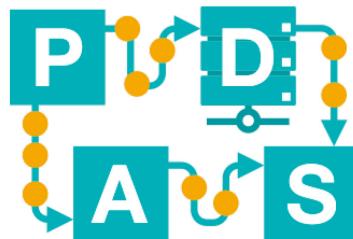


# Introduction to Data Science

Lecture 1

# IDS-L1



Chair of Process  
and Data Science

RWTH AACHEN  
UNIVERSITY

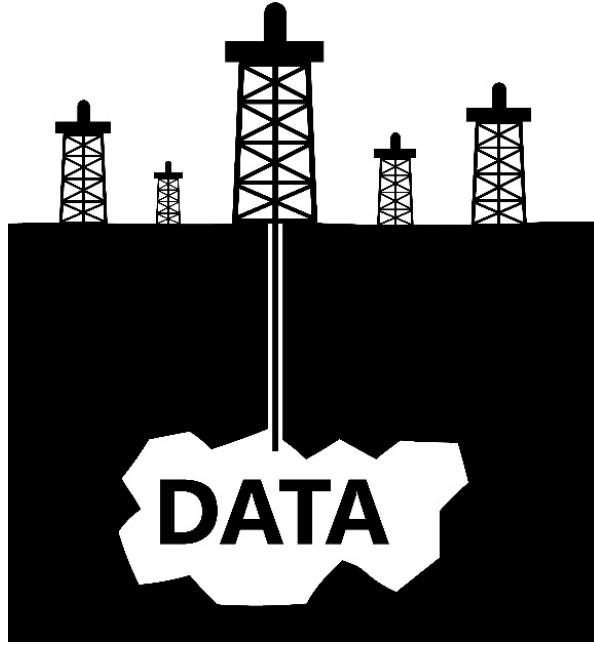
# Outline of Today's Lecture

- Outline of course
- Organization and exam regulations
- Motivation (Big Data & Skills)
- DS pipeline
- Types of data
- Terminology
- Data science process
- Challenges

# Outline of course



# Our society is driven by data, the new oil



- At all levels: personal, device, system, system of systems, organization, nation, world.
- Data volume grows exponentially.
- Allowing for new products and services.

- **exploration** (locating the data),
- **extraction** (how to get it),
- **transform** (clean and filter data)
- **storage** (Big data)
- **transport** (getting it to the right person)
- **usage** (analysis, actions, etc.)



- **data can be copied, oil not**
- **data is specific, oil is not**
- **if small, data storage and transport are cheap**

A woman with dark hair and bangs, wearing a black turtleneck and a colorful, patterned skirt with black fringe. She has her arms extended wide, showing off multiple bracelets and rings. Her nails are painted white. The background is a solid red.

# Four generic data science questions

#1



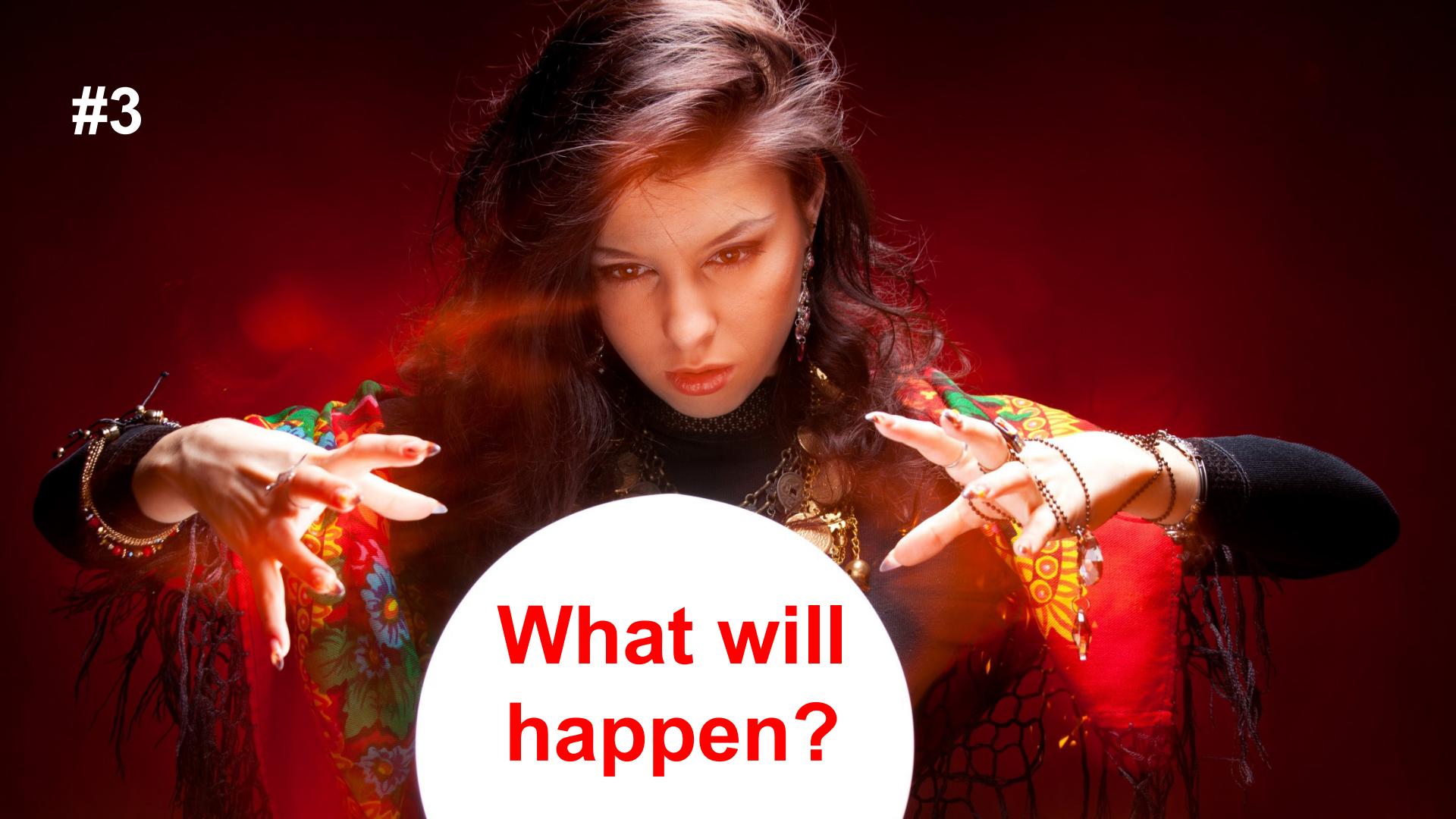
What  
happened?

#2



Why did  
it happen?

#3



What will  
happen?

#4



What is  
the best that  
can happen?

# Why now?



Example:



Sensors and actuators to align  
the digital and physical world!

LUCKY LUKE

Our digital  
shadows are  
catching up!



# Dimensions

- **Different types of data** (structured, unstructured, text, images, events, etc.).
- **Different types of tasks** (supervised, unsupervised).
- **Human versus machine** (who does what?).
- **Algorithm versus visualization.**
- **Flexibility versus usability** (e.g., Python versus Disco).
- **Scalability versus quality.**
- **Responsibility** (fairness, privacy, transparency, etc.).

# Lectures

1. Introduction
2. Crash course in Python
3. Basic data visualization/exploration
4. Decision trees (IG)
5. Regression
6. Support vector machines
7. Neural networks (1/2)
8. Neural networks (2/2)
9. Evaluation of supervised learning problems
10. Clustering (K-means)
11. Frequent items sets (a priori and FP growth)
12. Association rules
13. Sequence mining
14. Process mining (unsupervised)
15. Process mining (supervised)
16. Text mining (1/2)
17. Text mining (2/2)
18. Data preprocessing, data quality, binning, etc.
19. Visual analytics & Information visualization
20. Responsible data science: Fairness (discrimination-aware data mining, ethics, etc.)
21. Responsible data science: Confidentiality (various types of encryption and anonym.)
22. Big data (1/2): MapReduce & Distribution
23. Big data (2/2): Systems (Scala, Hadoop, etc.?)
24. Closing

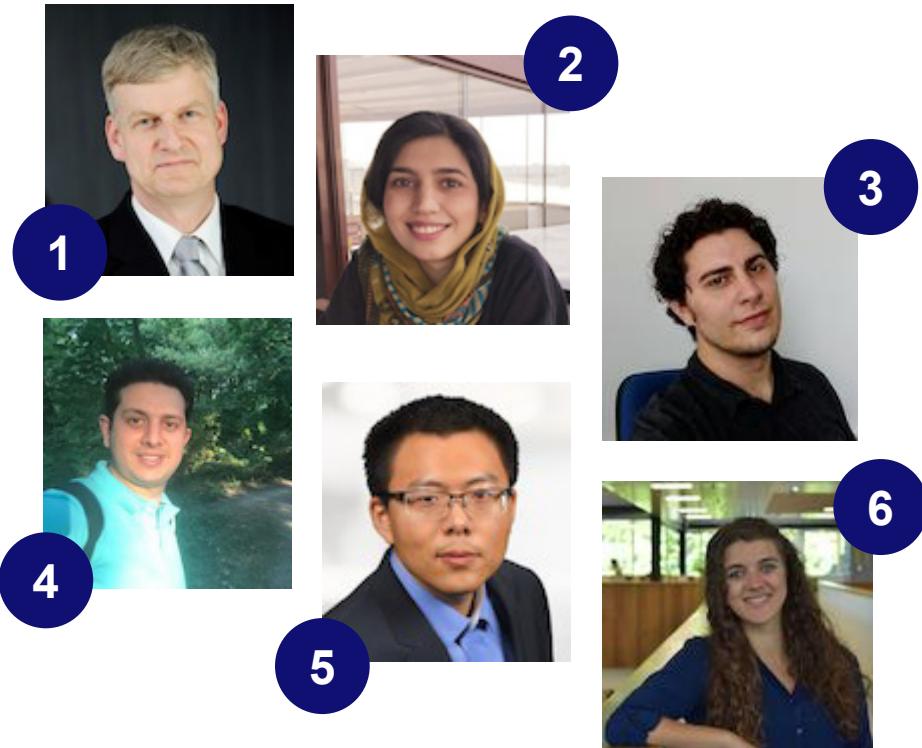


# Organization and exam regulations



# People involved

1. Wil van der Aalst
2. Mahsa Bafrani
3. Marco Pegoraro
4. Majid Rafiei
5. Yaguang Sun
6. Anja Syring

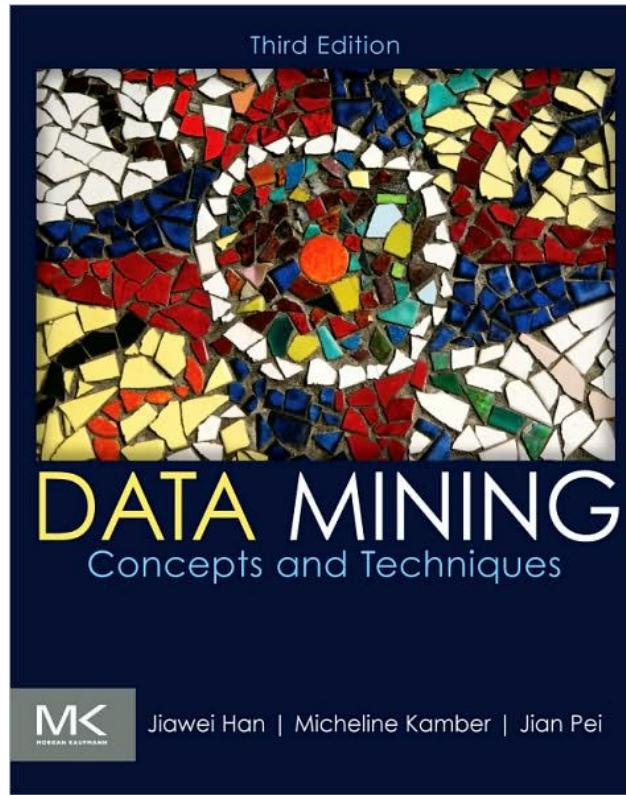
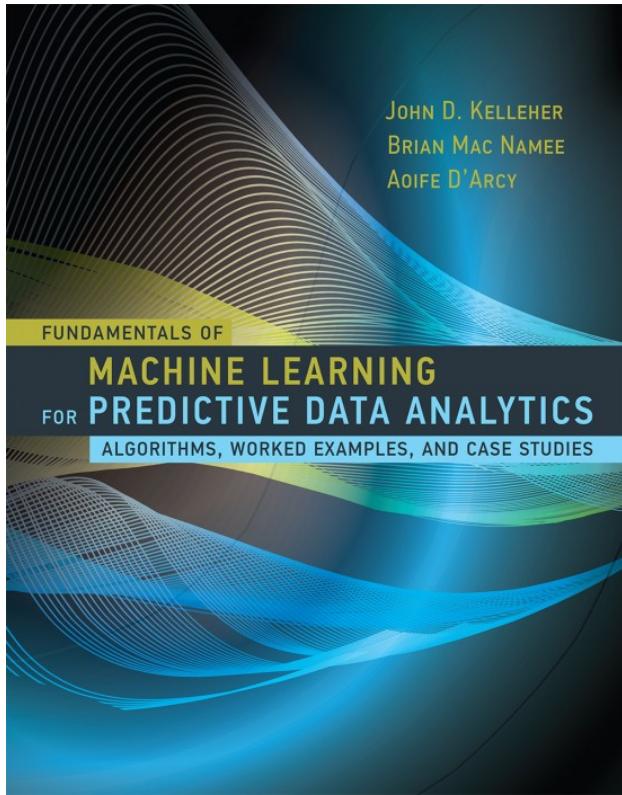


# Lectures

- The course starts on 10-10-2018.
- Lectures are on normally Wednesdays and Thursdays from 8.30 to 10.00 in Aula 2 (2352|021).
- Instructions are on Fridays from 8.30 to 10.00 also in Aula 2 (2352|021).
- Note that there are a few exceptions!

#	Lecture	date	day
	Lecture 1	10/10/2018	Wednesday
	Lecture 2	11/10/2018	Thursday
	Instruction 1	12/10/2018	Friday
	Lecture 3	17/10/2018	Wednesday
	Lecture 4	18/10/2018	Thursday
	Instruction 2	19/10/2018	Friday
	Lecture 5	24/10/2018	Wednesday
	Lecture 6	25/10/2018	Thursday
	Instruction 3	26/10/2018	Friday
	Lecture 7	31/10/2018	Wednesday
	Instruction 4	02/11/2018	Friday
	Lecture 8	07/11/2018	Wednesday
	Lecture 9	08/11/2018	Thursday
	Instruction 5	09/11/2018	Friday
	Lecture 10	14/11/2018	Wednesday
	Lecture 11	15/11/2018	Thursday
	Lecture 12	21/11/2018	Wednesday
	Lecture 13	22/11/2018	Thursday
	Instruction 6	23/11/2018	Friday
	Lecture 14	28/11/2018	Wednesday
	Lecture 15	29/11/2018	Thursday
	Instruction 7	30/11/2018	Friday
	Lecture 16	05/12/2018	Wednesday
	Instruction 8	06/12/2018	Thursday !!
	Lecture 17	12/12/2018	Wednesday
	Lecture 18	13/12/2018	Thursday
	Lecture 19	19/12/2018	Wednesday
	backup	20/12/2018	Thursday
	Instruction 9	21/12/2018	Friday
	Lecture 20	09/01/2019	Wednesday
	Lecture 21	10/01/2019	Thursday
	Instruction 10	11/01/2019	Friday
	Lecture 22	16/01/2019	Wednesday
	Lecture 23	17/01/2019	Thursday
	Instruction 11	18/01/2019	Friday
	Lecture 24	23/01/2019	Wednesday
	backup	24/01/2019	Thursday
	Instruction 12	25/01/2018	Friday
	backup	30/01/2019	Wednesday
	backup	31/01/2019	Thursday
	extra	01/02/2019	Friday

# Material



- ❑ Fundamentals of Machine Learning for Predictive Data Analytics: Algorithms, Worked Examples, and Case Studies by John D. Kelleher, Brian Mac Namee and Aoife D'Arcy. MIT Press. ISBN: 9780262029445, 624 pages, July 2015 (<http://machinelearningbook.com/>).
- ❑ Data Mining: Concepts and Techniques (3<sup>rd</sup> edition) by Jiawei Han , Micheline Kamber , Jian Pei. The Morgan Kaufmann Series in Data Management Systems, Elsevier. ISBN: 9780123814807, 744 pages, 2011 (<http://hanj.cs.illinois.edu/book>).

# Material

- Additional background material will be provided as the course progresses.
- Given the nature of the course there is not one textbook (many superficial books or books that focus on one aspect; a rapidly developing area).
- Distribution of slides, assignments, etc. via **RWTHmoodle**.

# Software (see next lecture)

- During the first instruction of the course, we will give detailed installation procedures for all the software.
- You will also find the procedures on RWTHmoodle.
- The study guide describes the required software. For the course, you will need a working installation of **Python 3.6.x** with some common Data Science software packages like:
  - numpy
  - scipy
  - pandas
  - scikit-learn
  - jupyter
- More packages will need to be installed later in the course.
- Anaconda Python is the advised Python distribution.  
Pip is the advised package manager.  
JetBrains Pycharm Community Edition is the advised IDE.

# Examinations and Assignments

- The exam consists of three parts: **two assignments** (Schriftliche Hausarbeit) each counting for 20% of the final result, and the **final written test** which counts for remaining 60% of the final result.
- Participation in the assignments is required for participation in the final test.
- Only the final test can be retaken in this semester (there will be one re-exam). Assignments can only be redone in the next academic year.

# Examinations and Assignments

- Final written test (60%) Questions to test the theoretical knowledge of the algorithms and techniques learned:
  - First option (PT1): **25-02-2019 09:00 – 11:00 in Aula 2**
  - Second option (PT2): **25-03-2019 09:00 – 11:00 in Aula 2**
- Schriftliche Hausarbeit / DS Assignment 1 (20%): Analysis of a real-life and/or synthetic data sets using the techniques and tools provided in the course. This assignment is used to test the understanding of the material in lectures 1-10. Deadline Sunday **09-12-2018 23:59**.
- Schriftliche Hausarbeit / DS Assignment 2 (20%): Analysis of more complex data sets using various data science techniques. This includes the interpretation of the results and creatively using multiple views on the data. The focus is on the lectures 11-21. Deadline Sunday **20-01-2019 23:59**.

Important: participation in both Schriftliche Hausarbeiten / Assignments is a prerequisite for taking the written exam. The three parts form a whole and it is not possible to retake parts of the course, i.e., the results of the assignments expire after the exam.

# Make sure

- to register for the course and final exam, and
- to hand in the assignments in time.

(sounds simple, but there are always sad cases)

# Who can take this course?

- As long as capacity allows we welcome all students.
- However, the course is a master course. It is mandatory for students taking the Data Science master (it is listed as a Wahlpflichtfach, but a requirement for doing a master thesis). It is a Wahlpflichtfach for Informatik, Media Informatics, and Software Engineering.
- Other students are free to participate, but it is up to the management and rules of the corresponding programs to decide whether the course "counts". (Ask your Studienberater/Prüfungsamt.)
- If you have problems with RWTHonline, RWTHmoodle, etc., that are not specific for this course, please contact the persons responsible for these systems and not the lecturer.

# Warning

- Slides are intended to be self-contained, but ...
- Overlap with earlier / later courses is unavoidable.  
(This is an introduction course and gives a broad overview.)
- The diversity of the course will make it tricky, so stay synchronized.

# Motivation (Big Data & Skills)



# Impact and size of data are obvious

## 2018 *This Is What Happens In An Internet Minute*

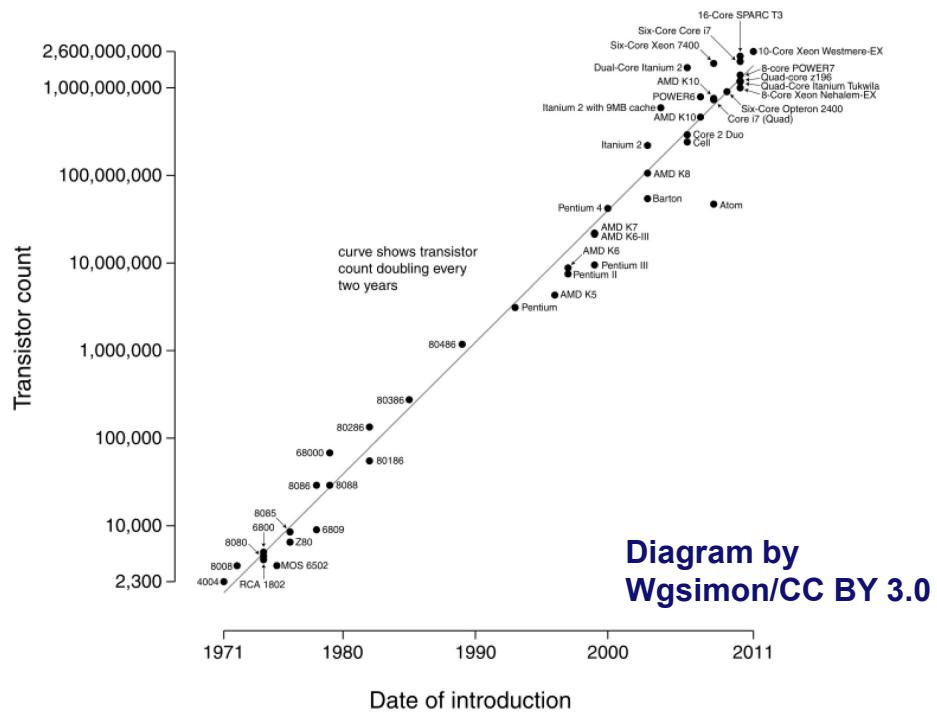
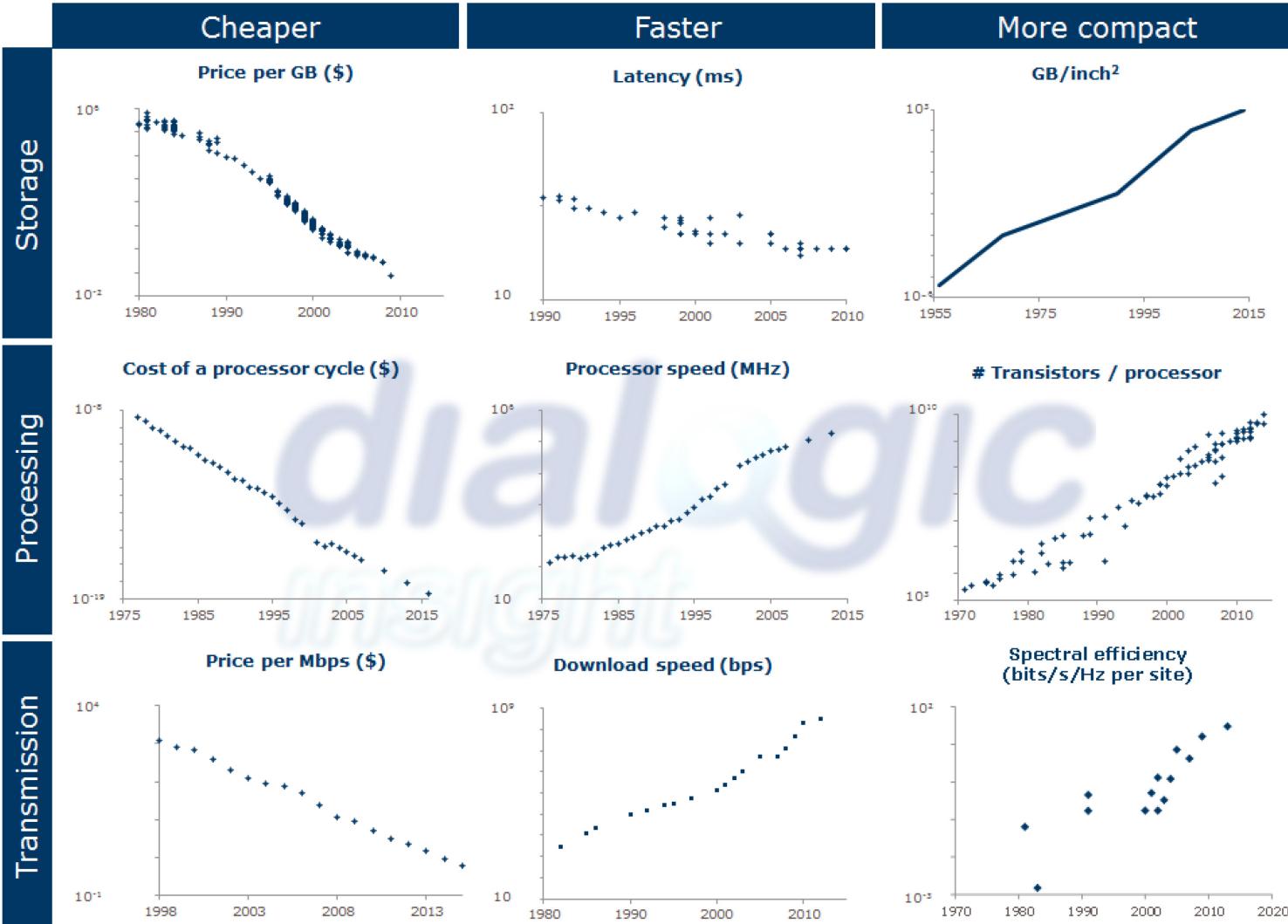


Diagram by  
Wgsimon/CC BY 3.0

$$2^{20} = 1.048.576 \text{ in 40 years}$$



Chair of Process  
and Data Science



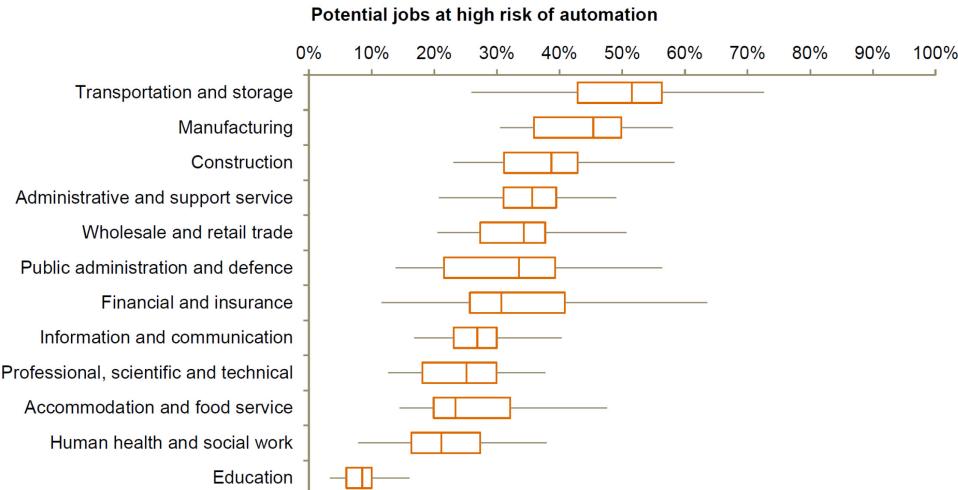
(Source: Dialogic, 2014. The impact of ICT on the Dutch economy based on Bronnen: [mkomo.com](http://mkomo.com) (linksboven), [storagenewsletter.com](http://storagenewsletter.com) (middenboven), IBM (rechtsboven), [singularity.com](http://singularity.com) (linksmidden + midden), Moore's law Wikipedia (rechts midden), [drpeering.net](http://drpeering.net) (linksonder), Nielsen's law (midden onder), Spectrale efficiëntie mobiele technologie 1G t/m 4G: Wikipedia.org (rechtsonder).

# Will robots really steal our jobs?

An international analysis of the potential long term impact of automation



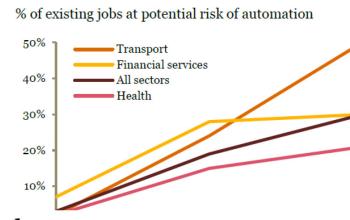
**Figure 4.1 – Share of jobs with potential high automation rates by industry**



## Key findings: impact of automation

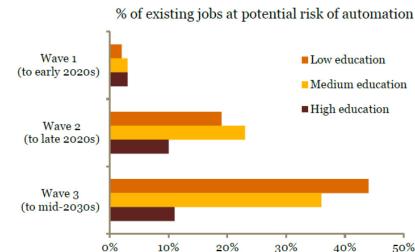
Financial services jobs could be relatively vulnerable to automation in the shorter term, while transport jobs are more vulnerable to automation in the longer term

**Figure 1 – Potential job automation rates by industry across waves**



In the long run, less well educated workers could be particularly exposed to automation, emphasising the importance of increased investment in lifelong learning and retraining

**Figure 2 – Potential job automation rates by education level across waves**



Source: PwC estimates based on OECD PIAAC data (median values for 29 countries)

Waves	Description and impact
<b>Wave 1:</b> <b>Algorithmic wave (to early 2020s)</b>	Automation of simple computational tasks and analysis of structured data, affecting data-driven sectors such as financial services.
<b>Wave 2:</b> <b>Augmentation wave (to late 2020s)</b>	Dynamic interaction with technology for clerical support and decision making. Also includes robotic tasks in semi-controlled environments such as moving objects in warehouses.
<b>Wave 3:</b> <b>Autonomous wave (to mid-2030s)</b>	Automation of physical labour and manual dexterity, and problem solving in dynamic real-world situations that require responsive actions, such as in transport and construction.

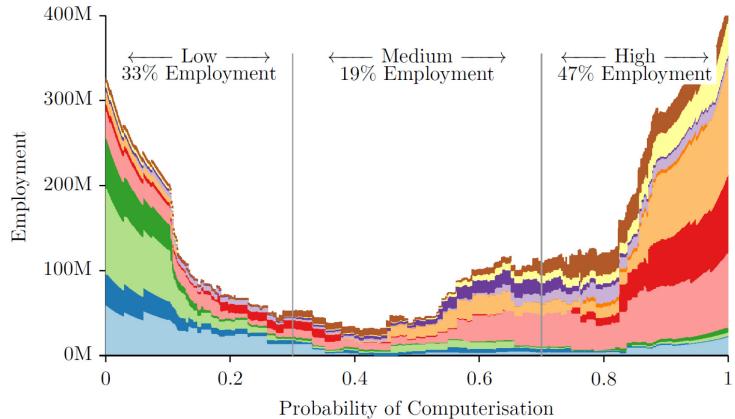
**"According to our estimates around 47 percent of total us employment is in the high risk category. We refer to these as jobs at risk – i.e. jobs we expect could be automated relatively soon, perhaps over the next decade or two."**



Published by the Oxford Martin Programme  
on Technology and Employment

### Occupations and probability of computerization (sample from 702 occupations):

- **Healthcare Social Workers 0.35%**
- **Firefighters 17%**
- **Statisticians 22%**
- **Accountants and Auditors 94%**
- **Bookkeeping and Accounting 98%**
- **Tax Preparers 99%**



# **Scope and impact of the growth of data are undeniable**



**data scientist**

# Data scientists are valuable



Cartoons where made  
for DSC/e were I was  
Scientific Director  
from 2013-2018

DSC/e



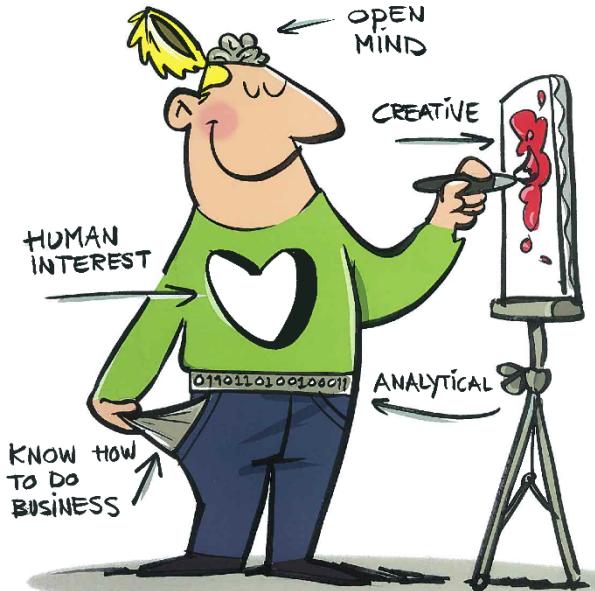
Chair of Process  
and Data Science

# Data scientists need to combine different skills

THE PERFECT DATA SCIENTIST

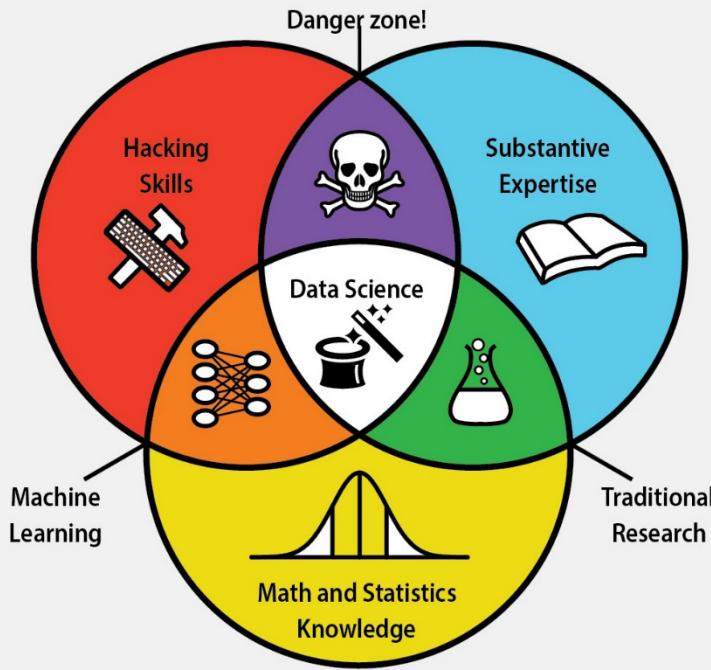


THE PERFECT DATA SCIENTIST



# Another characterization

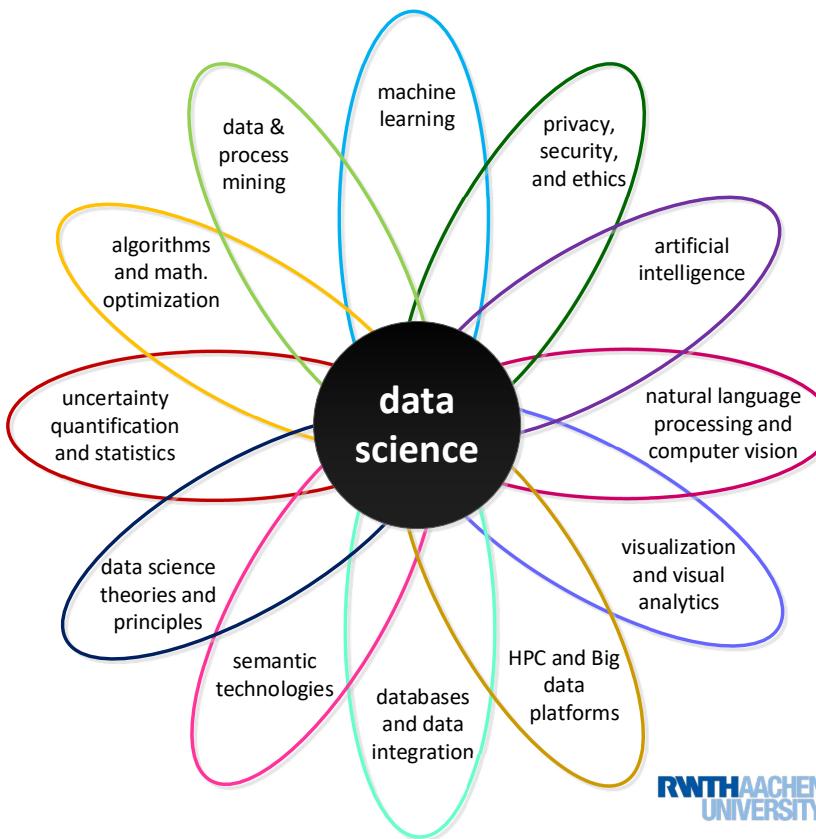
## DATA SCIENCE SKILLSET



	Data science, due to its interdisciplinary nature, requires an intersection of abilities: <b>hacking skills, math and statistics knowledge</b> , and <b>substantive expertise</b> in a field of science.
	<b>Hacking skills</b> are necessary for working with massive amounts of electronic data that must be acquired, cleaned, and manipulated.
	<b>Math and statistics knowledge</b> allows a data scientist to choose appropriate methods and tools in order to extract insight from data.
	<b>Substantive expertise</b> in a scientific field is crucial for generating motivating questions and hypotheses and interpreting results.
	<b>Traditional research</b> lies at the intersection of knowledge of math and statistics with substantive expertise in a scientific field.
	<b>Machine learning</b> stems from combining hacking skills with math and statistics knowledge, but does not require scientific motivation.
	<b>Danger zone!</b> Hacking skills combined with substantive scientific expertise without rigorous methods can beget incorrect analyses.

Picture by Natalia Bilenko, Drew Conway, et al.

# RWTH Data Science Flower



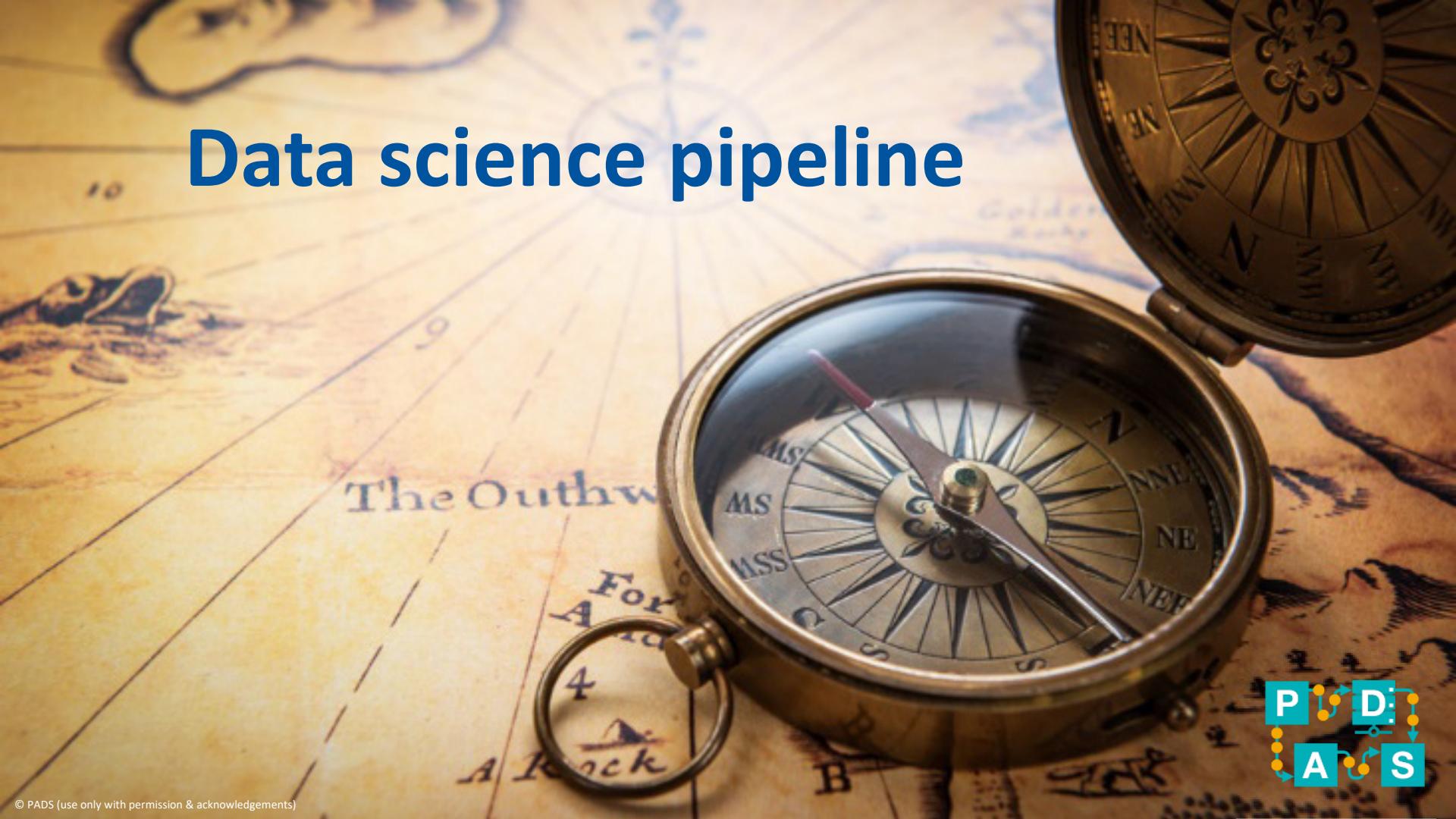
RWTH AACHEN  
UNIVERSITY

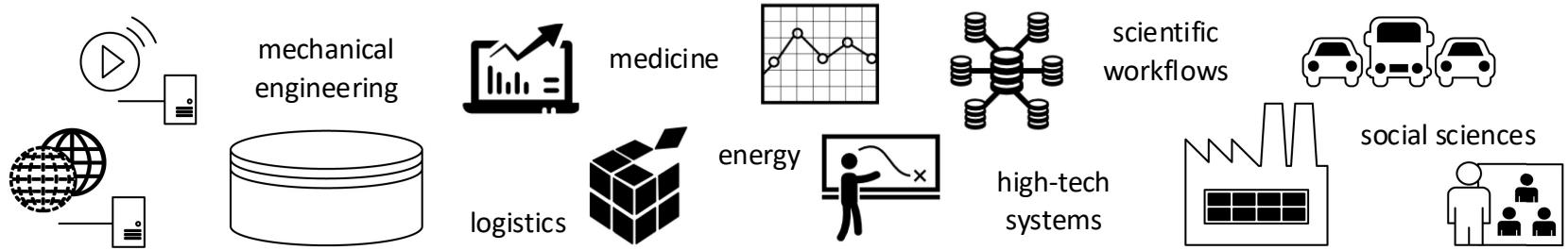


# It is not easy to turn data into value



# Data science pipeline





## infrastructure

“volume and velocity”

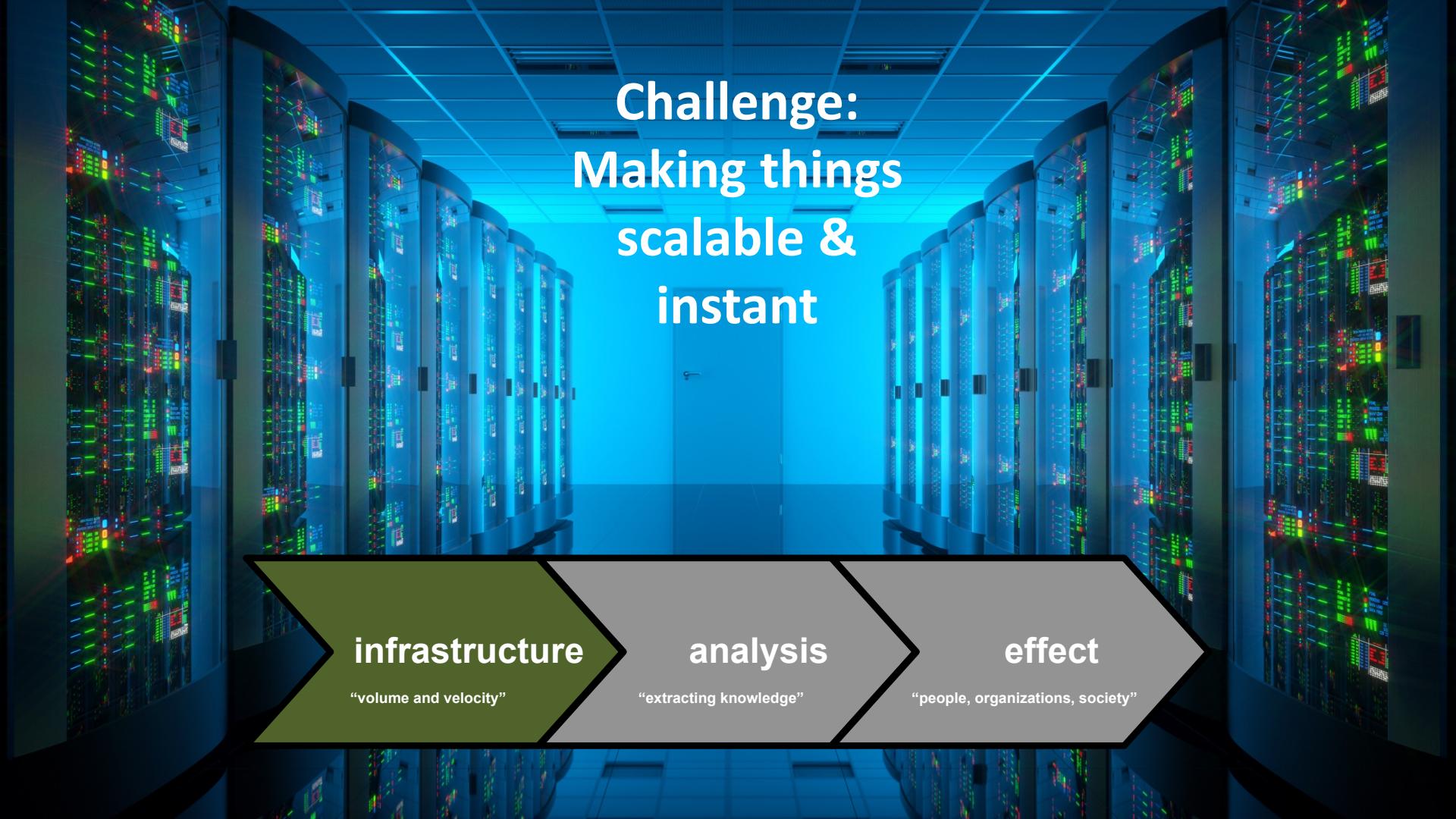
## analysis

“extracting knowledge”

## effect

“people, organizations, society”

- big data infrastructures
- distributed systems
- data engineering
- programming
- security
- ...
- statistics
- data/process mining
- machine learning
- artificial intelligence
- visualization
- ...
- ethics & privacy
- IT law
- operations management
- business models
- entrepreneurship
- ...



# Challenge: Making things scalable & instant

infrastructure

“volume and velocity”

analysis

“extracting knowledge”

effect

“people, organizations, society”

# Challenge: Providing answers to known and unknown unknowns



infrastructure

"volume and velocity"

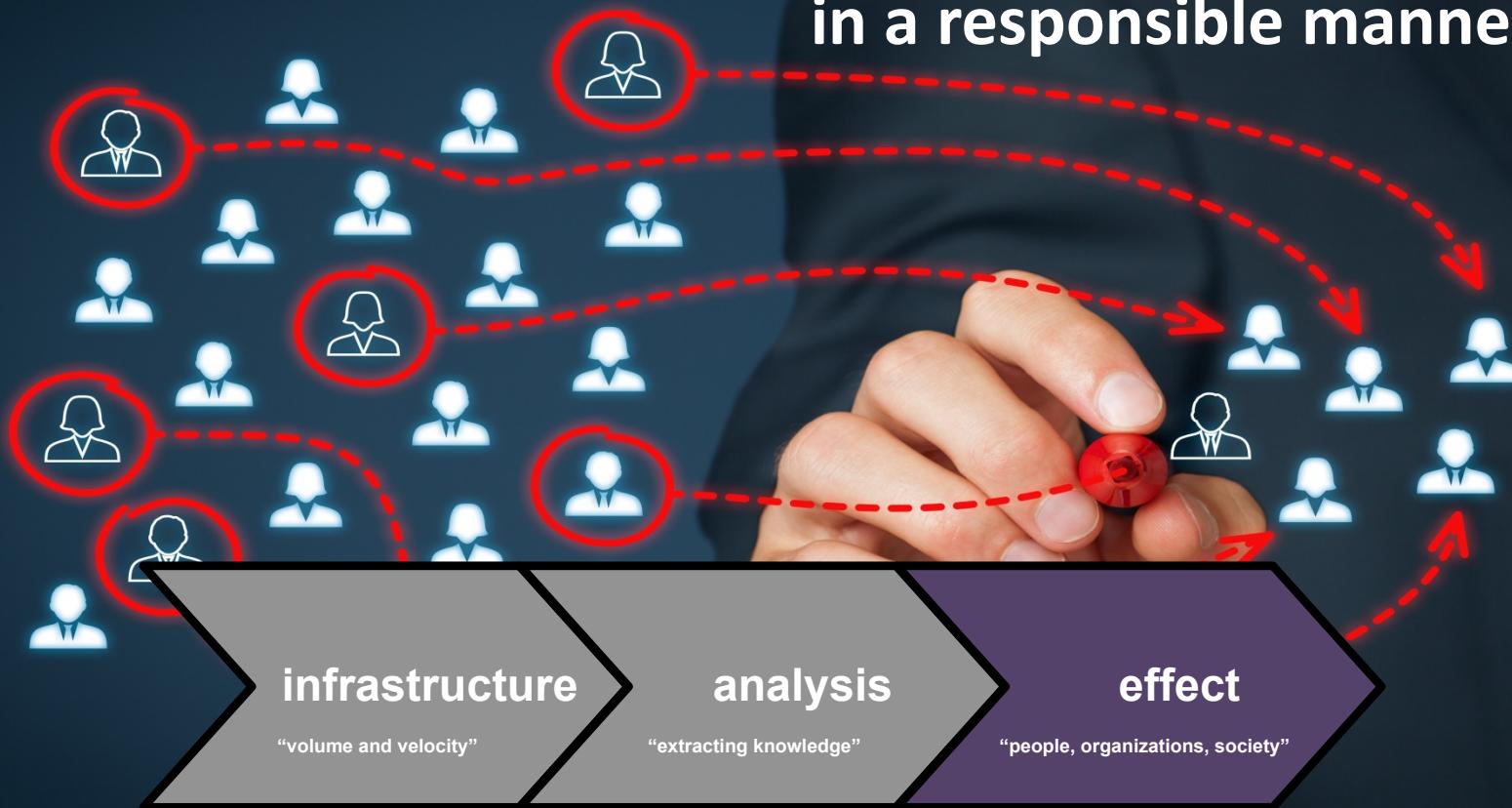
analysis

"extracting knowledge"

effect

"people, organizations, society"

# Challenge: Doing all of this in a responsible manner!



# Types of data



# Types of data

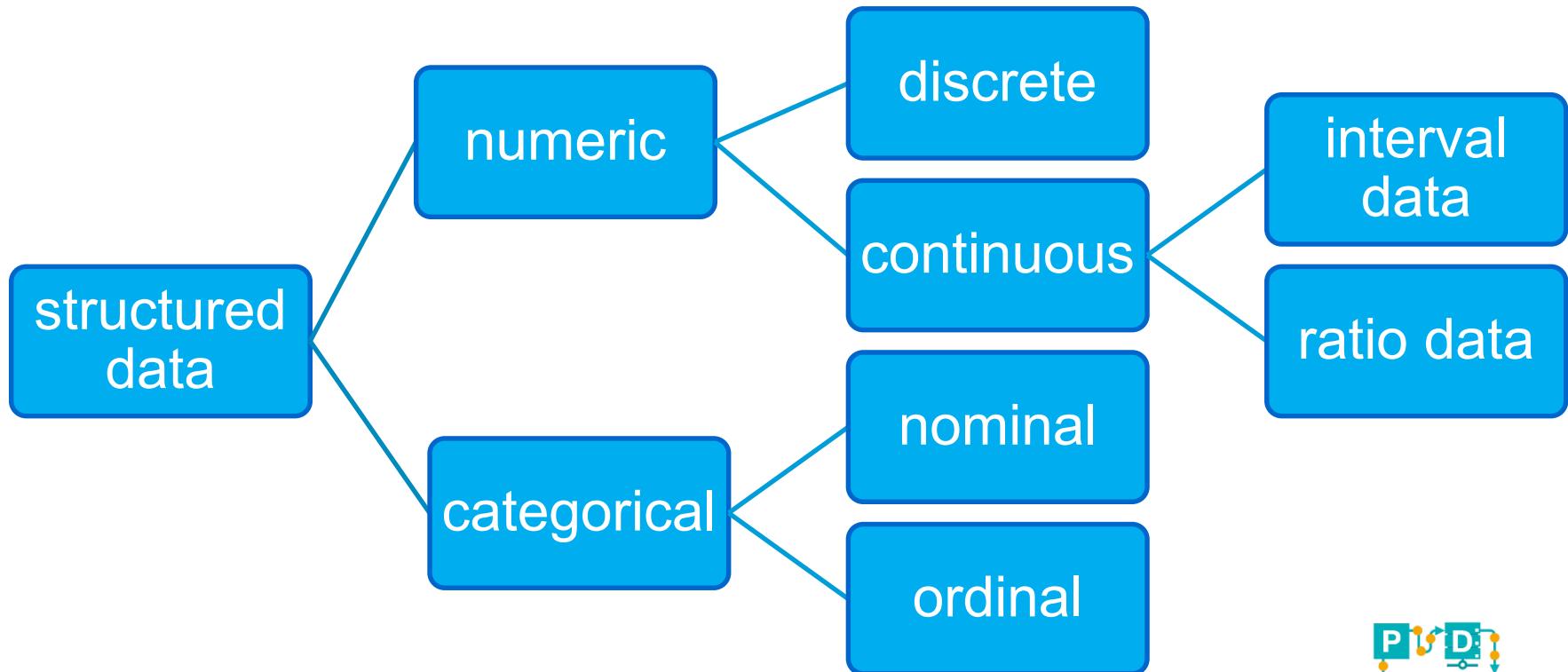
## Structured data

- Numerical data  
(age, time, temperature)
- Categorical data  
(gender, color, country, class)

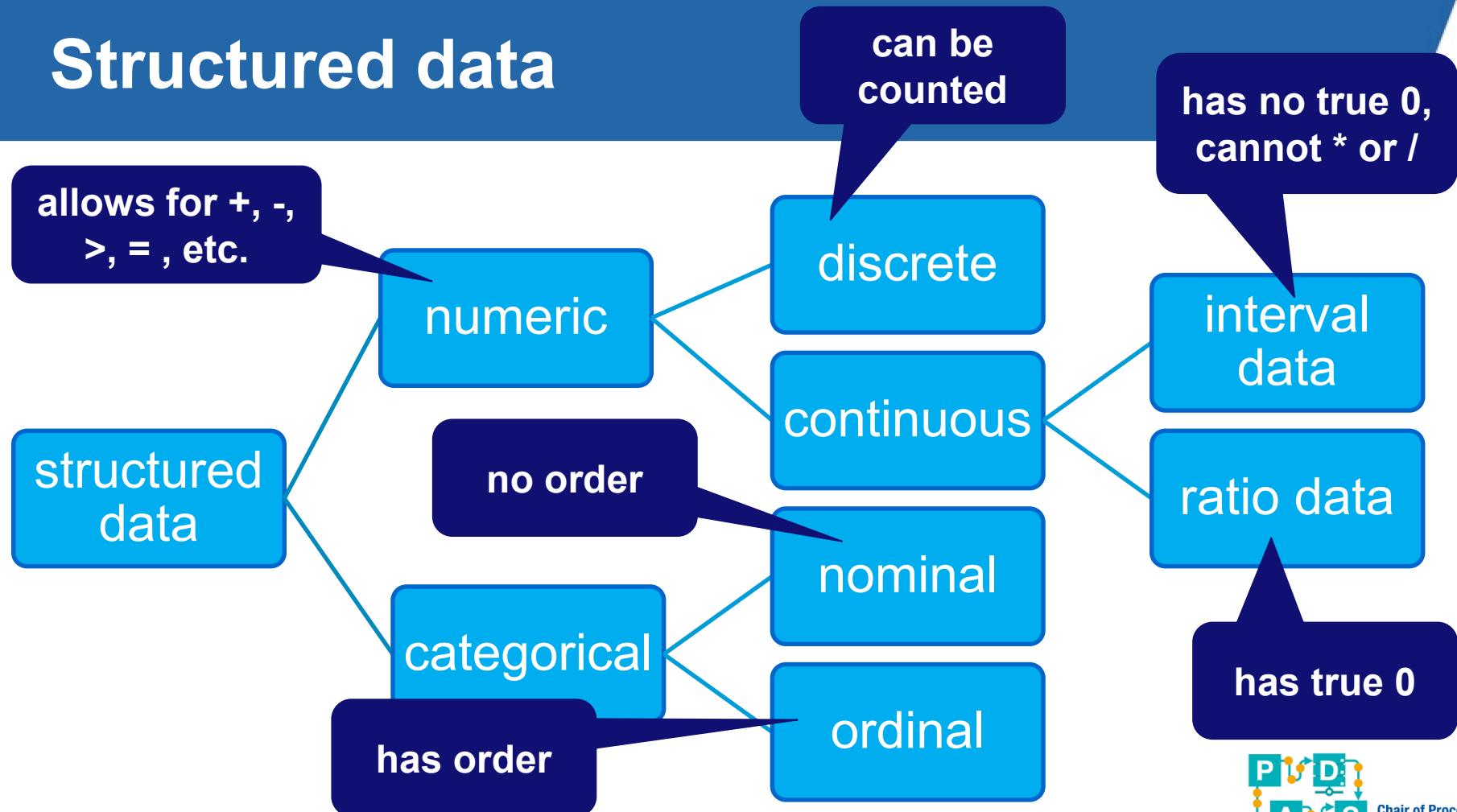
## Unstructured data

- Text
- Audio
- Image
- Signal
- Video
- etc.

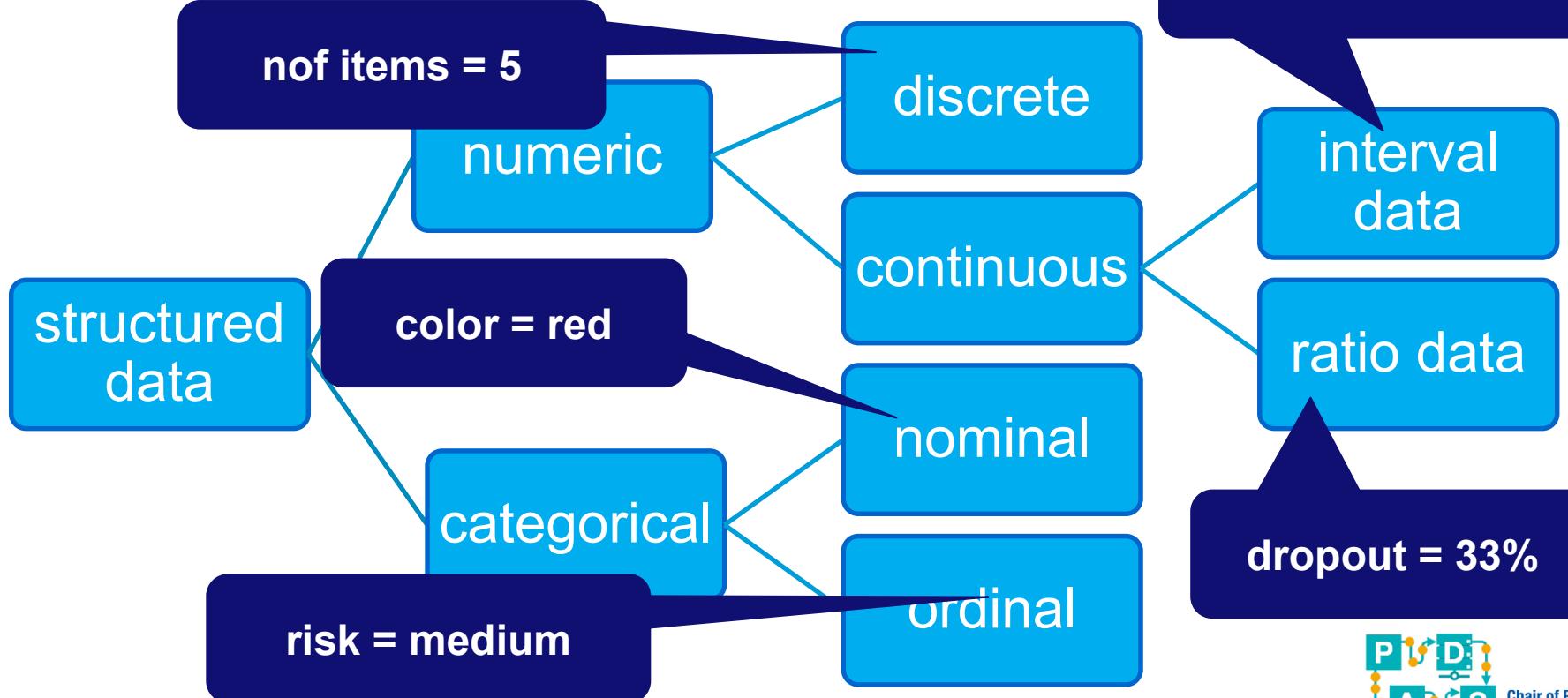
# Structured data



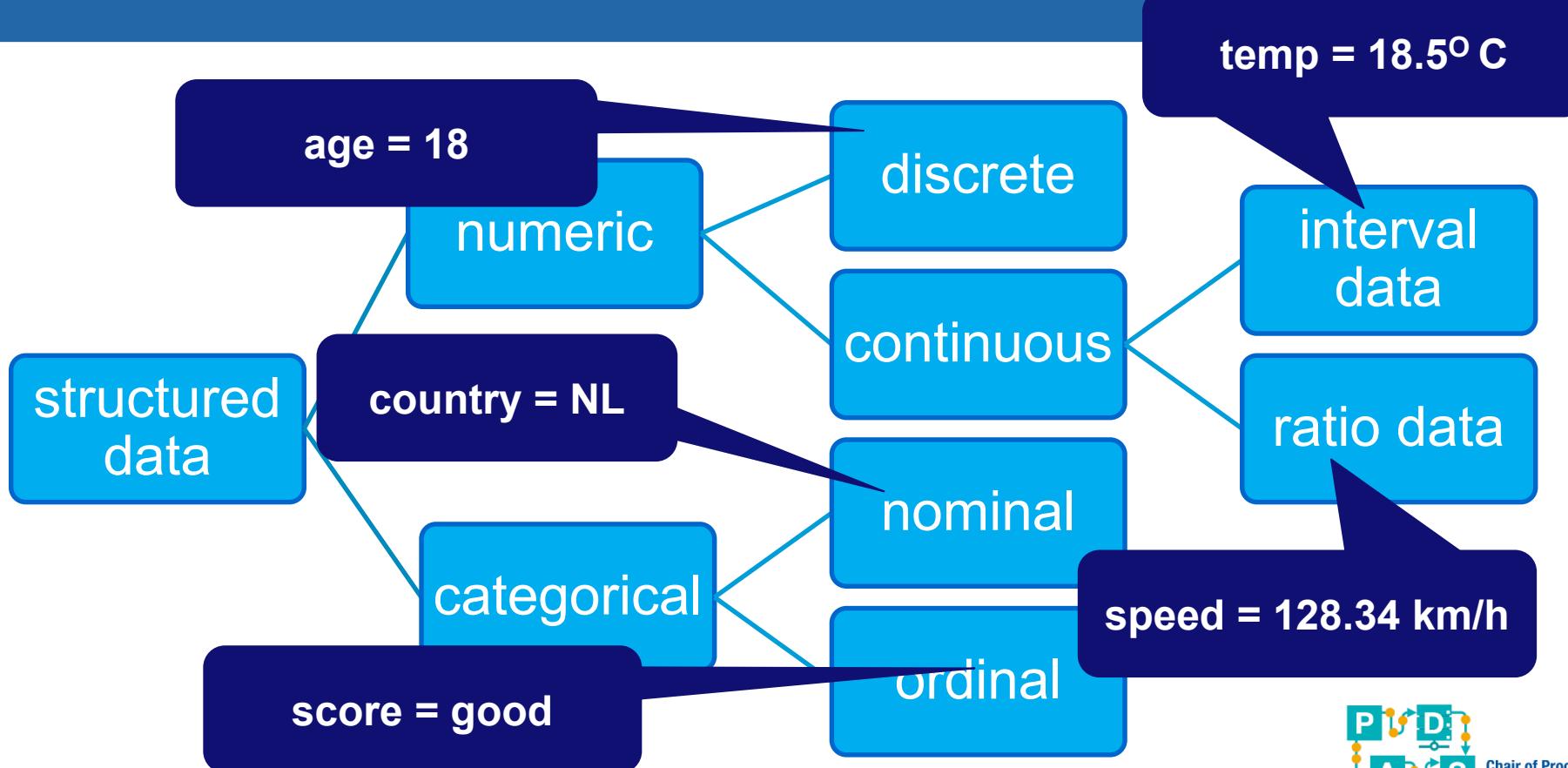
# Structured data



# Structured data

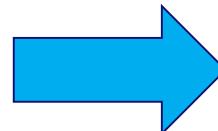


# Structured data

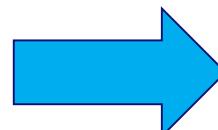
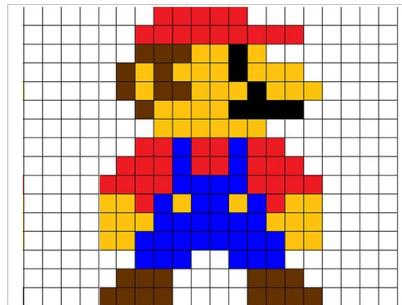


# Unstructured data

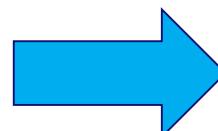
- Text
- Audio
- Image
- Signal
- Video
- etc.



1111100100010011101



010010010001011101



110010010111010101



Chair of Process  
and Data Science

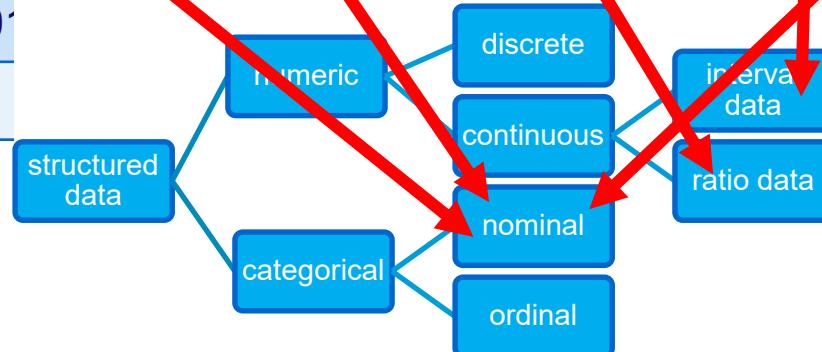
# Tabular data

order id	product	price	date	complaint
32424	718 Cayman	66.000	21-10-2018	no
34535	911 Carrera	102.000	22-10-2018	yes
43555	911 Turbo	154.000	24-10-2018	yes
64564	718 Cayman S	77.000	24-10-2018	no
23424	911 Targa	143.000	26-10-2018	yes
...	...	...	...	

columns are features and rows are instances

# Tabular data

order id	product	price	date	complaint
32424	718 Cayman	66.000	21-10-2018	no
34535	911 Carrera	102.000	22-10-2018	yes
43555	911 Turbo	154.000	24-10-2018	yes
64564	718 Cayman S	77.000	24-10-2018	no
23424	911 Carrera S	120.000	25-10-2018	yes
...				



# Tabular data

from	to	message	image	date
Sue	Pete	“How are you?”	😊	21-10-2018
Pete	Sue	“I’m busy!”	🎵	22-10-2018
Pete	Mary	“Let’s go out.”	♠	24-10-2018
Mary	Sue	“Pete joins us.”	☀️	24-10-2018
Mary	Kim	“We will go now.”	♫	26-10-2018
...	...	...	...	...

columns are features and rows are instances



# Features

- Features are **raw or derived** (max, min, average, rank, bin, etc.).
- **Time plays a special role:** time cannot decrease and often we want to predict the future based on the past.
- In case of **labeled data** there are **descriptive features** and a **target feature**.

# Labeled tabular data

order id	product	price	date	complaint
32424	718 Cayman	66.000	21-10-2018	no
34535	911 Carrera	102.000	22-10-2018	yes
43555	911 Turbo	154.000	24-10-2018	yes
64564	718 Cayman S	77.000	24-10-2018	no
23424	911 Targa	143.000	26-10-2018	yes
...	...	...	...	

target feature

descriptive features



Chair of Process  
and Data Science

# Labeled tabular data

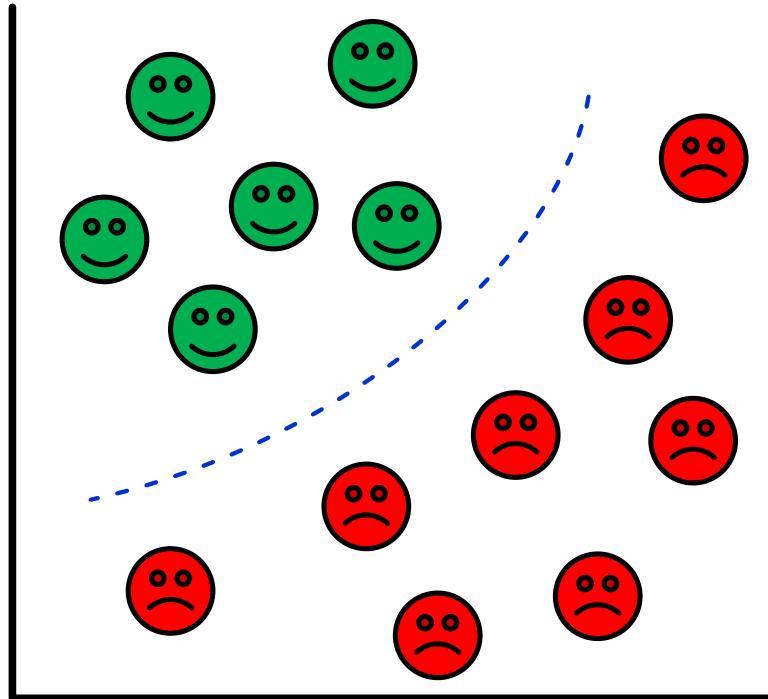
- Alternative names for descriptive features
  - predictor variables
  - independent variables
- Alternative names for target feature
  - response variable
  - dependent variable
- Alternative names for instances: individuals, entities, cases, objects, or records.

# Unlabeled tabular data

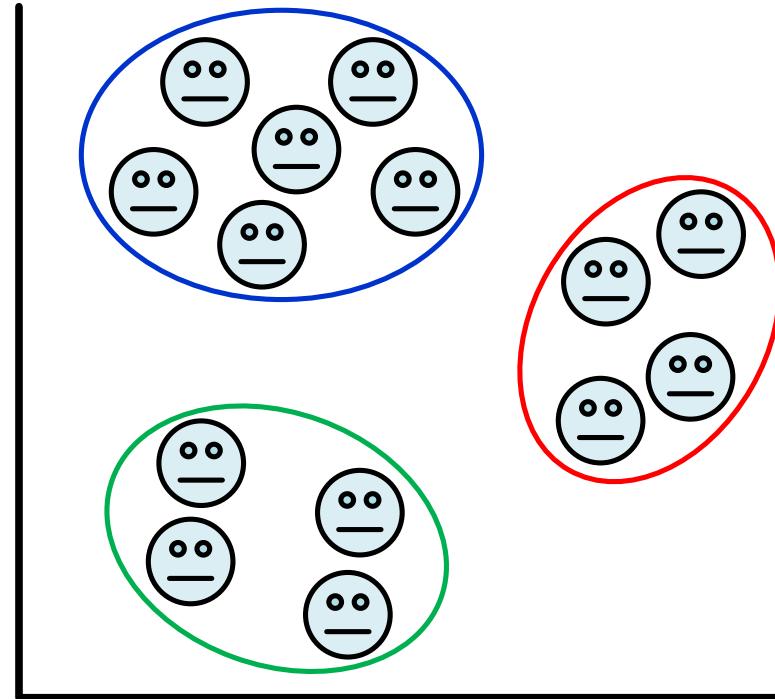
from	to	message	image	date
Sue	Pete	“How are you?”	😊	21-10-2018
Pete	Sue	“I’m busy!”	🎵	22-10-2018
Pete	Mary	“Let’s go out.”	♠	24-10-2018
Mary	Sue	“Pete joins us.”	☀️	24-10-2018
Mary	Kim	“We will go now.”	♫	26-10-2018
...	...	...	...	...

No target feature has been selected

# Supervised versus unsupervised learning

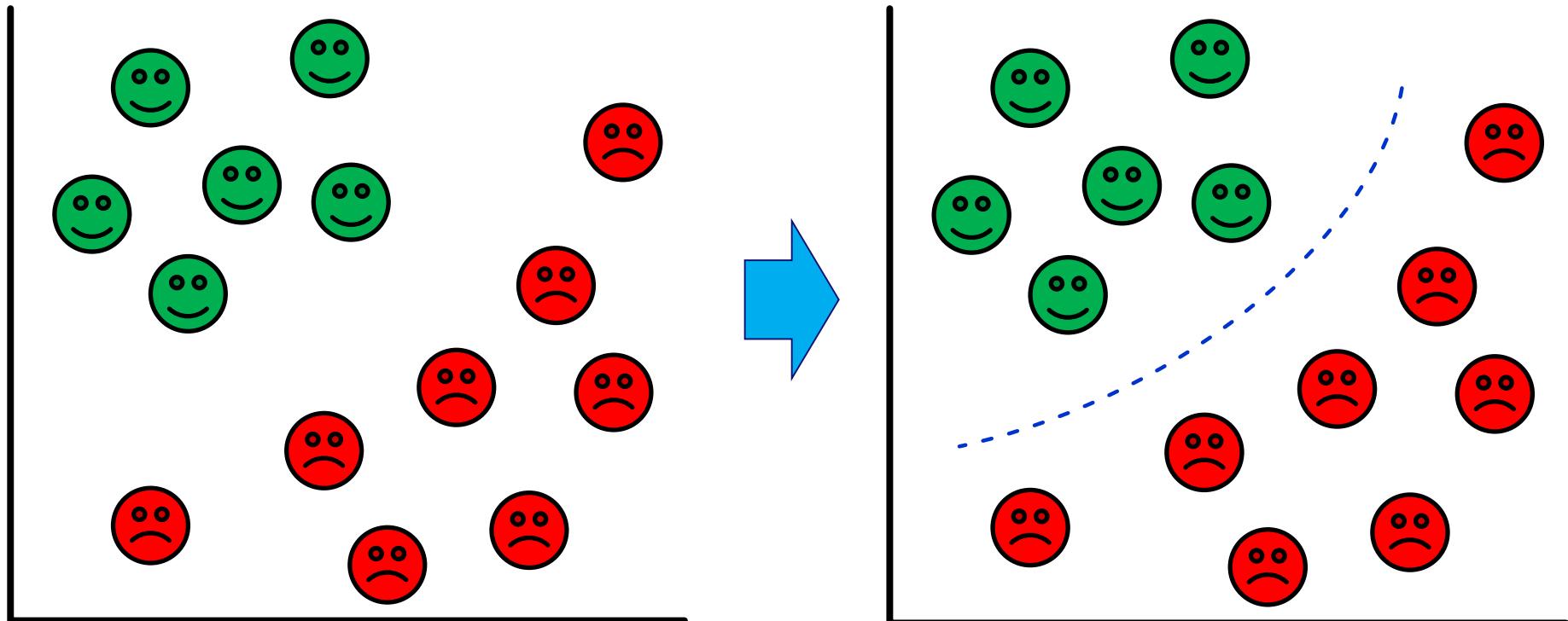


supervised



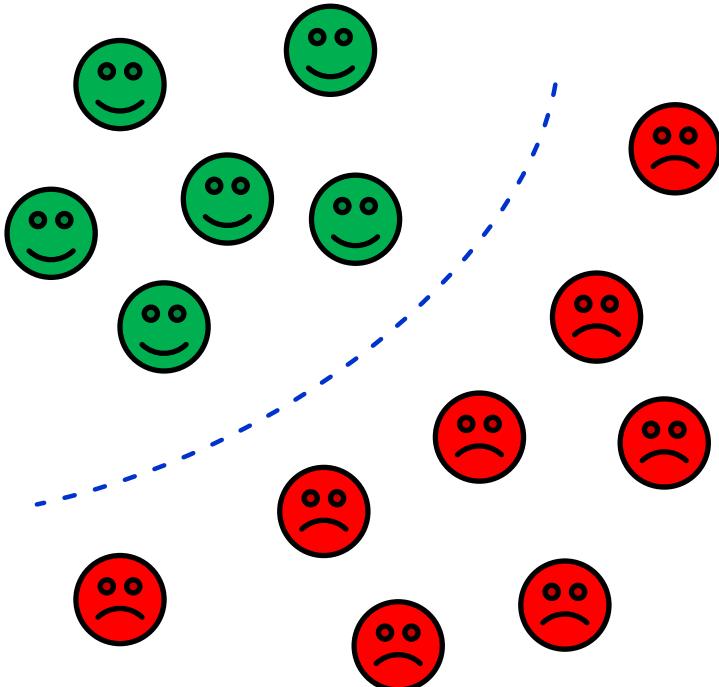
unsupervised

# Supervised learning using labeled data



All instances have a target label (here color).

# Supervised learning using labeled data

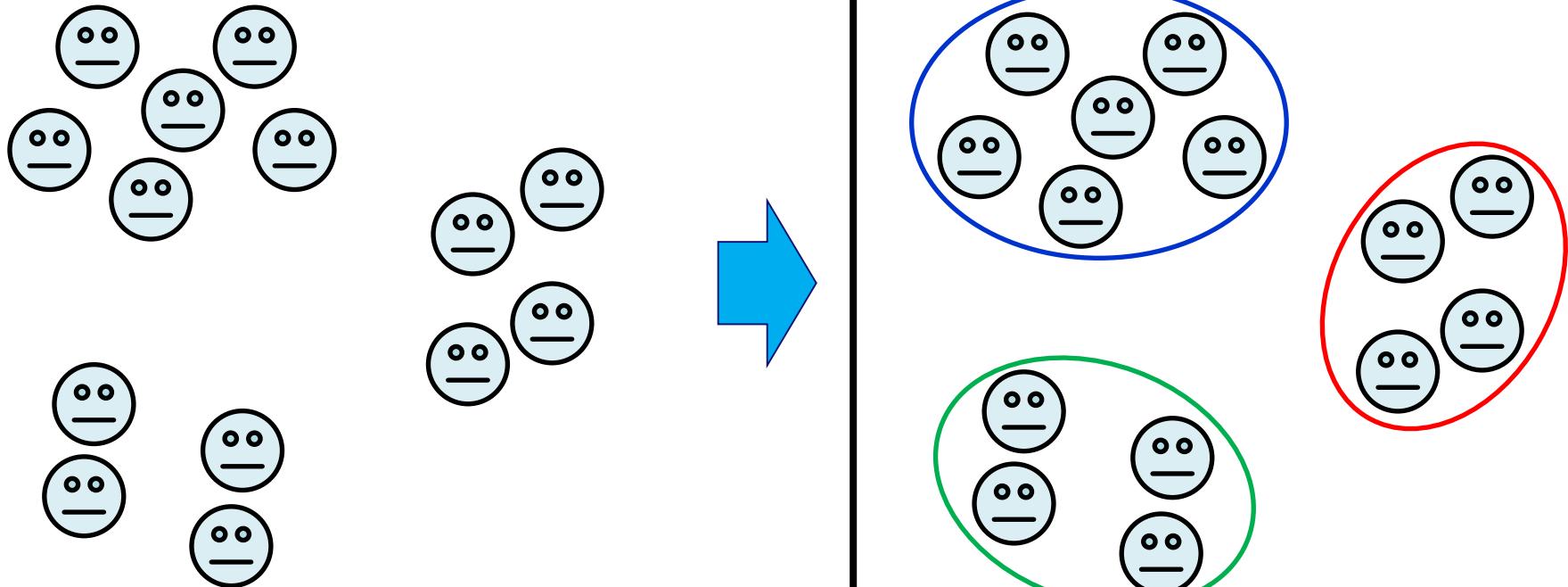


The goal is to find a “rule” in terms descriptive features that explains the target feature as good as possible.

# Examples of supervised learning

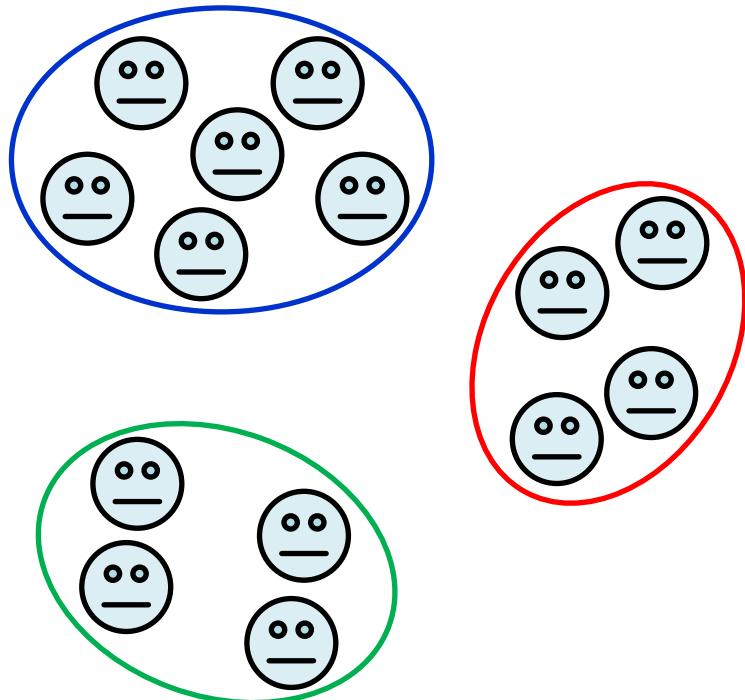
- Hospital:
  - **target variable = recover (Y/N)**
  - **descriptive variables = age, gender, smoking, ...**
- University
  - **target variable = drops out (Y/N)**
  - **descriptive variables = mentor, prior education, ...**
- Production
  - **target variable = order is delivered in time (Y/N)**
  - **descriptive variables = product, agent, ...**

# Unsupervised learning



Instances do not have a target label.

# Unsupervised learning



The goal is to find clusters or patterns.

- Clusters are homogeneous sets of instances.
- Patterns reveal hidden structures in the data (unknown unknowns).

# Examples unsupervised learning

- Find similar groups of patients, students, customers, orders, cars, companies, etc.
- Find rules of the form (unknown unknowns):
  - Customers that buy bread and butter typically pay by cash.
  - Patients that drink and smoke typically pay the hospital bill earlier than others.
  - Products produced by team A on Monday tend to be returned more frequently by the customer.

# Terminology



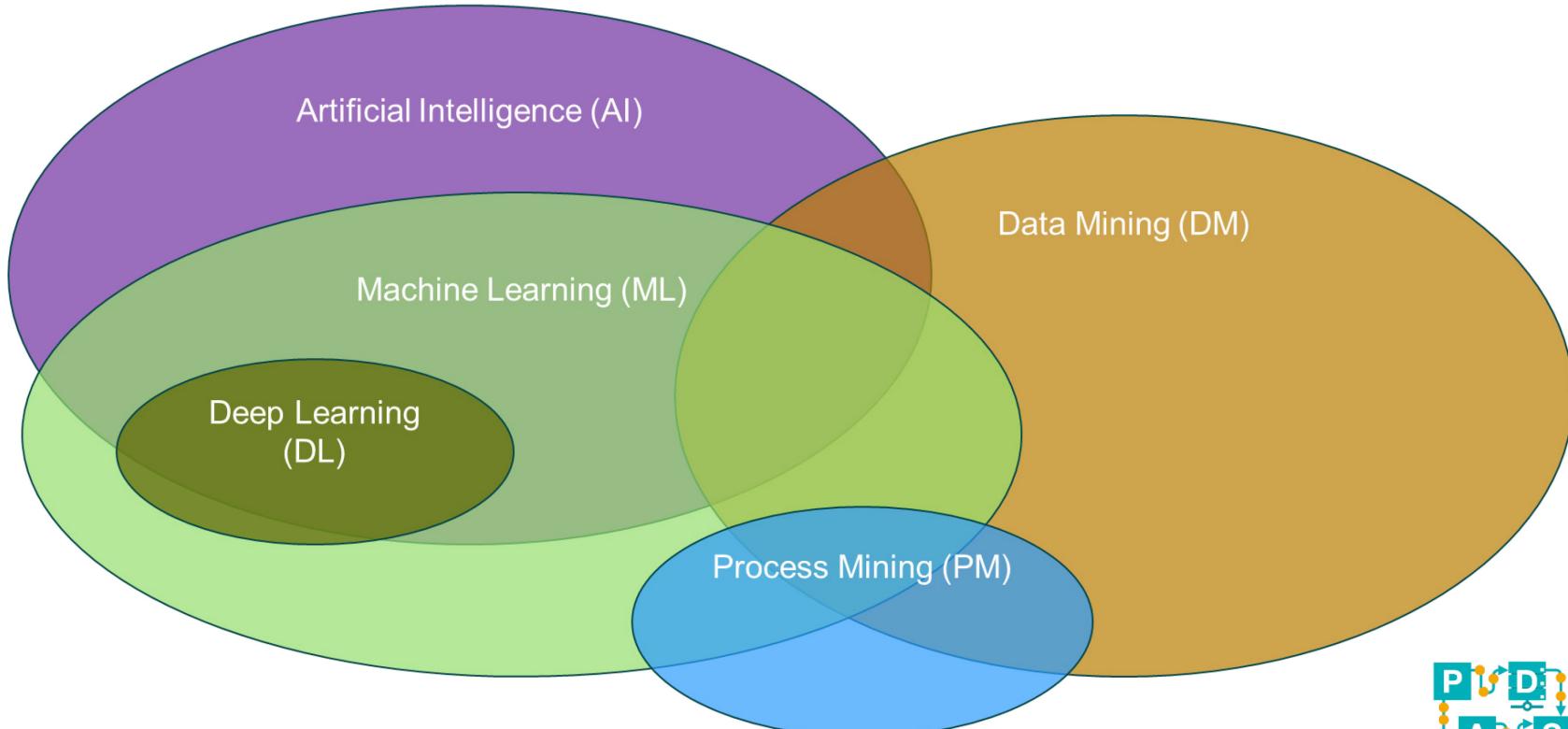
# Terminology

- **Many different names** (statistics, data analytics, data mining, machine learning, artificial intelligence, predictive analytics, process mining, etc.) **are used to refer to the key disciplines that contribute to data science.**
- **Unfortunately, the areas these names describe are heavily overlapping and context dependent.**

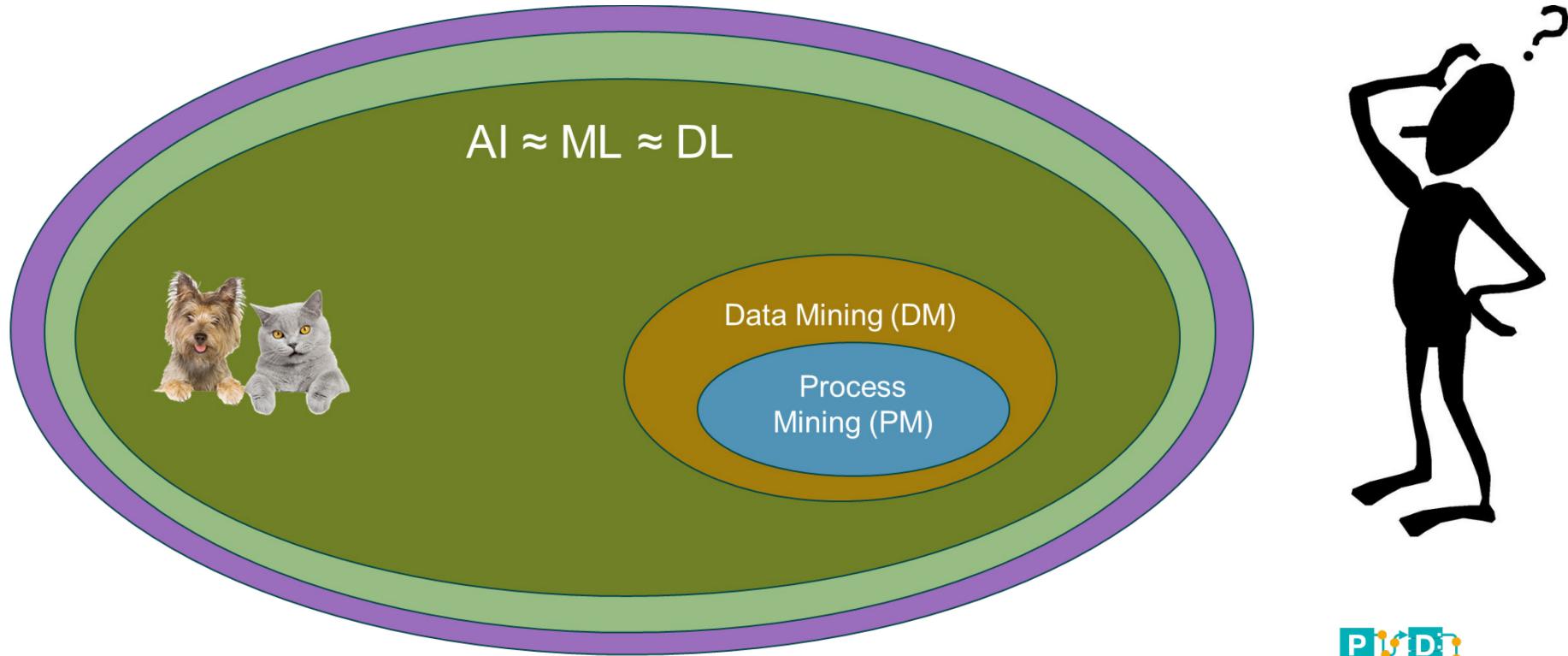
# Example: Scoping machine learning

- Sometimes **machine learning** is used as a synonym for **deep neural networks** and sometimes to cover the **entire spectrum of learning techniques** (including mining, etc.).
- The fact that a neural network can be used as a classifier does not imply that the numerous classification techniques developed in data mining are part of machine learning (in the narrow sense).

# A more balanced view



# Perception by outsiders



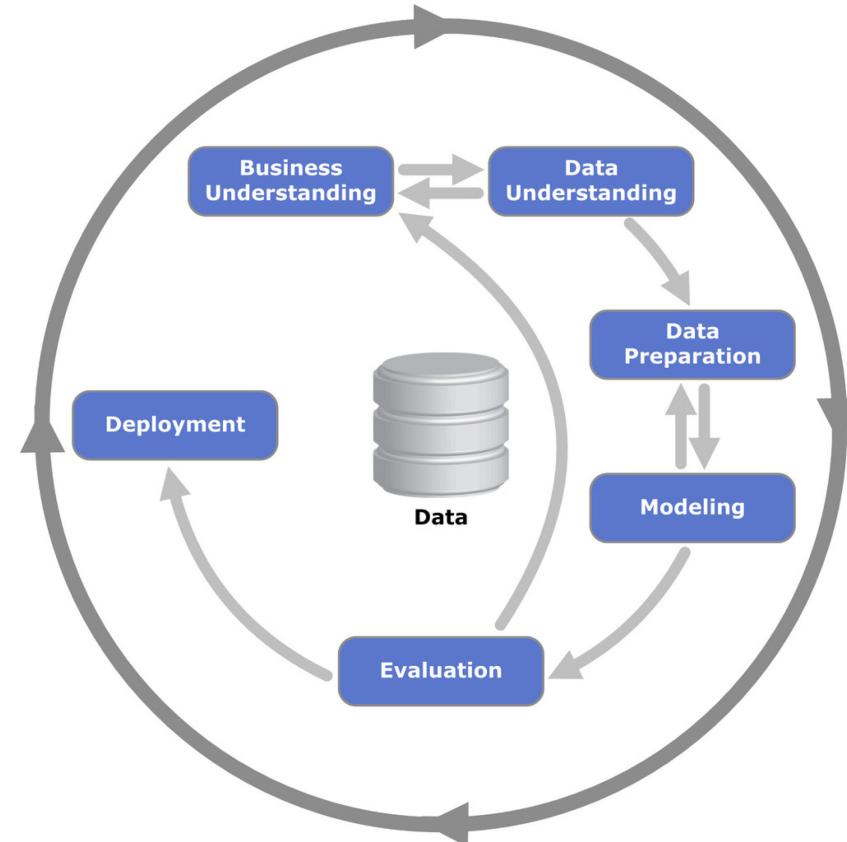
# Data science process



# CRISP-DM

## Cross-industry standard process for data mining

- CRISP-DM was developed in the late 1990-ties involving SPSS, Teradata, Daimler AG, NCR Corporation and Ohra.
- Quite obvious

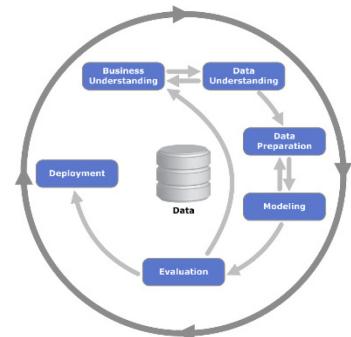


# CRISP-DM

## Cross-industry standard process for data mining

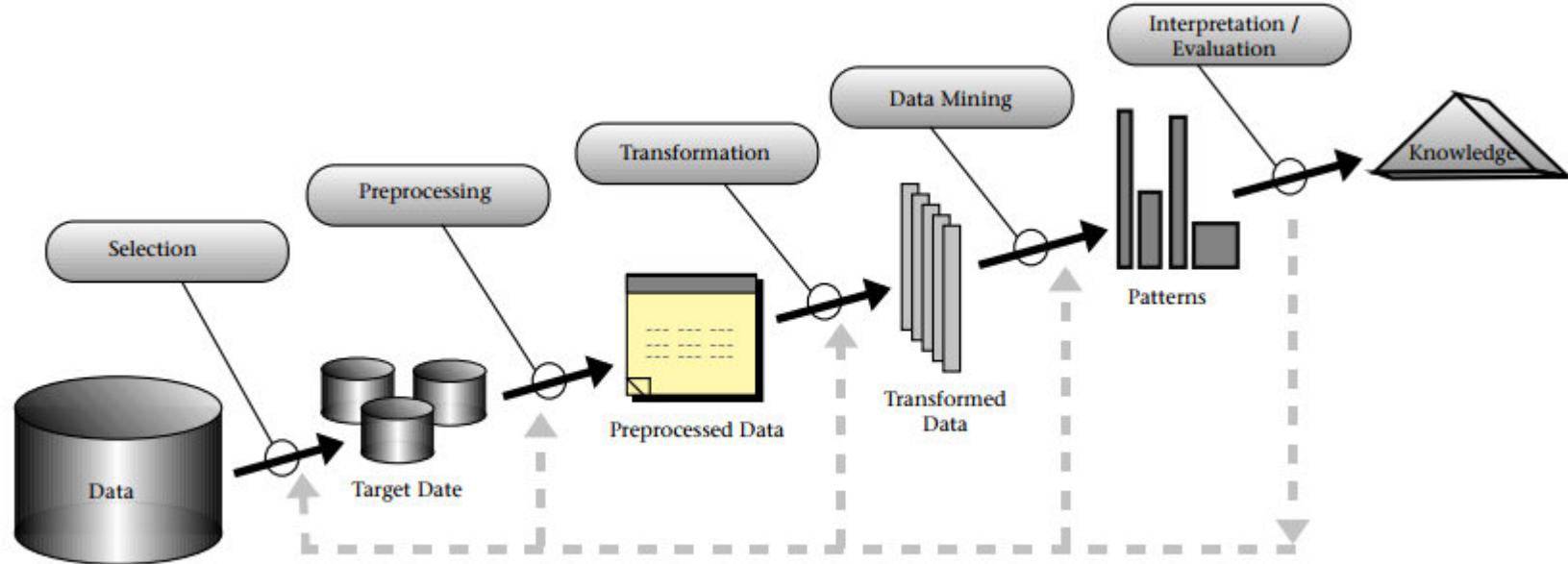
Business Understanding	Data Understanding	Data Preparation	Modeling	Evaluation	Deployment
Determine Business Objectives <i>Background</i> <i>Business Objectives</i> <i>Business Success Criteria</i>	Collect Initial Data <i>Initial Data Collection Report</i>	Data Set <i>Data Set Description</i>	Select Modeling Technique <i>Modeling Technique</i> <i>Modeling Assumptions</i>	Evaluate Results <i>Assessment of Data Mining Results w.r.t. Business Success Criteria</i> <i>Approved Models</i>	Plan Deployment <i>Deployment Plan</i>
Situation Assessment <i>Inventory of Resources Requirements, Assumptions, and Constraints</i> <i>Risks and Contingencies</i> <i>Terminology</i> <i>Costs and Benefits</i>	Describe Data <i>Data Description Report</i>	Select Data <i>Rationale for Inclusion / Exclusion</i>	Generate Test Design <i>Test Design</i>	Review Process <i>Review of Process</i>	Plan Monitoring and Maintenance <i>Monitoring and Maintenance Plan</i>
Determine Data Mining Goal <i>Data Mining Goals</i> <i>Data Mining Success Criteria</i>	Explore Data <i>Data Exploration Report</i>	Clean Data <i>Data Cleaning Report</i>	Build Model <i>Parameter Settings</i> <i>Models</i> <i>Model Description</i>	Determine Next Steps <i>List of Possible Actions</i> <i>Decision</i>	Produce Final Report <i>Final Report</i> <i>Final Presentation</i>
Produce Project Plan <i>Project Plan</i> <i>Initial Assessment of Tools and Techniques</i>	Verify Data Quality <i>Data Quality Report</i>	Construct Data <i>Derived Attributes</i> <i>Generated Records</i>	Integrate Data <i>Merged Data</i>		Review Project Experience Documentation
		Format Data <i>Reformatted Data</i>			
<p>The term “modeling” is sometimes a bit misleading: selection and assumptions (human) and automated learning by tool/algorithim are combined.</p>					

Taken from Pete Chapman (1999) The CRISP-DM User Guide.



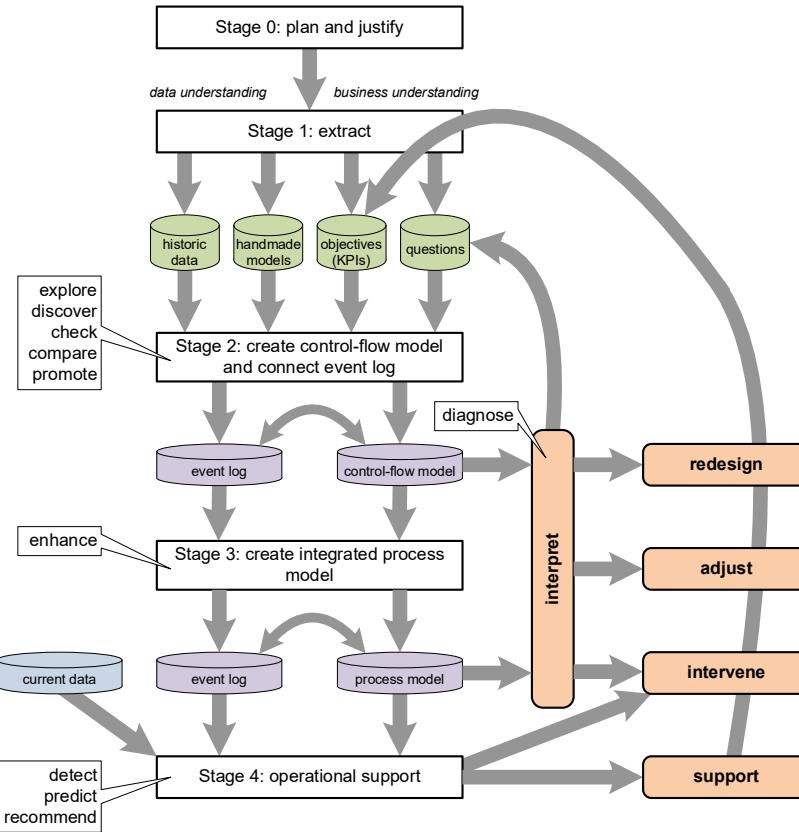
# KDD (Knowledge Discovery in Databases) Process

(by Fayyad, Piatetsky-Shapiro, and Smyth)

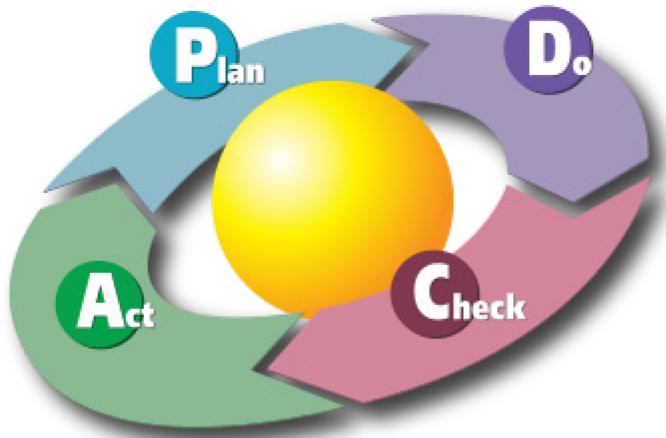


**SEMMA (with the phases Sample, Explore, Modify, Model, and Assess) is another process model developed by SAS Institute.**

# L\* lifecycle model (specific for process mining)



# Related to PDCA and DMAIC methodology

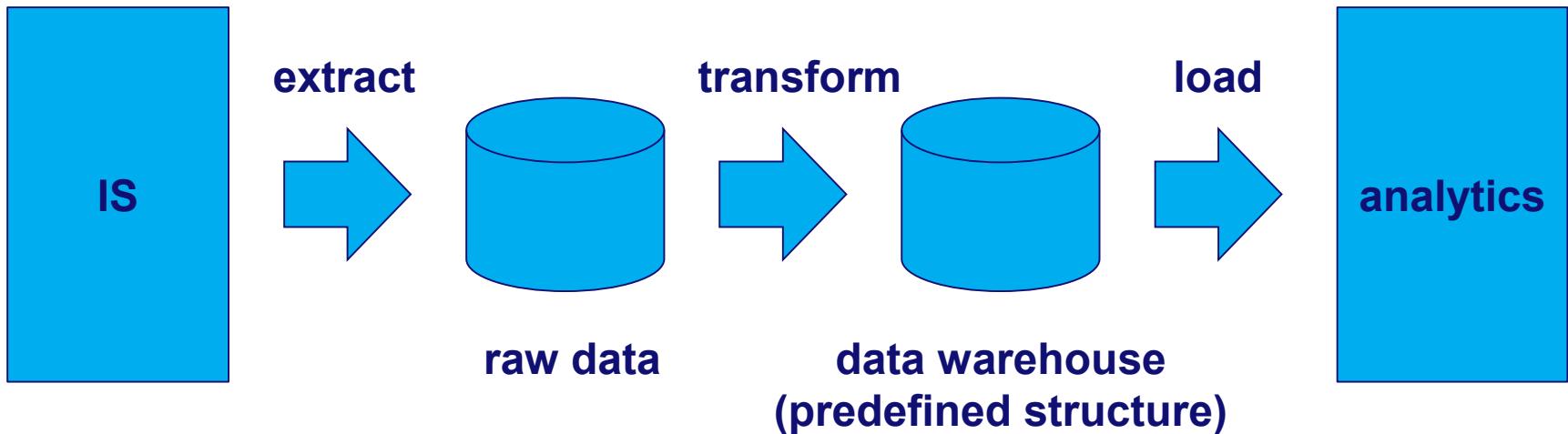


**PDCA (Plan–Do–Check–Act)**  
methodology by William Deming

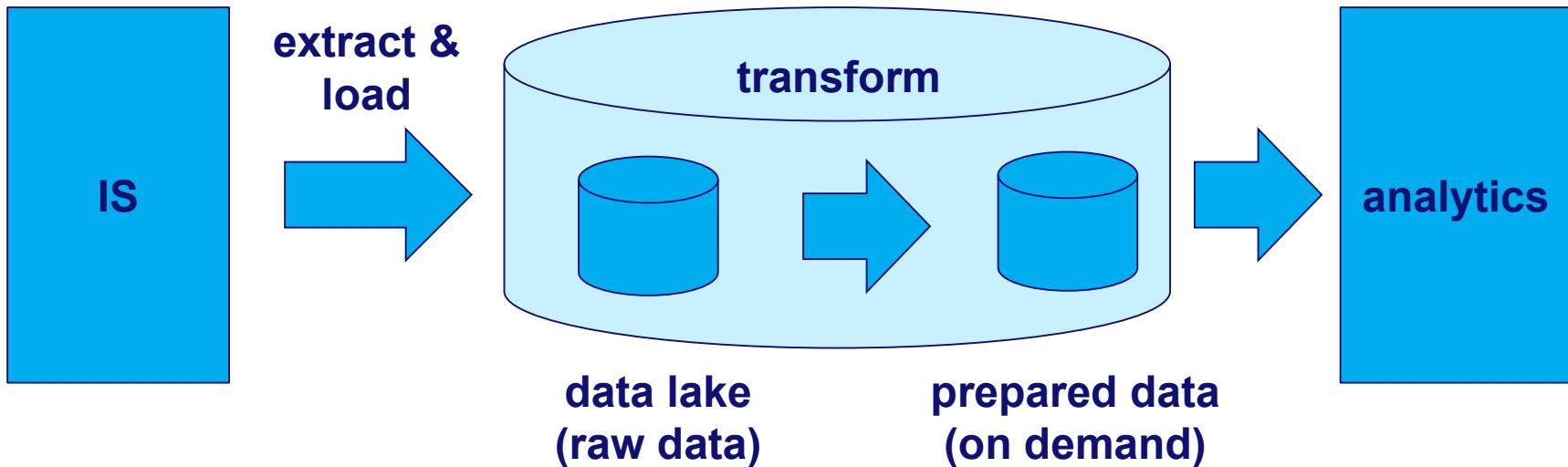


**DMAIC (Define, Measure, Analyze, Improve and Control) methodology used in Six Sigma projects**

# Extract, Transform, Load (ETL)



# Extract, Load, Transform (ELT)



A wide-angle photograph of a deep blue lake nestled among towering, forested mountains. The water is calm, reflecting the surrounding greenery and the clear, light blue sky above. In the foreground, the tops of several green trees are visible. On the left side of the lake, there's a small, rocky peninsula with some greenery. In the bottom left corner, a few buildings with brown roofs are visible, partially obscured by trees. The overall scene is a peaceful, natural landscape.

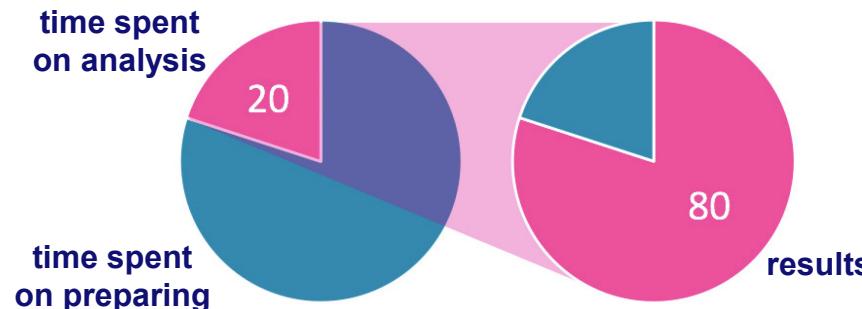
# data lakes (ELT)

A large warehouse interior is filled with numerous stacks of clear plastic water bottles. The bottles are arranged in long, low rows that recede into the distance. In the background, a red Toyota forklift is positioned near a white pillar. To the left, there are pallets of water bottles wrapped in blue plastic shrink-wrap. The floor is a polished concrete surface.

data warehouses  
(ETL)

# Another 80/20 rule

- 80 percent of a data scientist's time is spent on finding, cleaning, preprocessing, and organizing data, leaving only 20 percent to actually perform analysis.
- However, the 20 percent effort determines 80% of the results.



# Challenges



# Finding data

- There may be hundreds or thousands of tables.
- There many be many different entities that are less relevant.

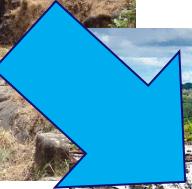


# Transforming data

- Reorganizing data.
- Extracting relevant features.



# Dealing with Big data



# Dealing with streaming data



# Data quality

- Data may be incomplete, invalid, inconsistent, imprecise, and/or outdated.
- Example timestamps:
  - Incomplete (missing event)
  - Invalid (14-14-2018)
  - Inconsistent (14-7-2018 => 7-14-2018)
  - Imprecise (2018-09-21'T'13:20:10)

# Overfitting / Underfitting

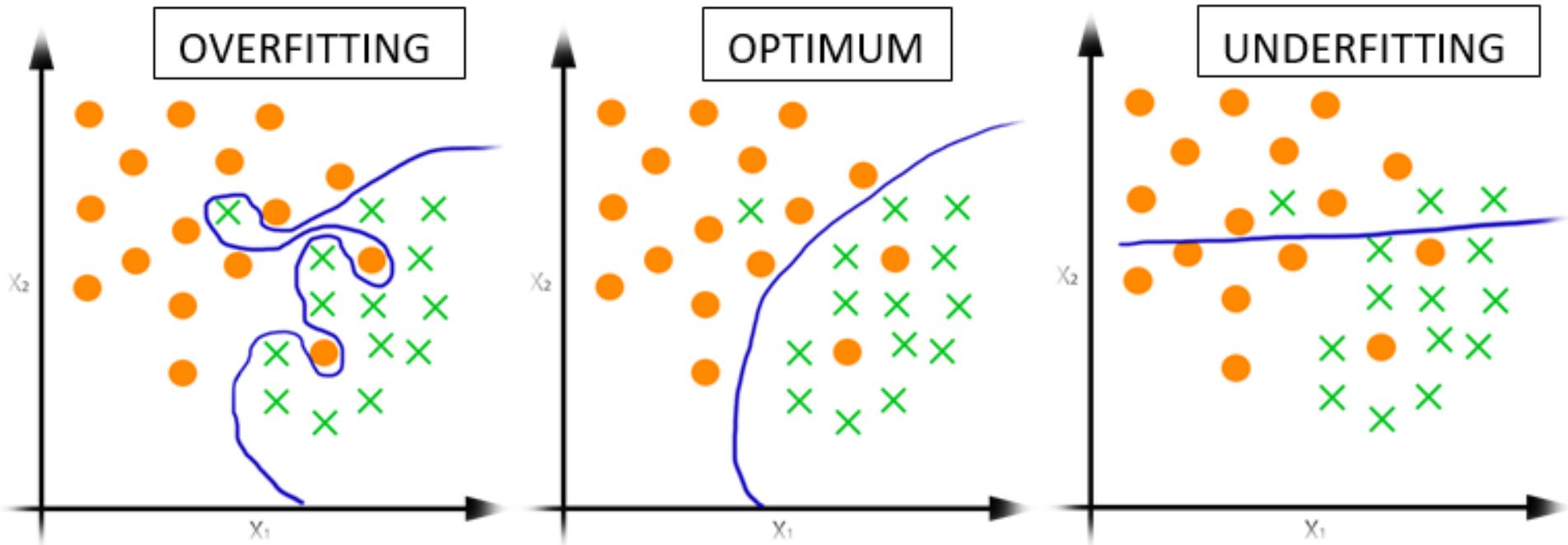
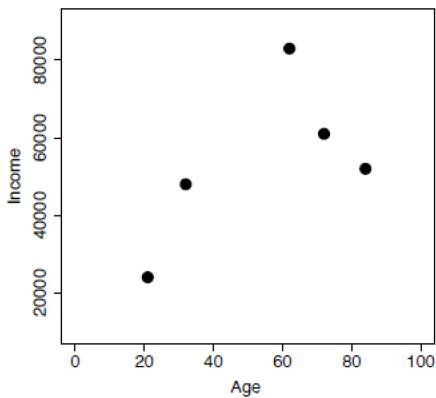
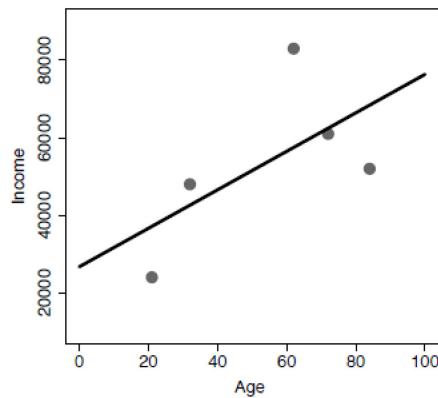


Diagram by Sachin Joglekar (Google).

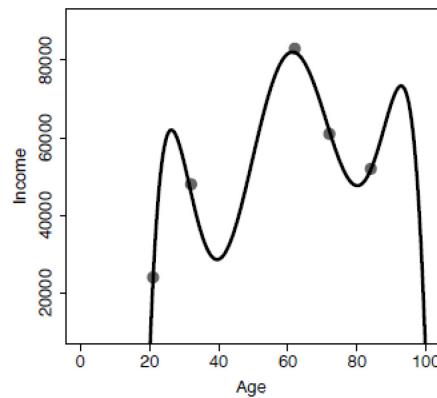
# Overfitting / Underfitting



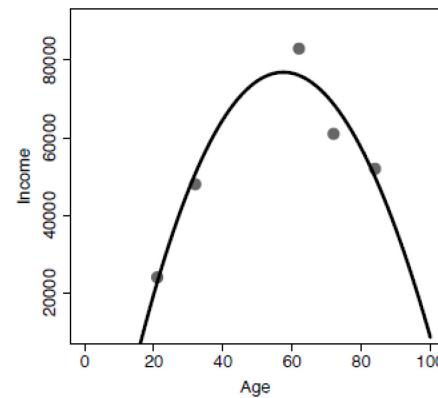
(a) Dataset



(b) Underfitting



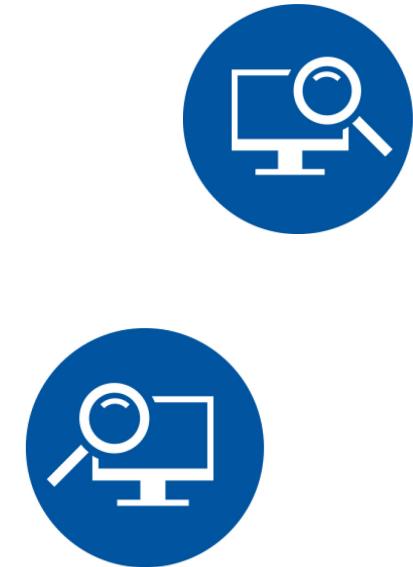
(c) Overfitting



(d) Just right

Diagram taken from Fundamentals of Machine Learning for Predictive Data Analytics by J. Kelleher, B. Mac Namee and A. D'Arcy.

# Dealing with concept drift



# Making results actionable

- **Analysis results need to be relevant, specific, novel and clear.**



# Ensuring fairness

**Fairness:** Data  
Science without  
prejudice: How to  
avoid unfair  
conclusions even  
if they are true?



Chair of Process  
and Data Science

# Ensuring accuracy

**Accuracy:** Data  
Science without  
guesswork: How to  
answer questions  
with a guaranteed  
level of accuracy?



# Ensuring confidentiality

**Confidentiality:**  
Data Science that  
ensures  
confidentiality:  
How to answer  
questions without  
revealing  
secrets?



# Ensuring transparency

**Transparency:**  
Data Science that  
provides  
transparency: How  
to clarify answers  
such that they  
become  
indisputable?



# III-posed problems

- A problem is well-posed if
  - a solution exists and
  - the solution is unique.
- Problems in data science are often ill-posed:
  - there may be many possible models explaining observed phenomena,
  - the (training) data set is just a sample,
  - there may be noise (exceptional or incorrectly recorded instances) in the data set, and
  - the result needs to generalize to have predictive/explanatory value.



# Conclusion



# Conclusion

- **Data science**
  - has many faces,
  - is “super relevant”, and
  - provides many challenges.
- **Next lecture:**
  - Crash course in Python

#	Lecture	date	day
	<b>Lecture 1</b> Introduction	10/10/2018	Wednesday
	<b>Lecture 2</b> Crash Course in Python	11/10/2018	Thursday
<i>Instruction 1</i>	<i>Python</i>	<i>12/10/2018</i>	<i>Friday</i>
	<b>Lecture 3</b> Basic data visualisation/exploration	17/10/2018	Wednesday
	<b>Lecture 4</b> Decision trees	18/10/2018	Thursday
<i>Instruction 2</i>	<i>Decision trees and data visualization/exploration</i>	<i>19/10/2018</i>	<i>Friday</i>
	<b>Lecture 5</b> Regression	24/10/2018	Wednesday
	<b>Lecture 6</b> Support vector machines	25/10/2018	Thursday
<i>Instruction 3</i>	<i>Regression and support vector machines</i>	<i>26/10/2018</i>	<i>Friday</i>
	<b>Lecture 7</b> Neural networks (1/2)	31/10/2018	Wednesday
<i>Instruction 4</i>	<i>Neural networks and supervised learning</i>	<i>02/11/2018</i>	<i>Friday</i>
	<b>Lecture 8</b> Neural networks (2/2)	07/11/2018	Wednesday
	<b>Lecture 9</b> Evaluation of supervised learning problems	08/11/2018	Thursday
<i>Instruction 5</i>	<i>Neural networks and supervised learning</i>	<i>09/11/2018</i>	<i>Friday</i>
	<b>Lecture 10</b> Clustering	14/11/2018	Wednesday
	<b>Lecture 11</b> Frequent items sets	15/11/2018	Thursday
	<b>Lecture 12</b> Association rules	21/11/2018	Wednesday
	<b>Lecture 13</b> Sequence mining	22/11/2018	Thursday
<i>Instruction 6</i>	<i>Clustering, frequent items sets, association rules</i>	<i>23/11/2018</i>	<i>Friday</i>
	<b>Lecture 14</b> Process mining (unsupervised)	28/11/2018	Wednesday
	<b>Lecture 15</b> Process mining (supervised)	29/11/2018	Thursday

<i>Instruct</i>	<b>Lecture 1</b> Introduction	10/10/2018	Wednesday
<i>Instruct</i>	<b>Lecture 2</b> Crash Course in Python	11/10/2018	Thursday
<i>Instruct</i>	<b>Instruction 1</b>	<i>Python</i>	<i>12/10/2018</i>
<i>Instruct</i>	<b>Lecture 3</b> Basic data visualisation/exploration	17/10/2018	Wednesday
<i>Instruct</i>	<b>Lecture 4</b> Decision trees	18/10/2018	Thursday
<i>Instruct</i>	<b>Instruction 2</b>	<i>Decision trees and data visualization/exploration</i>	<i>19/10/2018</i>

	<b>Lecture 23</b> Big data (2/2)	17/01/2019	Thursday
<i>Instruction 11</i>	<i>Big data</i>	<i>18/01/2019</i>	<i>Friday</i>
	<b>Lecture 24</b> Closing	23/01/2019	Wednesday
	backup	24/01/2019	Thursday
<i>Instruction 12</i>	<i>Example exam questions</i>	<i>25/01/2018</i>	<i>Friday</i>
	backup	30/01/2019	Wednesday
	backup	31/01/2019	Thursday
extra	<i>Question hour</i>	<i>01/02/2019</i>	<i>Friday</i>