outcome: Helps you fine-tune training hyperparameters.**

Robert Haase,
Robert Haase,
Max Joas
@haesleinhuepf
Intro to ML & DL
AI4Seismology
AI4Seismology
May 5th 2025
May 5th 2025

# Learning Curve Analysis



**High bias, high variance:**
**=> More complex model to reduce bias.**

**Low bias, high variance:**
**=> Collect more data.**

Ng, A. (2018). *Machine Learning Yearning* (Draft). DeepLearning.AI.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# Active learning allows to train better models with less labeled data



m (training set size)

Dev Error
Training Error
Desired Performance
Error

**Uncertainty sampling** based on softmax output

**Diversity sampling** based on neuron activations

**Cluster sampling** based on input data similarity

Robert Haase,
s
ML & DL
nology
May 5th 2025

UNIVERSITÄT
PZIG

# Uncertainty sampling - an intuitive explanation



$$E(p) = - \Sigma\ p_i \log p_i$$

```
[0.26, 0.23, 0.28, 0.23]
[0.24, 0.27, 0.22, 0.27]
[0.29, 0.21, 0.25, 0.25]
[0.22, 0.26, 0.24, 0.28]
```

```
[0.92, 0.03, 0.03, 0.02]
[0.01, 0.94, 0.02, 0.03]
[0.03, 0.01, 0.95, 0.01]
[0.02, 0.04, 0.01, 0.93]
```

Robert Haase,
Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
AI4Seismology
May 5th 2025
May 5th 2025
@haesleinhuepf

# Neural Network Architectures

## Robert Haase

TECHNISCHE
UNIVERSITÄT
DRESDEN

UNIVERSITÄT
LEIPZIG

# NN Architectures: Recurrent Neural Networks

Introducing some form of memory through additional connections and nodes.



Hidden layer with
self-feedback

Hidden layer with
context nodes
(Elman network, 1990)

Fully
Recurrent NN

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Figure sources: Staudemeyer & Rothstein Morris 2.0
https://arxiv.org/pdf/1909.09586

# Training Recurrent Neural Networks

- Backpropagation through time

- Computationally expensive

- Unfolding through time

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

https://en.wikipedia.org/wiki/Backpropagation_through_time

# NN Architectures: Long Short-Term Memory (LSTM)

Differentiation between updating short-term memory (all the time) and updating long-term memory ([not] forgetting) thanks to separate input- and forget-gates.

# Traditional architecture: Encoder-Decoder Networks

Related: „Auto-encoder", „Variational Auto-Encoder", „U-Net"



Input: noisy image

Encoder

Decoder

Skip-connections (optional)

"Bottleneck", "Embedding", "Latent space"

Input: denoised image

# Traditional architecture: Encoder-Decoder Networks



"Embedding"

Figure source: Ronneberger et al 2015
https://arxiv.org/pdf/1505.04597

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# NN Architectures: Transformers



Decoder

Encoder

Vectors and matrices

Attention

Embedding

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Source: Vaswani et al (2017)
https://arxiv.org/abs/1706.03762

Slide 82

# Scaled dot-product attention

Attention score: How much related are two words?

Query: For which word are we calculating attention?

Key: To which word are we calculating attention

Value: Relevance of the query-key relationship

The cat is black and white.

Relevance value: 0.1

attention score

The cat is meowing.

Relevance value: 0.9

attention score

Scaled Dot-Product Attention

# Multi-head attention

## Multiple aspects represented by multiple attention heads

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 84

# NN Architectures: Transformers

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Source: Vaswani et al (2017)
https://arxiv.org/abs/1706.03762

# NN Architectures: Transformers

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# NN Architectures: Transformers

Related terms:

- Generative Pretrained Transformer (GPT)
- Large Language Models
- Next word-prediction

Decoder

Output

„[...] Microscope"

Encoder

„Die Katze sitzt neben dem Mikroskop."

**1**

Input

Masking

**2**

**3**

„The cat sits next to the [...]"

[Shifted] Output

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Slide 87

ScaDS.AI
DRESDEN LEIPZIG

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# NN Architectures: Vision Language Models

VLMs use combinations of traditional neural network architectures and transformers.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# NN Architectures: DNA Language Models

DNA-LMs use a variation of the transformer architecture.



Encoder

Decoder

TAGA GAC GAGG

TAGA GAC GAGG

Input embedding

Positional encoding

Multihead attention

LayerNorm

12×

Feed forward

LayerNorm

"Embedding", "Latent space"

Output

Masking

Output

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Figure cropped from Sanabria et al (2024), licensed
CC-BY 4.0
https://www.nature.com/articles/s42256-024-00872-0

TECHNISCHE UNIVERSITÄT DRESDEN

UNIVERSITÄT LEIPZIG

# Multi-modal Language Models

MMLMs use combinations and/or variations of traditional neural network architectures and transformers.

Robert Haase,
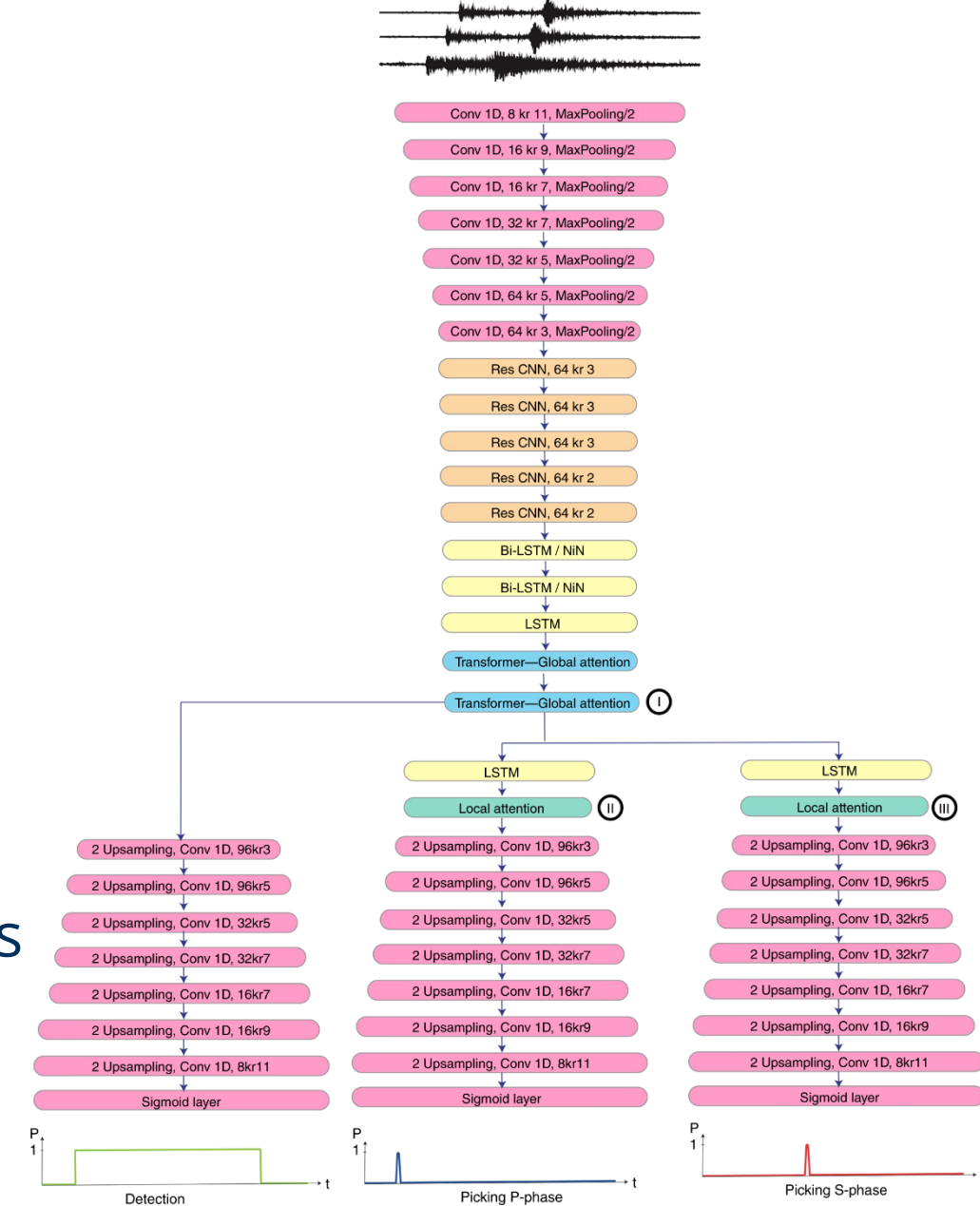Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

# NN Architectures

Modern NN architectures combine techniques quite freely. Example, for large earthquake detection:

- **LSTMs**

- **Transformers**

- **Convolutional**

- **Attention**

Combining architectures sometimes appears *more art than science*. Computer scientists world-wide struggle comparing different architectures.

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025

Figure cropped from Mousavi et al, licensed CC-BY 4.0
https://www.nature.com/articles/s41467-020-17591-w

# Summary

Unsupervised ML: Explorative data science , Embeddings

Supervised ML / DL: Preduction: classification / regression , Embeddings

Explainability: SHapleys Additive exPlanations (SHAP-Analysis)

Neural networks

- Many hidden layers -> *deep* learning, Embeddings

- Training: Drop-out, batch-size, epochs, active learning

- RNNs / LSTMs -> Memory

- Transformers -> Attention, Embeddings

Good scientific practice

- Train-test-split

- Overfitting / underfitting

Robert Haase,
Max Joas
Intro to ML & DL
AI4Seismology
May 5th 2025