

Data Analytics and Artificial Intelligence

Summer School 2025

Dr. Manuel Dömer
Programme Director BSc. Data Science

doem@zhaw.ch
+41 (0) 58 934 43 30

30.06.2025

Manuel Dömer

Zürcher Hochschule
für Angewandte Wissenschaften



- Programme Director BSc. Data Science
- Co-Head ZHAW Datalab
- EELISA ZHAW Academic Coordinator

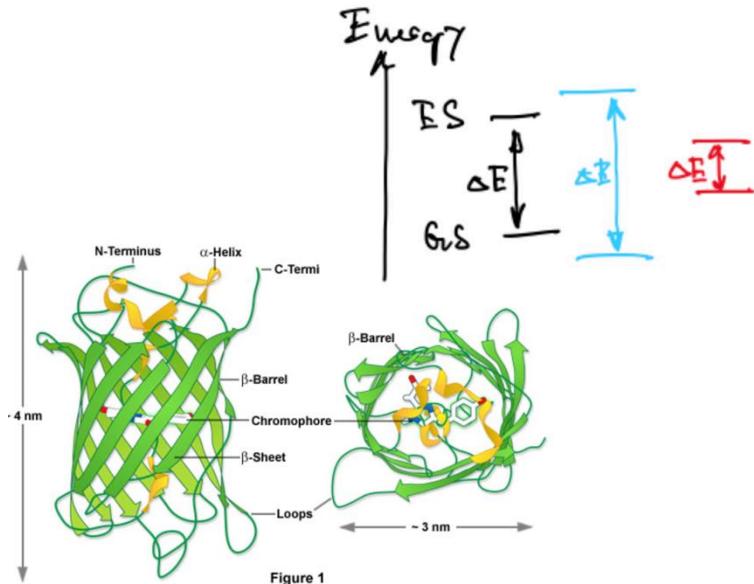
- Lead Modelling & Analytics in Digital Markets
- Data Scientist
- Credit Portfolio Modelling
- Natural Language Processing
- Junior Corporate Actuary
- Lecturer in Data Science
Lucerne University of Applied Sciences and Arts
- Co-Organiser
NLP Zürich-Meetup



- PhD in Computational Chemistry
- BSc and MSc in Chemistry



Research Interests



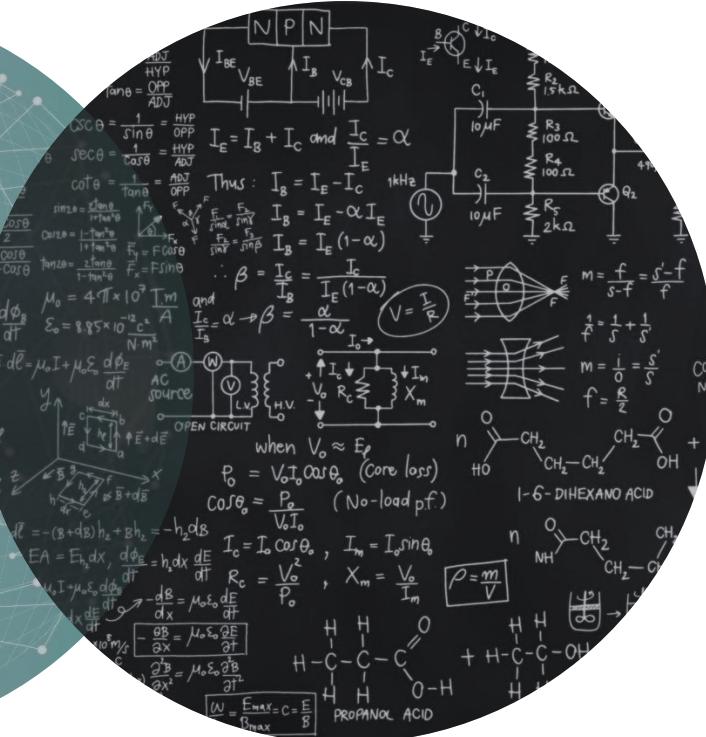
Schedule

	Mo, June 30	Tue, July 1	Wed, July 2	Thur, July 3
9-12 am	<ul style="list-style-type: none">• Introduction• Data Science-Process	<ul style="list-style-type: none">• Data Engineering Part 1	<ul style="list-style-type: none">• DE Part 2	<ul style="list-style-type: none">• DE Part 3
Lunch				
1-5 pm	<ul style="list-style-type: none">• Projects	<ul style="list-style-type: none">• Projects	<ul style="list-style-type: none">• Projects	<ul style="list-style-type: none">• Project Plan Presentations

	Mo, July 7	Tue, July 8	Wed, July 9	Thur, July 10
9-12 am	<ul style="list-style-type: none">• Supervised Machine Learning	<ul style="list-style-type: none">• Association Rules• Recommender Systems	<ul style="list-style-type: none">• Projects	<ul style="list-style-type: none">• Projects
Lunch				
1-5 pm	<ul style="list-style-type: none">• Unsupervised ML	<ul style="list-style-type: none">• Projects	<ul style="list-style-type: none">• Projects	<ul style="list-style-type: none">• Project Status Updates

Data Science

Data



Zurich UAS Racing – Formula Student

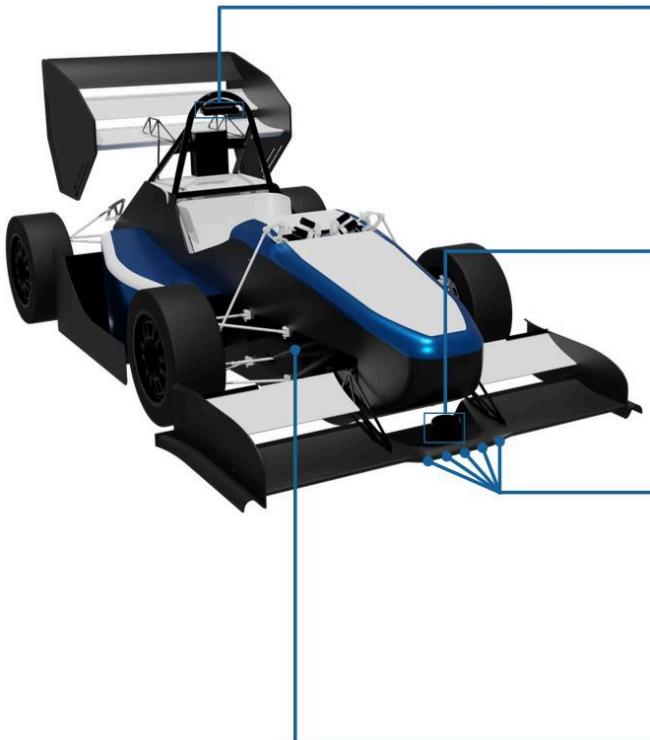
Zürcher Hochschule
für Angewandte Wissenschaften



<https://zurichuasracing.ch>



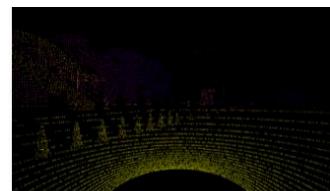
Sensor Data of an Autonomous Vehicle



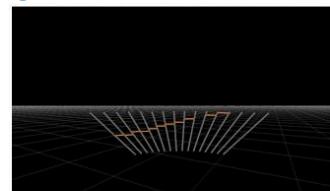
Stereo camera with depth perception



LIDAR sensors



Ultrasonic sensors



Position and acceleration sensors

Data in the context of «Data Science»

Zürcher Hochschule
für Angewandte Wissenschaften



**Data implies digital data in binary form,
which is stored and processed by electronic computer systems
in form of files of given formats that are interpretable by specific software
applications, e.g. as text, numbers, tables, images sound or executable programs...**

Structured vs. unstructured Data

Structured data

- pre-defined data model
- different datasets can follow the same data model

Examples

- tabular data

File formats

- csv
- ods
- xlsx
- HDF (Hierarchical Data Format)
- Apache Parquet

Unstructured data

- not a fixed data model
- no explicit structure, but might feature implicit structure, such as grammatical rules in text

Examples and typical file formats

Content	Formats
Text	txt, odt, docx, html
Images	jpeg, png
Audio	wave, aiff, flac, mp3, aac
Video	mp4, mov, wmv

Structured vs. unstructured Data

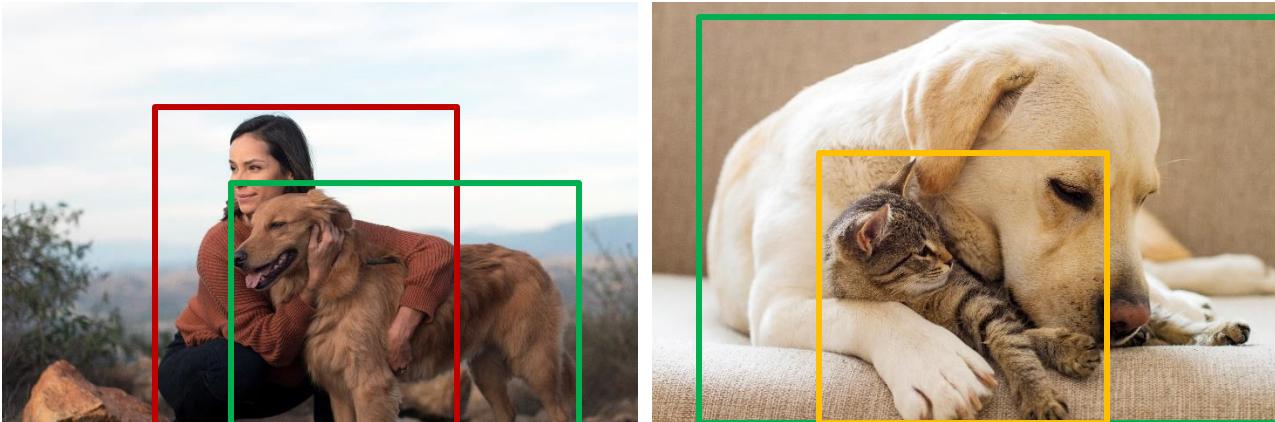
... Roche disappoints analysts and drags the SMI into the red. The daily loss is 4 percentage points higher than that of its competitor Novartis, which closes only just in the red at 0.1%. Financial stocks were able to protect the index from even greater losses, led by Credit Suisse, which ...

Ticker	Last	Change	Volume
CSGN	14.19	+2.5 %	8'142'755
NOVN	83.06	-0.1 %	4'875'390
PGHN	1433	+1.5 %	98'060'000
ROG	224.15	-4.1 %	2'7680'115
SREN	92.22	+0.7 %	2'298'221



Information extraction

Structured vs. unstructured Data



Object detection



Women	1	0
Men	0	0
Dogs	1	1
Cats	0	1
Environment	Outdoor	Indoor

Does not have a fixed data model, but can contain structural information as **tags** or **markers**, so that the schema can be (partially) (re-)constructed.

Entities of the same class might contain **different attributes** and the **order** is arbitrary.

Popular file formats: **XML** and **JSON** – both machine and human readable

Applications/Technologies:

- Communication in **Web-Applikationen** - JSON (REST) nowadays more common than XML
 - **E-Mail**
 - Document-based (**No SQL**) databases such as MongoDB and Elasticsearch (fulltext search)
-
- + Objects can be stored **without pre-defined schema**
 - No performant query language such as **SQL**
 - **Error-prone** due to missing data model
 - + **Mapping** to a relational data model unnecessary
 - + Natural representation of nested objects or hierarchical dependencies

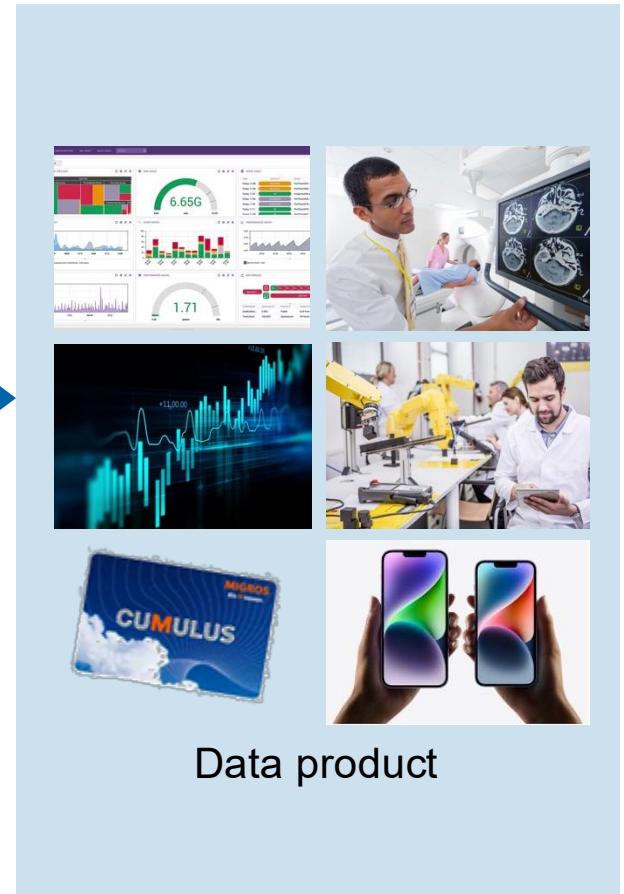
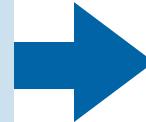
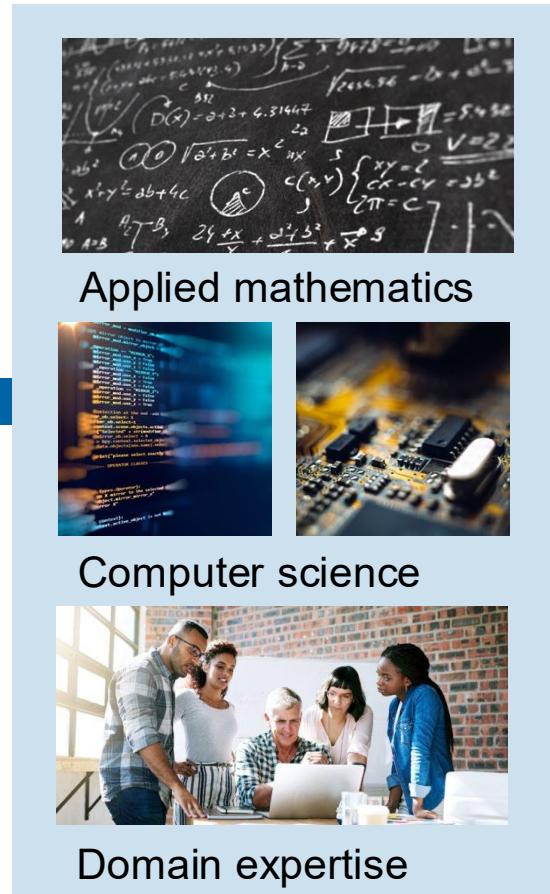
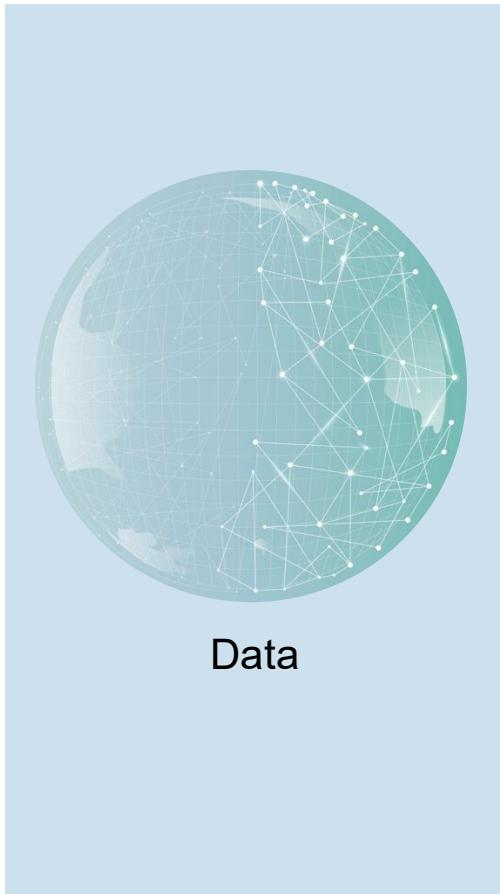
Metadata

«Data about data»: Contains additional information about specific datasets and supports their interpretation - e.g. description of the data model and the origin of the data

File Source	https://commons.wikimedia.org/wiki/File:Cat03.jpg
License	CC BY-NC
Attribution	Fir0002/Flagstaffotos
File Type	JPEG
File Size [kB]	152
X Dimension	1200
Y Dimension	1199
Resolution [pixels/inch]	72
Color Encoding	sRGB

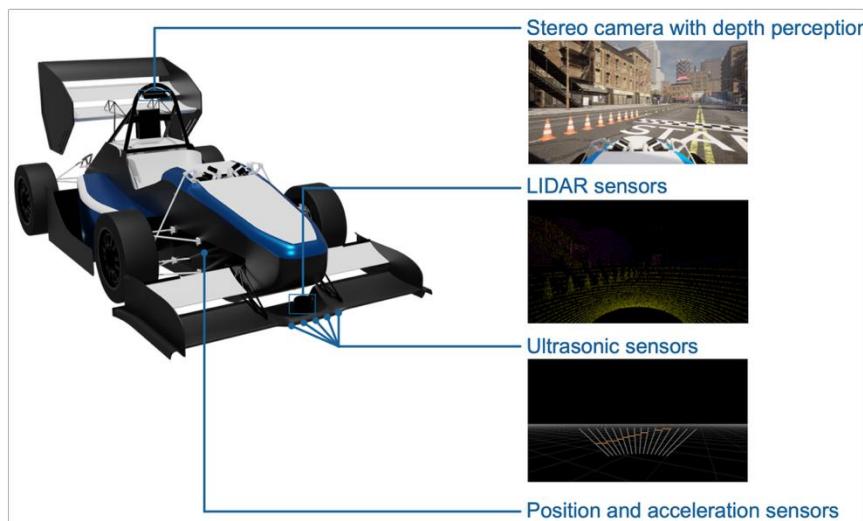


Data Science is data product development

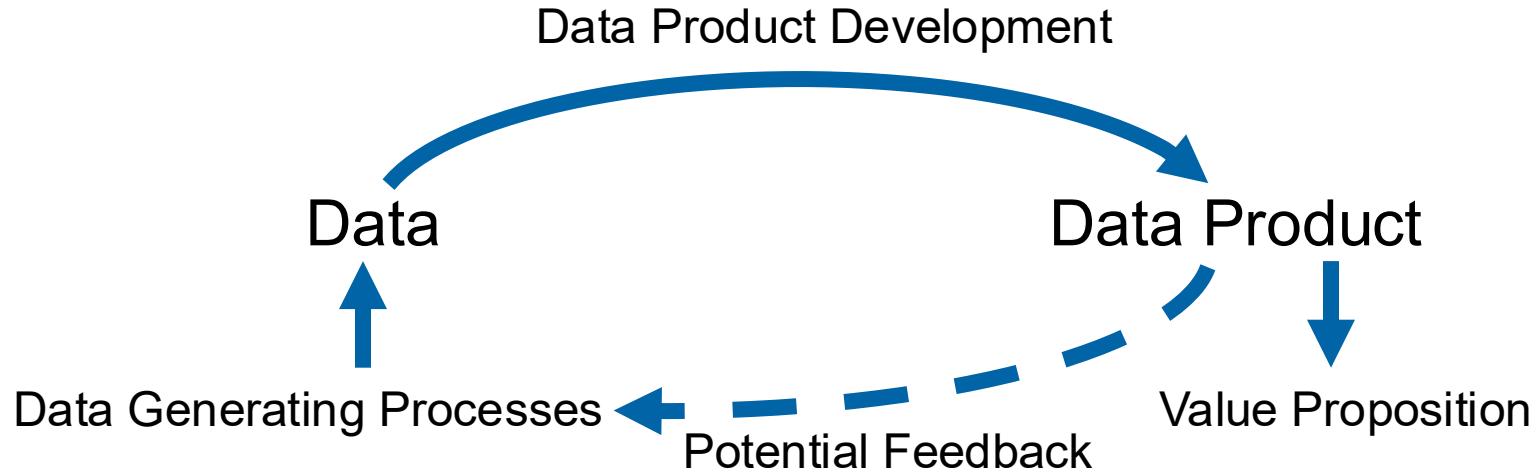


Use Case: Autonomous Vehicle

Cars can use cameras and other sensors to recognise their surroundings and navigate independently.



Data product

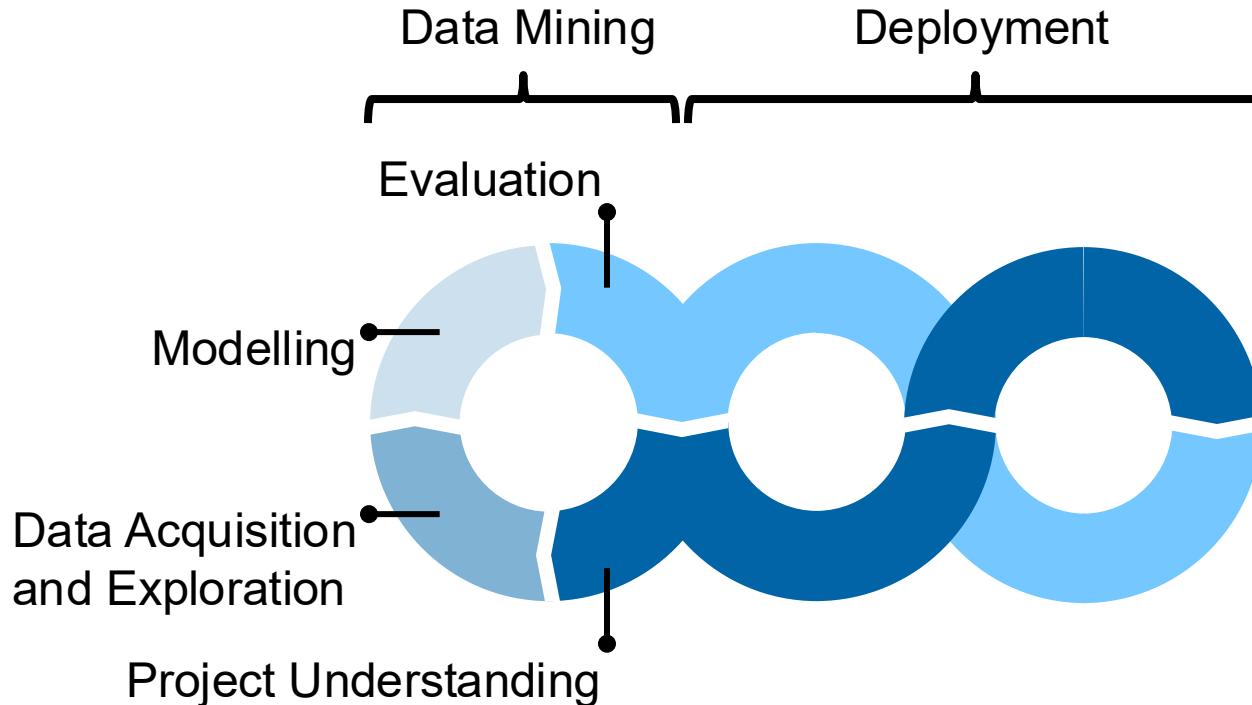


Value Proposition Canvas

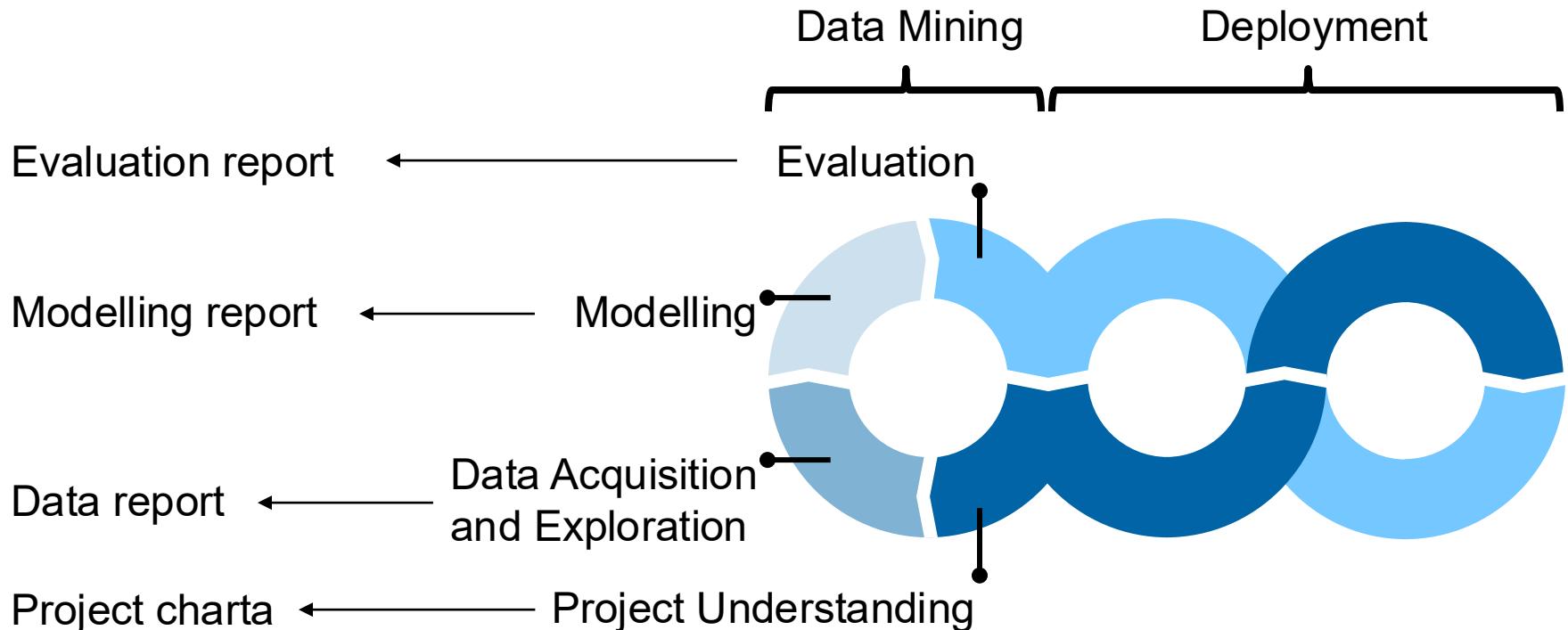


Osterwalder, A., Pigneur, Y., Bernarda, G., Smith, A., & Wegberg, T. A. A. E. (2015). *Value Proposition Design: Entwickeln Sie Produkte und Services, die Ihre Kunden wirklich wollen*. John Wiley & Sons, Incorporated.

Data Product-Development



Documentation Artefacts



Tooling



Local Development Environment

Suggestion:



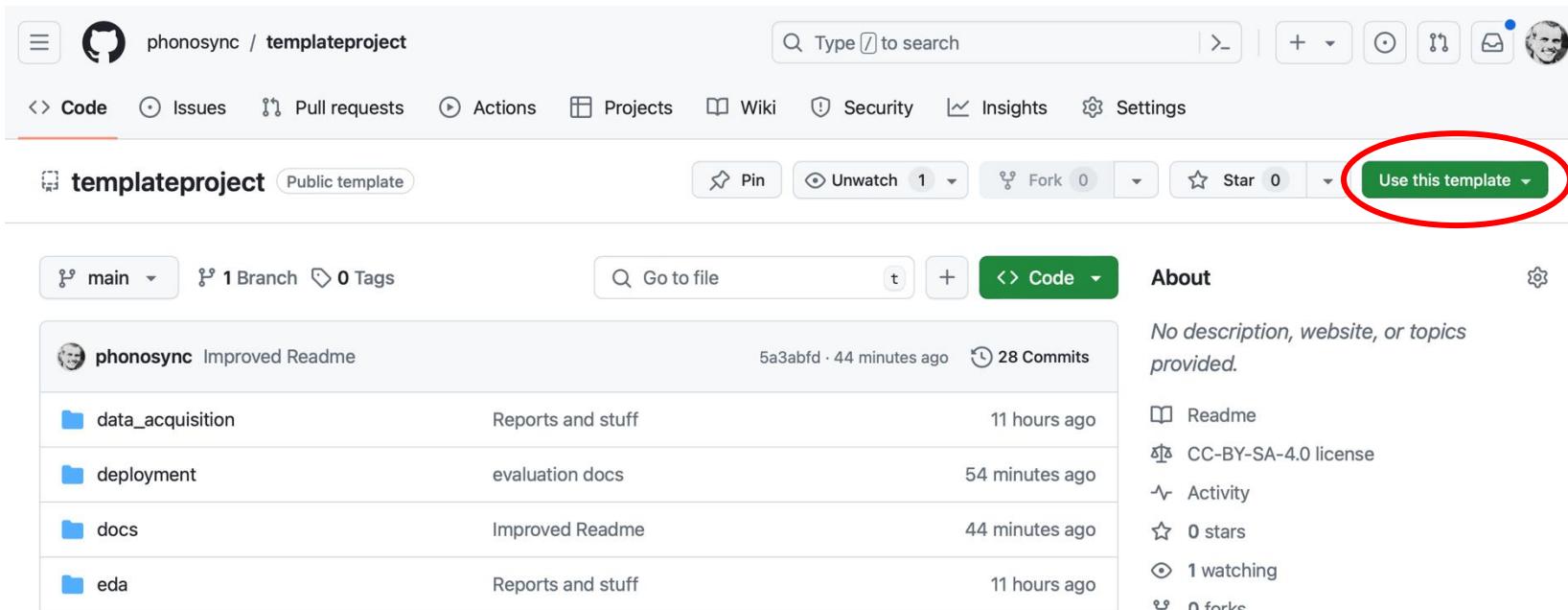
<https://code.visualstudio.com>

Extensions:

- Python (Microsoft)
- Jupyter (Microsoft)
- GitHub Copilot (Microsoft)
- Live Preview (Microsoft)
- Quarto

Structure and Version Control in a Data Science Project

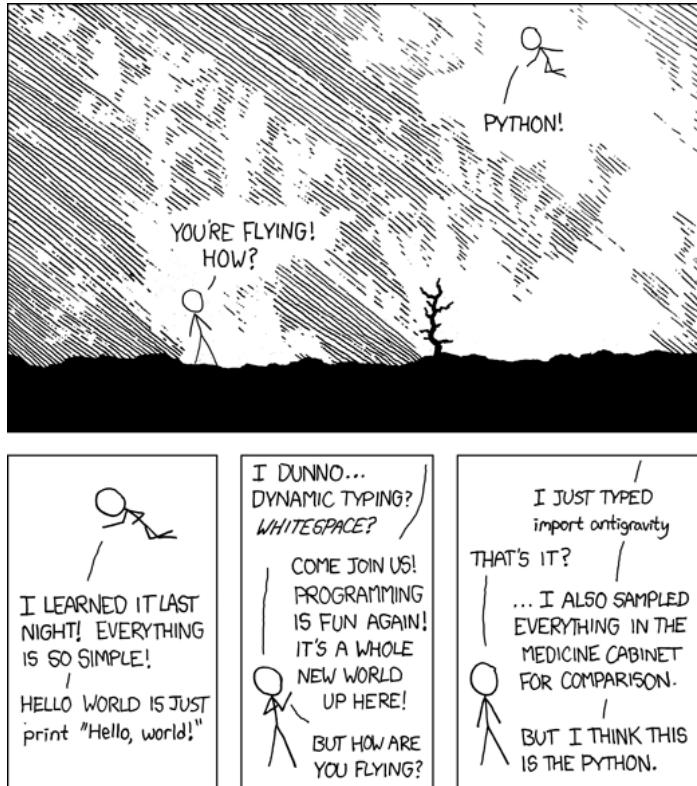
<https://github.com/phonosync/templateproject>



The screenshot shows a GitHub repository page for 'phonosync / templateproject'. The repository is a public template. The top navigation bar includes links for Code, Issues, Pull requests, Actions, Projects, Wiki, Security, Insights, and Settings. Below the navigation is a search bar and a 'Use this template' button, which is circled in red. The main content area displays the repository's structure and recent activity. The repository has 1 branch and 0 tags. The 'About' section notes that there is no description, website, or topics provided. The repository contains several files and folders, including 'data_acquisition', 'deployment', 'docs', and 'eda', each with their respective descriptions and last modified times.

File/Folder	Description	Last Modified
data_acquisition	Reports and stuff	11 hours ago
deployment	evaluation docs	54 minutes ago
docs	Improved Readme	44 minutes ago
eda	Reports and stuff	11 hours ago

Python - The power comes with the packages

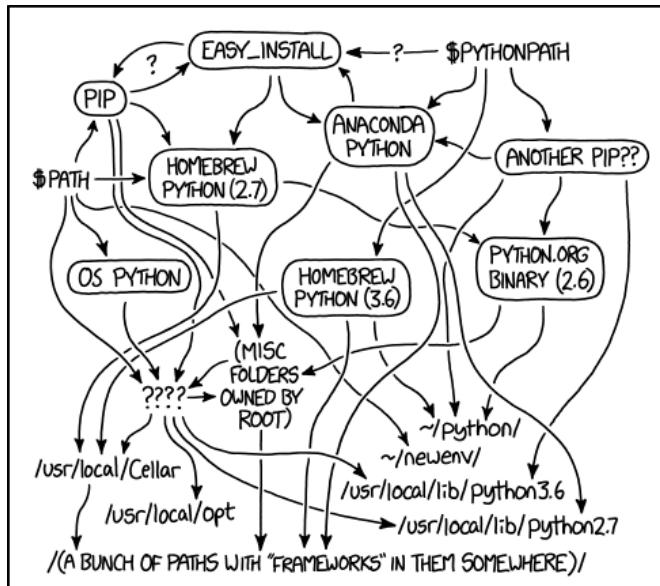


Popular registries:

- <https://pypi.org/>
- <https://anaconda.org/anaconda/repo>

Package Management

Different projects require different version of packages and dependencies...



Python Environments

... let you isolate project specific interpreter and package installations.

- venv
- virtualenv
- Anaconda
- uv

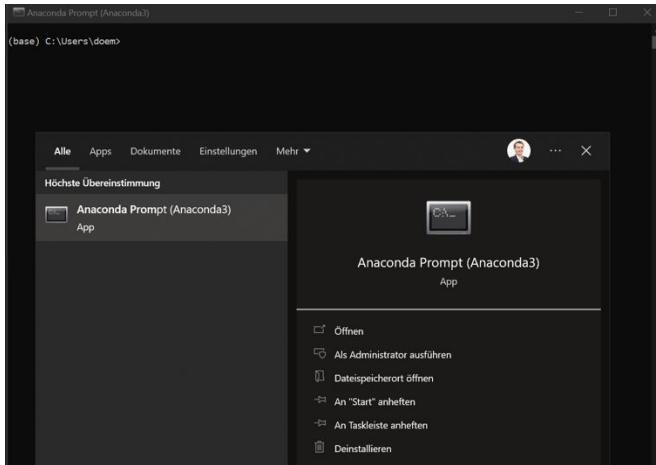
We go with **Anaconda** - it provides:

- a package index: <https://anaconda.org/anaconda/repo>
- a package manager: > conda install pandas
- an environment manager: > conda create --name myproject

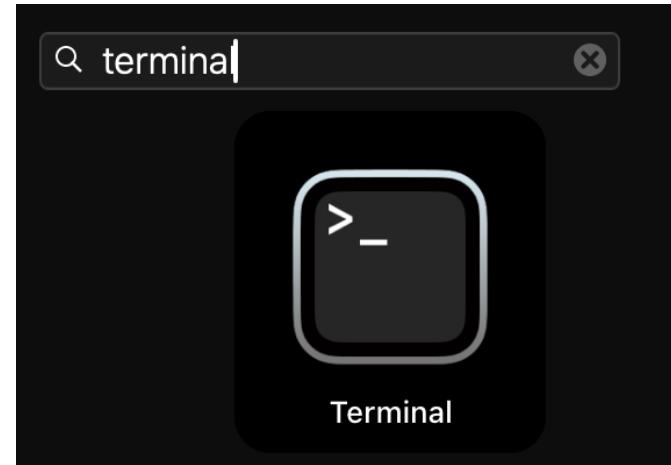
Download and install Anaconda for your system: [installation instructions](#)

Using conda

Windows



Mac OS X



- 2.In the console navigate to your project folder, e.g. cd C:\Users\doem\projects\demo
- 3.execute conda env create --file conda.yml
- 4.confirm installation

Using conda

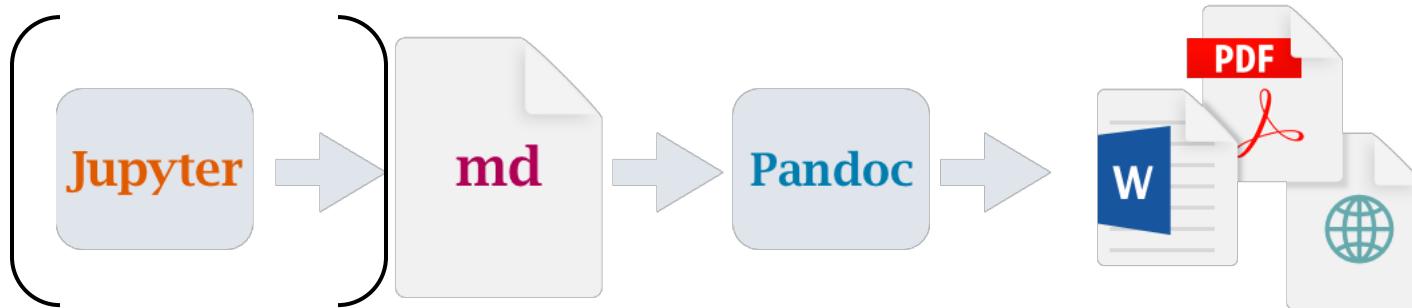
To activate the environment execute in the console:
`conda activate demo`

Remember to activate whenever you work on the project!

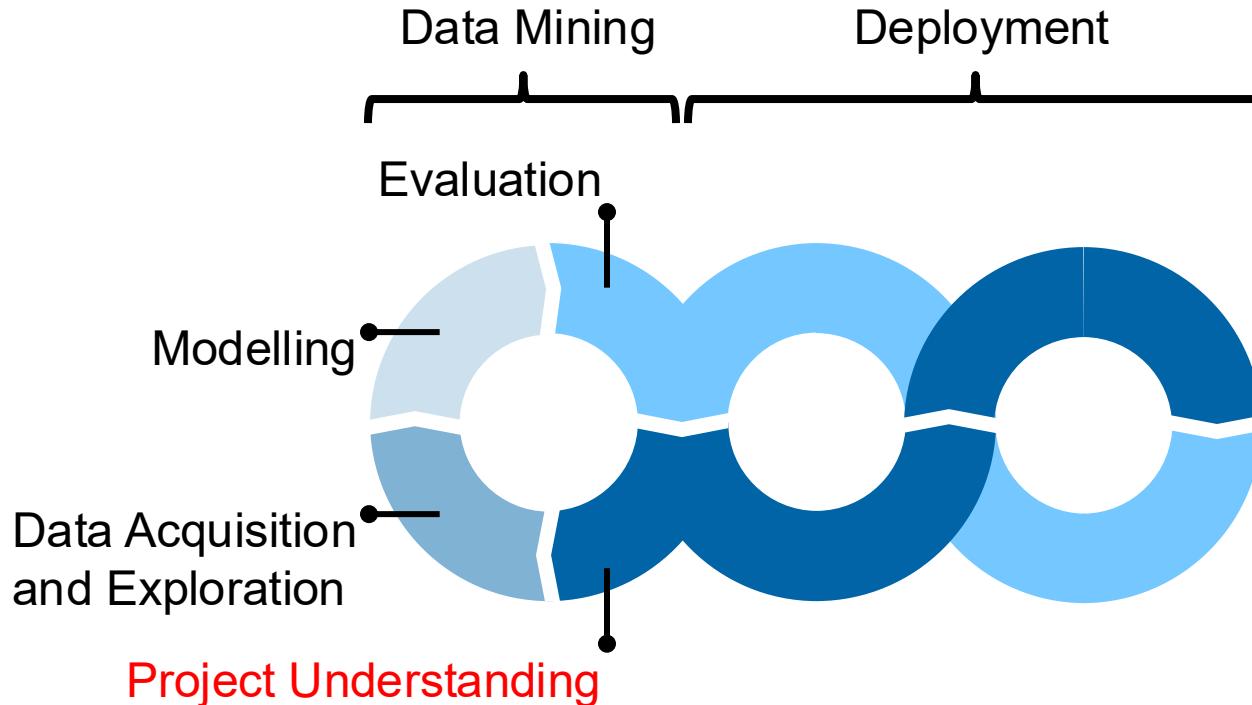
More information on managing the conda environments in the template project Readme

Project Documentation with Quarto

- Download from <https://quarto.org> and -> <https://quarto.org/docs/get-started/>



Data Product-Development



Project Understanding with Albert Einstein



„If I had an hour to solve a problem, I'd spend 55 minutes thinking about the problem and 5 minutes thinking about solutions“

Project Understanding Overview

	Activities	Output/Artifact
Project Understanding	<ul style="list-style-type: none">• Formulate problem statement• Analyse initial situation• Set project goals and success criteria• Derive data analysis goals• Make project plan	<ul style="list-style-type: none">• Project charta

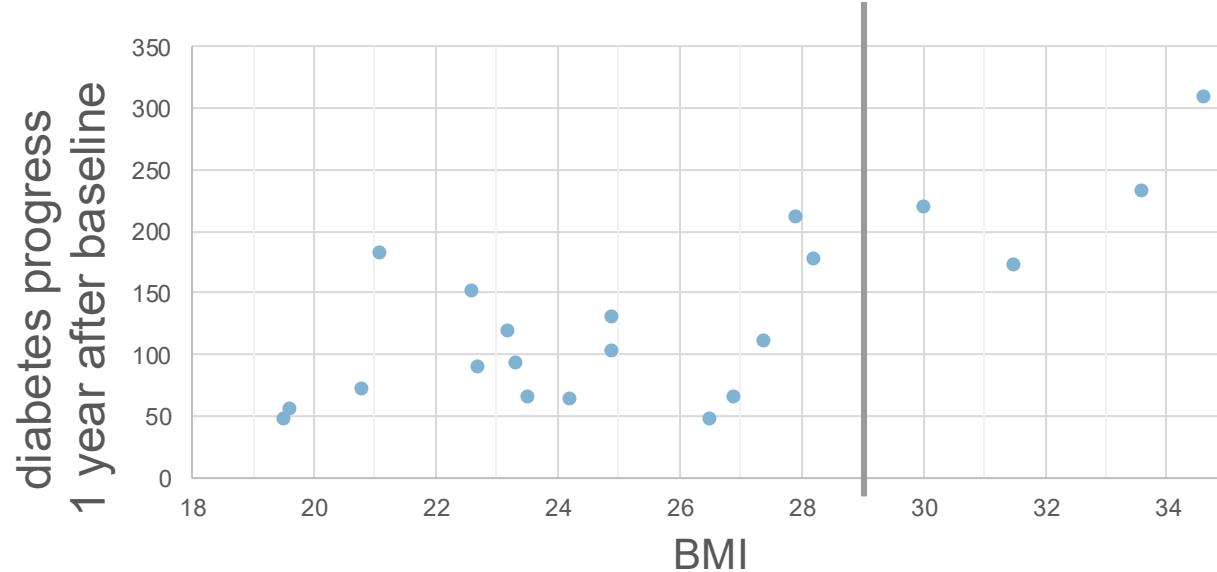
The target values of the key metrics provide transparency on progress and a quantitative basis for decision-making at project milestones.

Map the Business Problem onto a Data Mining Task

- Regression
- Classification
- Clustering
- Outlier Detection
- Association rule learning (market basket analysis)
- Recommender System

Regression

Example: Predict diabetes progress based on BMI



Regression aims at predicting a continuous target variable based on a set of dependent variables.

Classification

Example: Muffin/Chihuahua



Classification groups objects into distinct categories (classes).

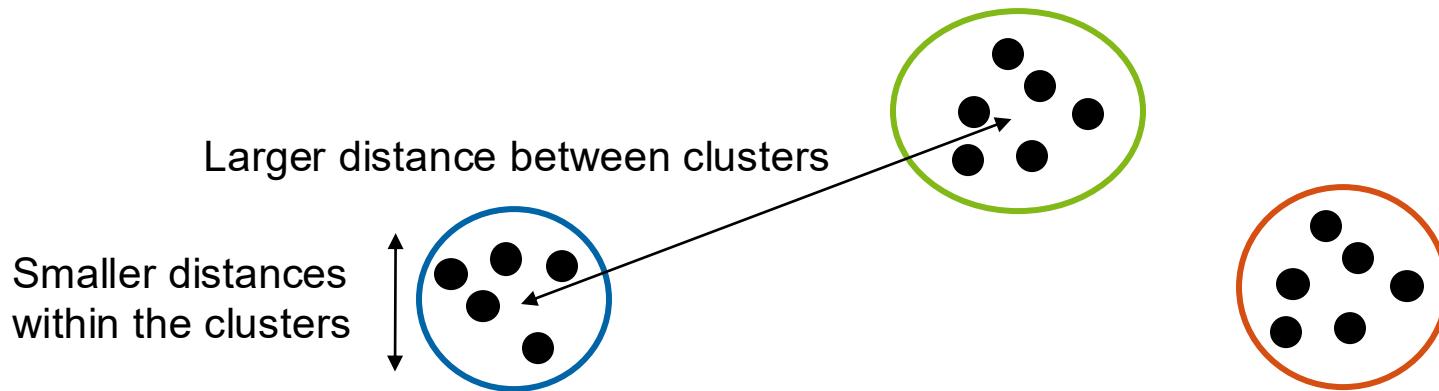
Clustering



Clustering identifies subgroups of datapoints that are more similar to each other than to the elements in other subgroups.

Clustering

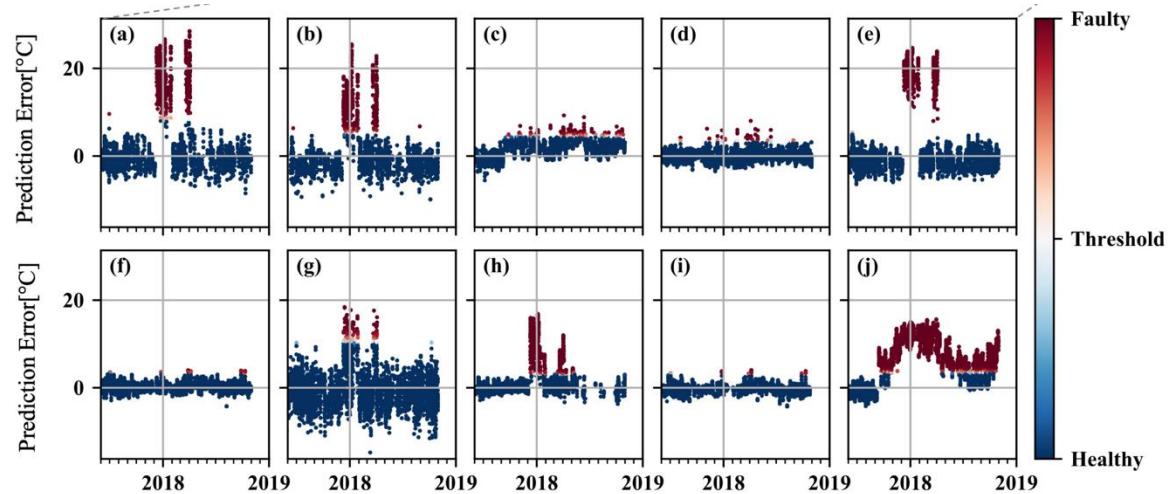
Goal: **Identify subgroups** of datapoints that are more similar to each other than to the elements in other subgroups.



→ Needs metric to quantify similarities.

Outlier Detection

Example: Detection of malfunctions from the sensor data of a wind turbine



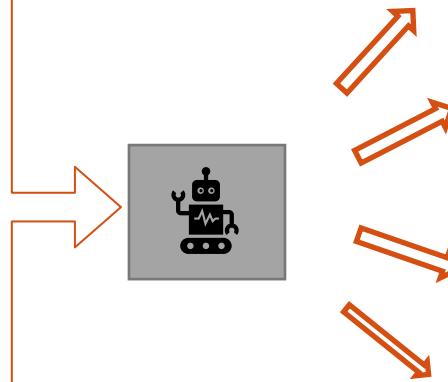
Outlier detection identifies datapoints that do not belong to the same distribution as „normal“ observations.

Recommender Systems

Example: Online shop – personalised recommendations
based on shopping behaviour

MIGROS | Operations

Bought	Sales	Feature XYZ
	True	CHF 12 23.23
	True	CHF 5 56.23
	False	CHF 0 12.34
	True	CHF 65 45.89
	True	CHF 12 23.54
	True	CHF 34 63.55
	False	CHF 0 28.23



Personalised recommendations

	1	
	2	
	3	
	:	
	N	

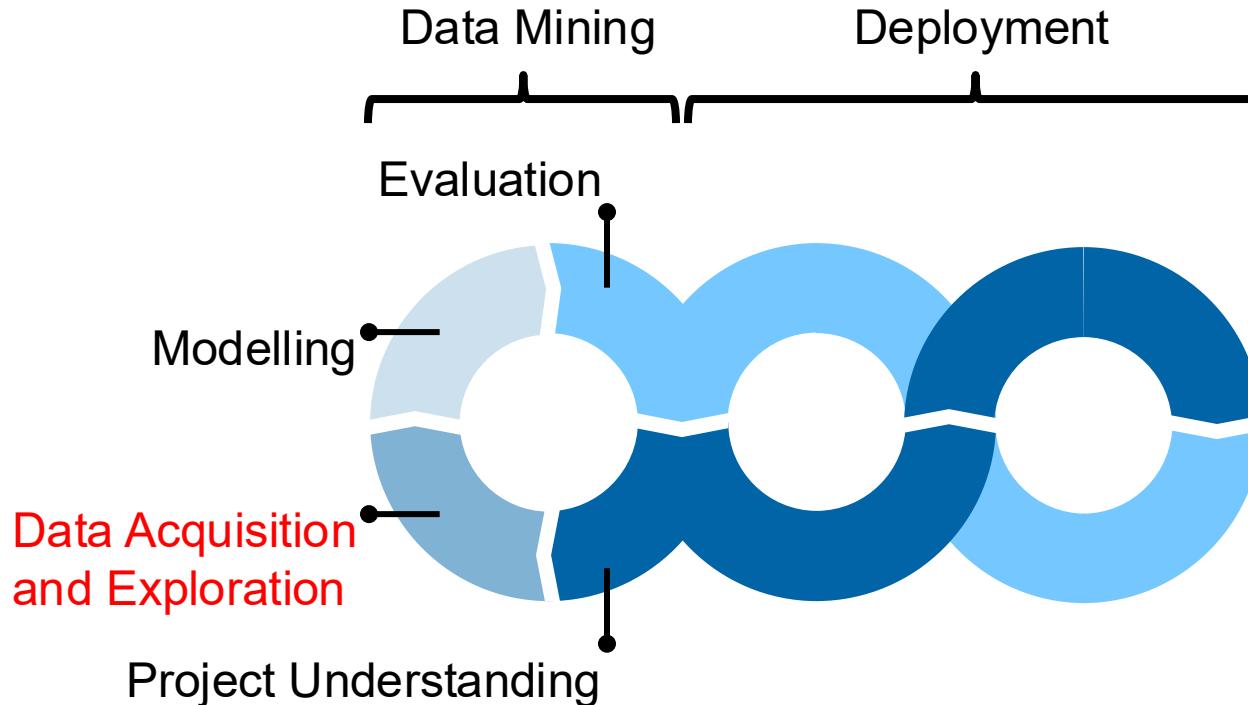
	1	
	2	
	:	
	N	

Recommender systems provides suggestions for items that are most pertinent to a particular user.

Identify the Data Mining Task

Develop an automatic system for a real estate agent that determines the sales price of a flat based on data (location, size, year of construction, type of construction).

Data Product-Development



Data Acquisition and Exploration

	Activities	Output/Artefact
Data Acquisition	<ul style="list-style-type: none">• Identify data sources• Acquire data• Store and organize data• Compile data catalogue	<ul style="list-style-type: none">• Data report
Data Exploration	<ul style="list-style-type: none">• Inject data in analytics environment• Explore data, identify main characteristics	

Setup programmatic data pipelines and document them to support traceability, reproducibility and allow for a systematic expansion of the data pool.

Exploratory data analysis

Data science-projects often start with an **incomplete knowledge** of the data creation process and different **data sources** are **combined**.

- Data Science «in the wild»
- «Observational Data»

A thorough understanding of the characteristics and data quality issues is required for the subsequent analysis and modelling steps.

Objectives of an exploratory data analysis

Identify...

- the main characteristics of the dataset
- data quality issues
 - missing attributes and/or values
 - duplicates
 - noise: erroneous entries and outliers
 - inconsistencies
- necessary pre-processing steps
 - cleaning
 - transformations/feature engineering, e.g. scaling, normalisation,...
- the need for more or different data
- properties that impact the subsequent modelling steps

Which issues can you find in the data?

id	title	date	rating	episode	season	length
1	live free or die	01.11.2012	9.3	1	5	43
2	madrigal	06.12.2013	89	2	5	48
3	sunset	25.10.2011	9.3	6	3	47
4	grilled	20.12.2009	9.3	2	2	46
5	down	42.12.2009	8.3	4	2	333
6	cancer man	19.10.2008	8.3	4	1	48
7	live fre or die	01.11.2012	9.3	1	5	43

firstname	lastname	bdate	age	gender	phone
bryan	cranston	03-07-1956	65	1	999-9999
aaron	paul	27-08-1979	44	m	777-53474
anna, gunn	gunn	08-11-1968	52	0	040-15627

(Near) duplicate detection

id	title	date	rating	episode	season	length
1	live free or die	01.11.2012	9.3	1	5	43
2	madrigal	06.12.2013	89	2	5	48
3	sunset	25.10.2011	9.3	6	3	47
4	grilled	20.12.2009	9.3	2	2	46
5	episode title	25.10.2011	9.3	6	3	46
6	cancer man	19.10.2008	8.3	4	1	48
7	live fre or die	01.11.2012	9.3	1	5	43

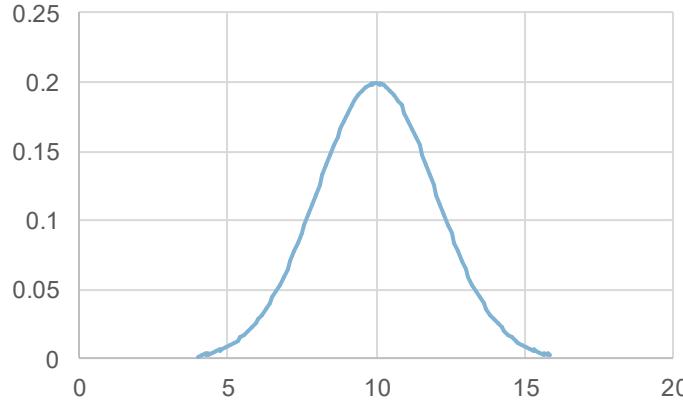
- For numeric values: Distance/similarity of feature vectors
- For text: pairwise Levenshtein distance between the texts

The tools in EDA: Descriptive statistics and visualisation

The data:

```
...  
7.60.09709303  
7.72 0.1041539  
7.840.11132675  
7.960.11856598  
8.080.12582217  
8.20.13304262  
8.320.14017191  
8.440.14715251  
8.560.15392563  
8.68 0.1604319  
8.8 0.1666123  
8.92 0.172409  
9.040.17776626  
9.160.18263134  
9.28 0.1869553  
9.40.19069391  
9.520.19380831  
9.640.19626574  
9.760.19804011  
9.880.19911242  
10.0.19947114  
10.120.19911242  
10.240.19804011  
10.360.19626574  
10.480.19380831  
10.60.19069391  
10.72 0.1869553  
10.840.18263134  
10.960.17776626  
11.08 0.172409  
11.2 0.1666123  
11.32 0.1604319  
11.440.15392563  
11.560.14715251  
11.680.14017191  
...
```

Normal Distribution



The probability density of the normal distribution

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

is determined by 2 parameters:

- Expectation value (mean) μ
- Standard deviation σ

The description:

- $\mu = 10.0$
- $\sigma = 2.0$

Statistical variable types

Variable Type	Level of Measurement (Scale)	Allowed Operations	Allowed Statistics	Examples
Categorical	nominal	count	mode, Chi-squared	profession, sex
	ordinal	count, sort	+ rank statistics (Spearman)	ratings, school marks
Numerical • discrete • continuous	interval	count, sort, subtract	+ mean, median, standard dev., correlation	time, temperature in Celsius
	ratio	count, sort, subtract, divide	+ geometric mean, harmonic mean	cost, age, size, temperature in Kelvin

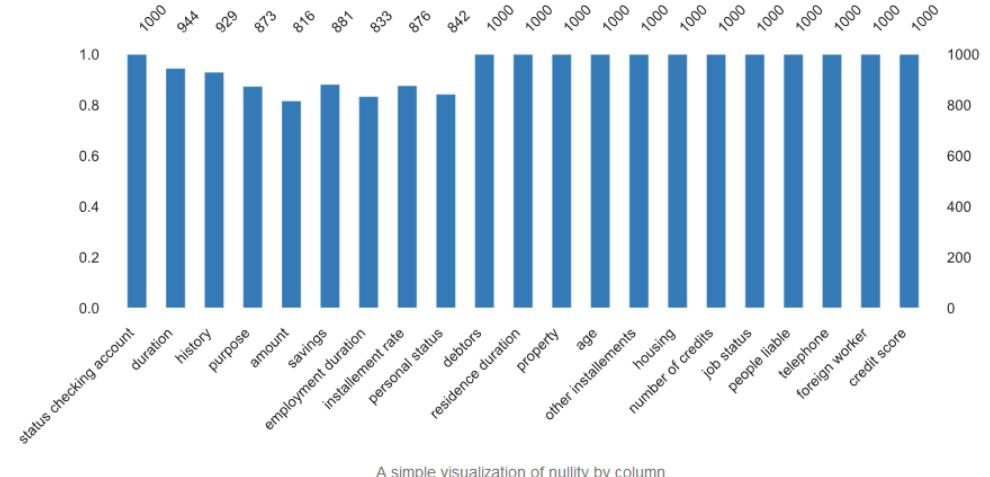
No strict connection between statistical variable types and data types / encoding schemes on the computer system.

Dataset overview

- number of variables per type, i.e. numerical/categorical
- number of observations
- missing values
- duplicates

Dataset statistics

Number of variables	21
Number of observations	1000
Missing cells	1006
Missing cells (%)	4.8%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	813.6 KiB
Average record size in memory	833.2 B



Univariate description of categorical variables

- number of unique values
- missing values
- frequency distribution of the values

purpose	Distinct	10
Categorical	Distinct (%)	1.1%
HIGH CORRELATION	Missing	127
MISSING	Missing (%)	12.7%
	Memory size	55.3 KIB

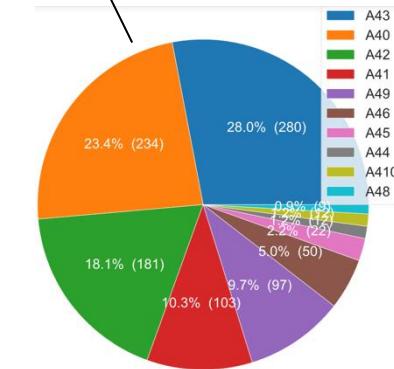
number of missing values
divided by number of observations

Common Values		
Value	Count	Frequency (%)
A43	248	24.8%
A40	201	20.1%
A42	160	16.0%
A41	90	9.0%
A49	80	8.0%
A46	45	4.5%
A45	19	1.9%
A44	11	1.1%
A410	11	1.1%
A48	8	0.8%
(Missing)	127	12.7%

relative frequency: number (Count) divided by total number of observations

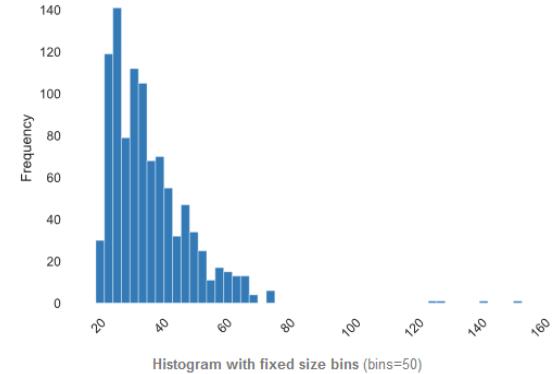
bar chart of the relative frequencies

pie chart: relative frequencies encode the area of the circle segments



Univariate description of numerical variables

- number of unique values
- frequencies
- missing values
- extreme values: minimum and maximum
- fraction of zero-values
- negative values
- histogramm



age

Distinct	57
Distinct (%)	5.7%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%
Mean	35.946

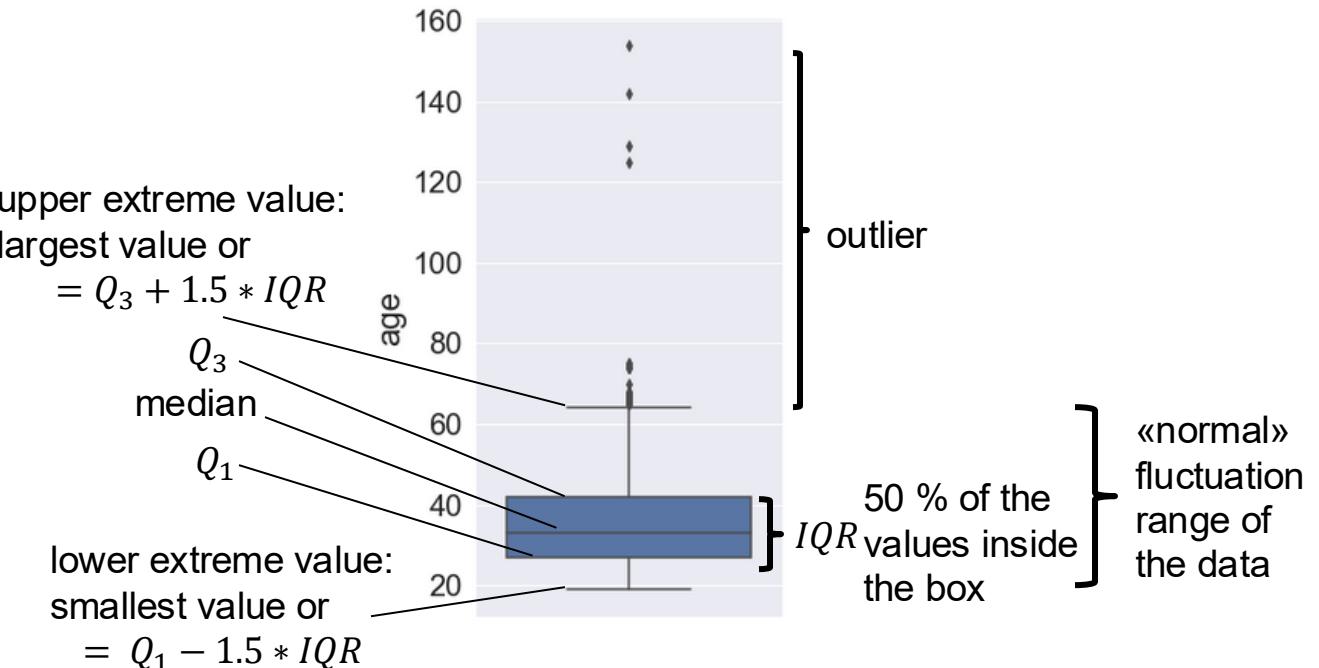
Minimum	19
Maximum	154
Zeros	0
Zeros (%)	0.0%
Negative	0
Negative (%)	0.0%
Memory size	7.9 KiB

Extreme values	
Minimum 5 values	Maximum 5 values
Value	Count Frequency (%)
154	1 0.1%
142	1 0.1%
129	1 0.1%
125	1 0.1%
75	2 0.2%
74	4 0.4%
70	1 0.1%
68	3 0.3%
67	3 0.3%
66	5 0.5%

Univariate description of numerical variables

- Central tendency:
 - arithmetic mean
 - median
 - mode
 - quantiles
- Dispersion metrics:
 - interquartile range IQR
 - variance
 - standard deviation

Five number-summary and boxplot:

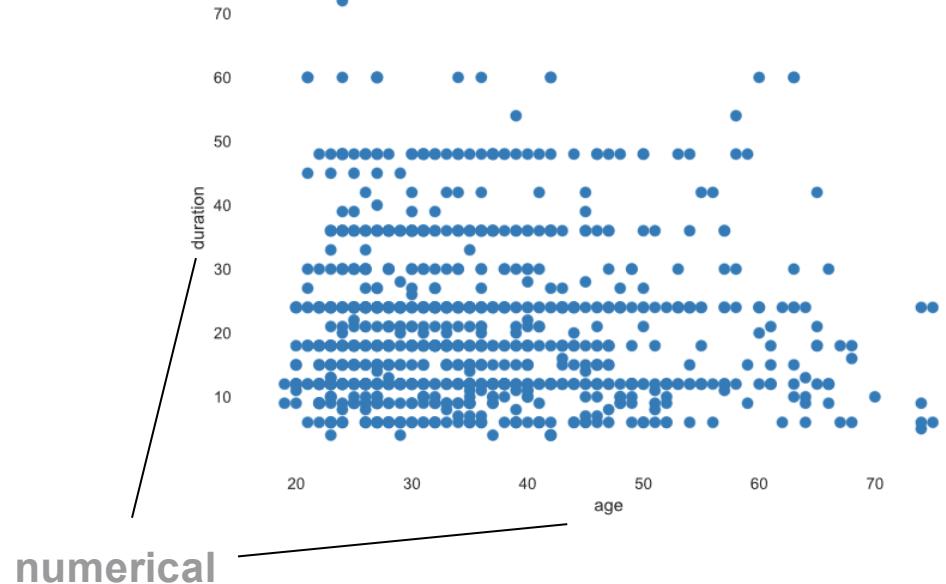
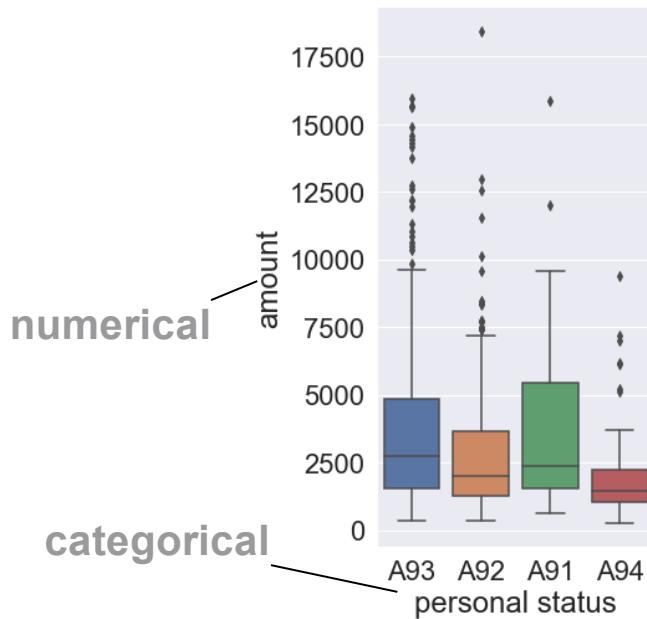


Outlier detection



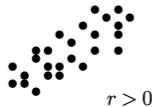
title	Date	...	length
live free or die	01.11.2012		43
madrigal	06.12.2013		48
sunset	25.10.2011		47
grilled	20.12.2009		46
down	22.12.2009		91
cancer man	19.10.2008		48
shotgun	09.11.2012		47

Assessing relationships between variables

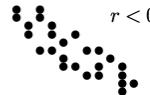


Measuring the strength of relationships between variables

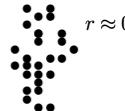
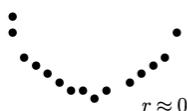
positive correlation,
values (x_i, y_i) group around
a line with positive slope



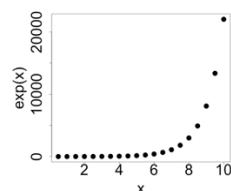
negative correlation,
inverse relationship,
values (x_i, y_i) group around
a line with negative slope



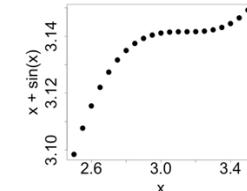
uncorrelated,
or no linear (or monotonic)
relationship



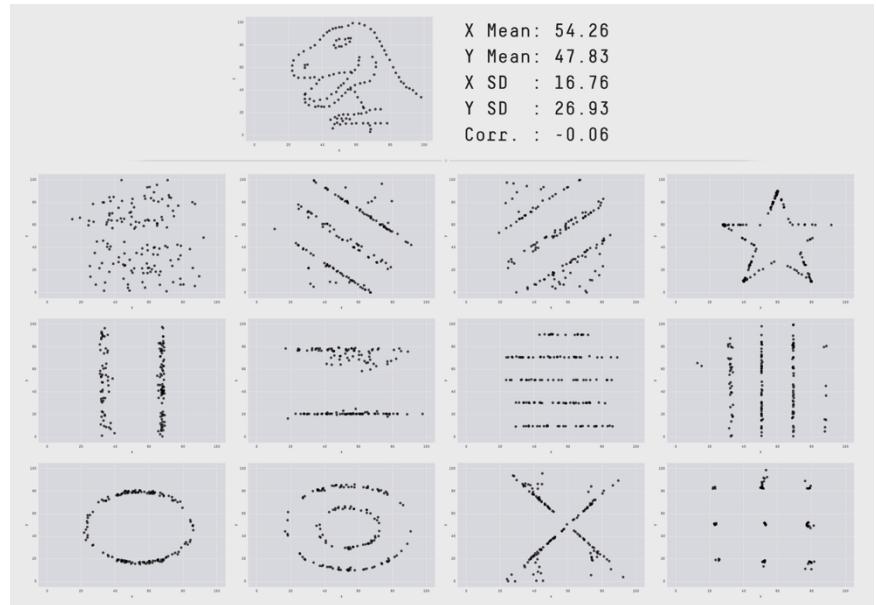
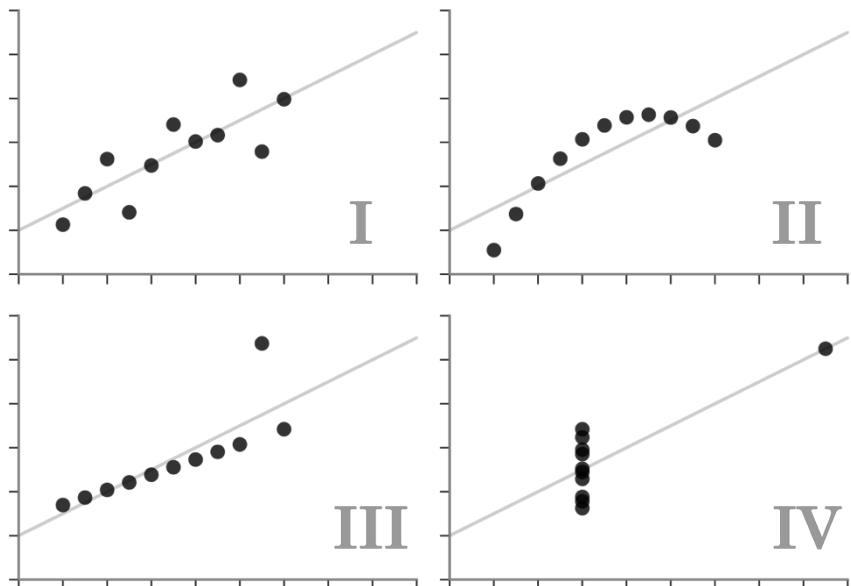
$$r = 0.7 \quad r_{SP} = 1.0$$



$$r = 0.86 \quad r_{SP} = 1.0$$



EDA: Combination of statistical metrics and visualisation



Matejka, J., & Fitzmaurice, G. (n.d.). *Same Stats, Different Graphs: Generating Datasets with Varied Appearance and Identical Statistics through Simulated Annealing*. <https://doi.org/10.1145/3025453.3025912>
<https://www.autodesk.com/research/publications/same-stats-different-graphs>

How do we deal with missing data?

title	date	rating	episode	season	length
live free or die	01.11.2012	9.3	1	5	43
madrigal	06.12.2013	8.9	2	5	48
sunset	25.10.2011		6	3	47
grilled	20.12.2009	9.3	2	2	46
down	22.12.2009	8.3	4	2	
cancer man	19.10.2008	8.3	4	1	48
shotgun	09.11.2012	8.7	5	4	47

How do we deal with missing or incorrect data?

Ignore the Tuple

- + can be easily done
- + no computational effort
- loss of information
- unnecessary if the attribute is not needed

Enter Value Manually

- + for small datasets effective
- + “real” value
- not for large datasets
- time consuming
- error-prone

Use Attribute Mean

- + simple to implement
- ignores correlations between attributes

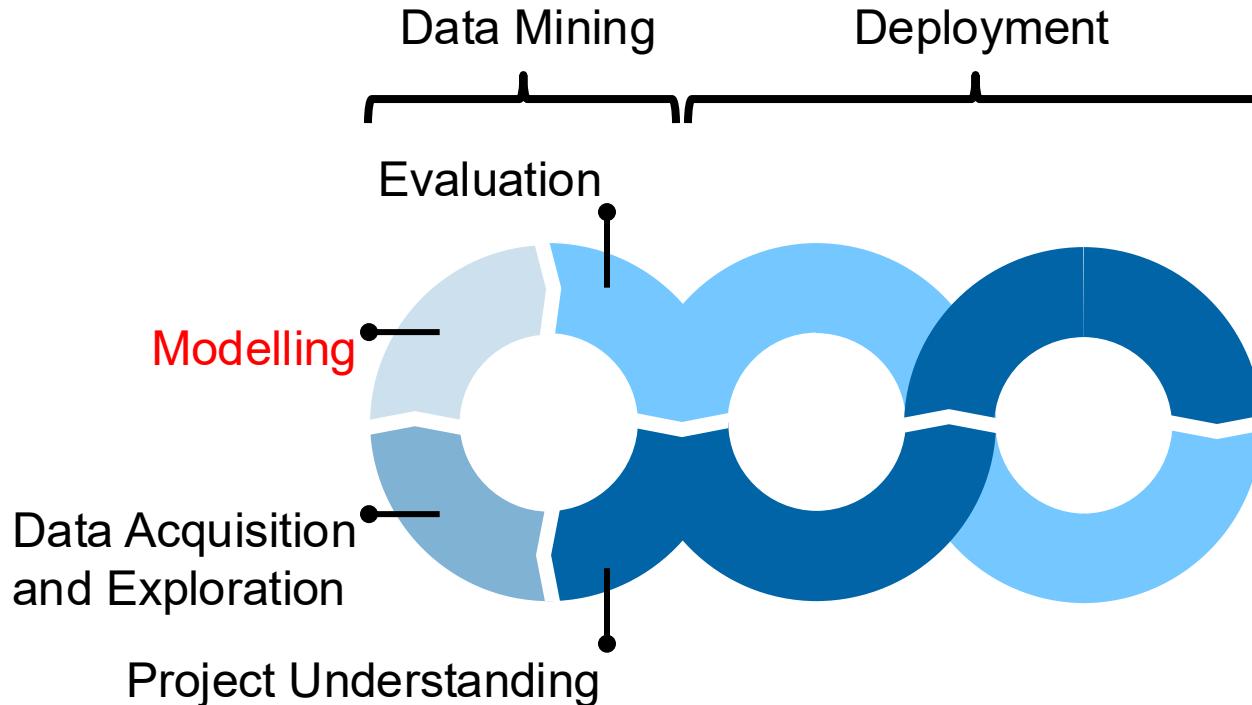
Use a Global Constant

- + can be easily done
- + missing values are marked
- values can not be used in algorithms

Use Most Probable Value

- + most accurate approximation of the value
- most computational effort

Data Product-Development



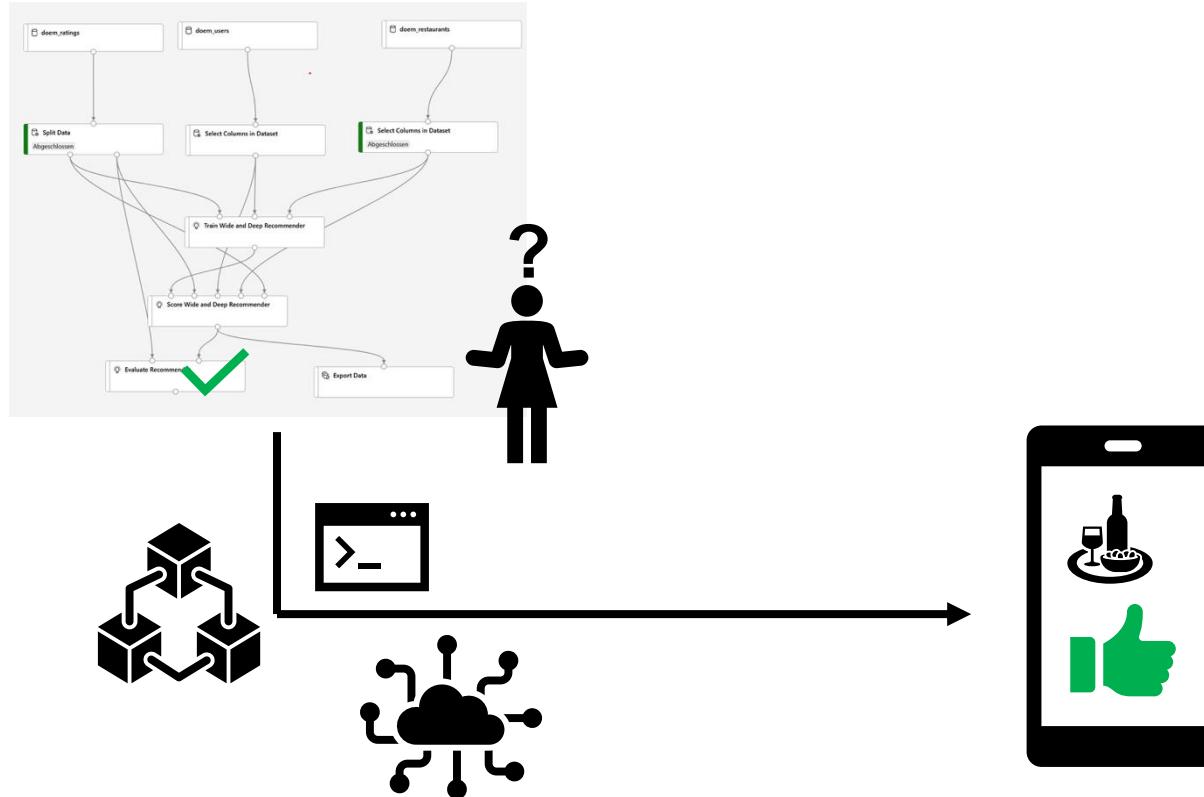
Modelling

	Activities	Output/Artefact
Modelling	<ul style="list-style-type: none">• Construct modelling pipeline• Transform data, engineer features• Train model• Evaluate model• Interpret results (predictions and model)	<ul style="list-style-type: none">• Modelling report

Evaluation

	Activities	Output/Artefact
Evaluation	<ul style="list-style-type: none">• Decide: Do results meet user needs?<ul style="list-style-type: none">→ planning of the deployment→ project stop→ Additional data mining iteration	<ul style="list-style-type: none">• Evaluation decision log

Deployment

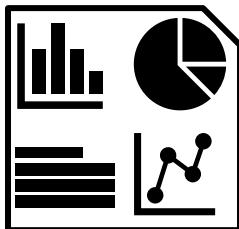


Drivers of Complexity

simple

complex

e.g. one –off report



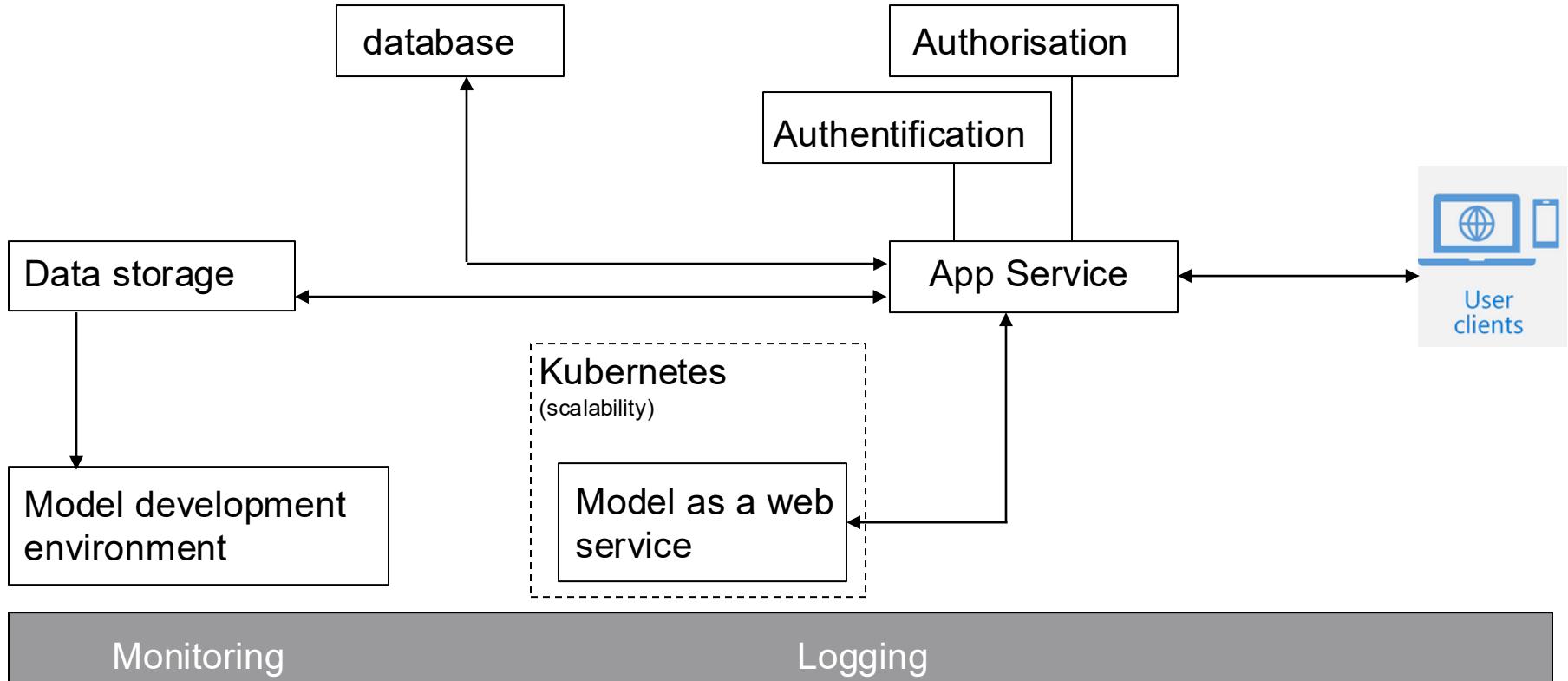
Drivers of complexity

- **Automated training** and deployment on changing data
- **Interactivity** and personalisation
- Integration in larger **software application**
- **External** users
- **Legal/regulatory** requirements

«continuously adapting»
recommender system in
online shop



Architecture Blueprint

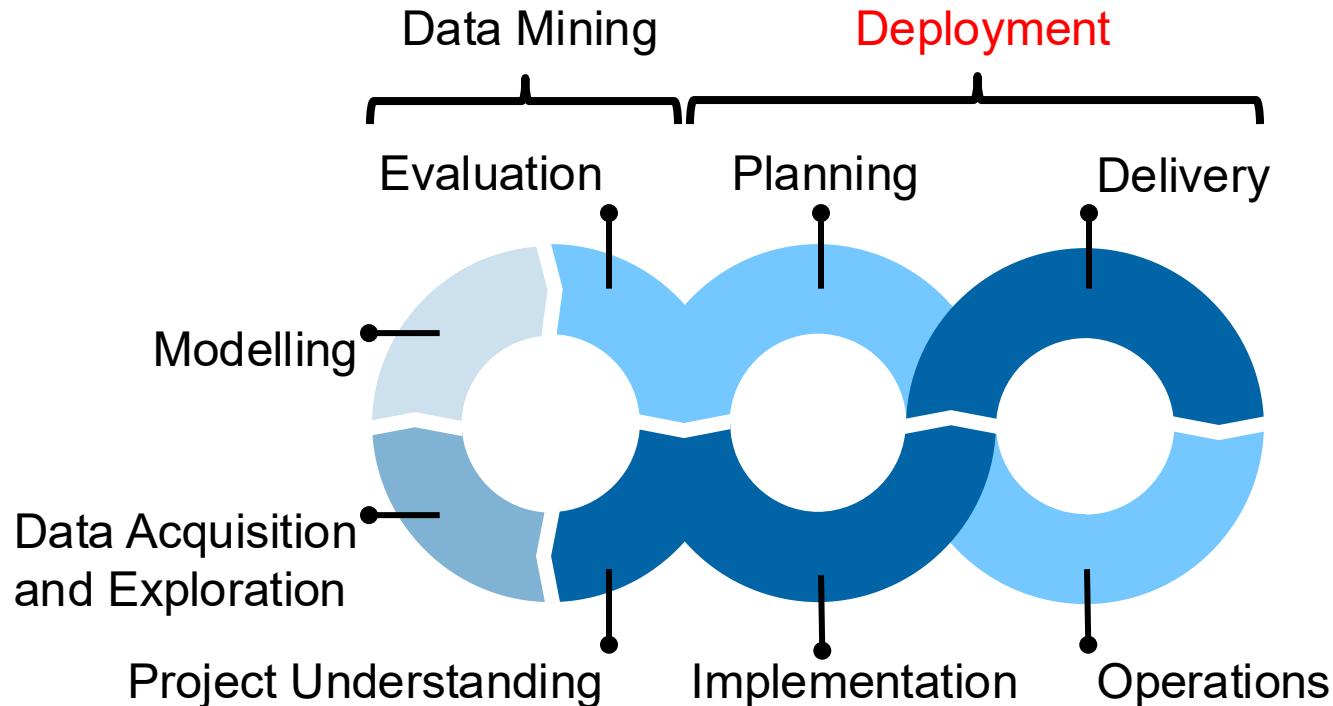


Deployment Beyond Technical Aspects

The data product only adds value if it is accepted by the user base. Aspects to consider are for example:

- Training and Change Management
- Adjustments to business processes
- Connection to other products and offers (platform)
- Continuous collection of user feedback for iterative refinements

Data Product-Development



Data Product-Development

References and Further Reading

- Kempf D., Dömer M. (2022). Is It Ops That Make Data Science Scientific? *Archives of Data Science*, Series A, vol 8, p. 12.
- Microsoft **Team-Data-Science-Process**:
<https://docs.microsoft.com/en-us/azure/architecture/data-science-process/overview>
<https://github.com/Azure/Azure-TDSP-ProjectTemplate>
- The “Cross-Industry Standard Process for Data Mining”: Shearer, C. (2000). The **CRISP-DM** Model: The New Blueprint for Data Mining. *Journal of Data Warehousing*, 5(4), 13–22.
(Datei auf Moodle)
- **KDD**: Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). From Data Mining to Knowledge Discovery in Databases. *AI Magazine*, 17(3), 37–54. <https://doi.org/10.1609/aimag.v17i3.1230>
(Datei auf Moodle)