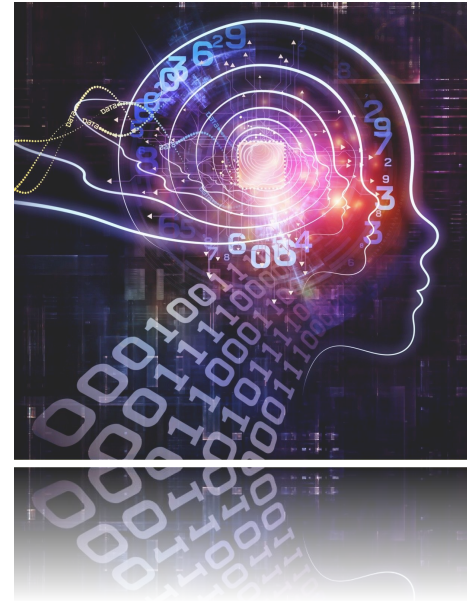# Data Quality for AI

Prof. Dr. Felix Naumann and Dr. Hazar Harmouch

WS 2021/2022

# Agenda

- ❏ Chair Introduction

- ❏ Organizational Information

- ❏ Data quality and AI

- ❏ Your Tasks

# Agenda

❏ **Chair Introduction**

❏ Organizational Information

❏ Data quality and AI

❏ Your Tasks

# Information Systems Team



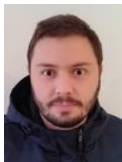Anna **Zobel**

Prof. Felix **Naumann**

Dr. Hazar **Harmouch**

Phillip **Wenig**

Leon **Bornemann**

Sebastian **Schmidl**

Tobias **Bleifuß**

Alejandro **Sierra-Múnera**

Nitisha **Jain**

Mazhar **Hameed**

Lan **Jiang**

Gerardo **Vitagliano**

Duplicate Detection
Data Change
Data Fusion
project **AKITA**
Entity Search
Data Profiling
Information Integration
project **DataKnoller**
Web Science
Data Scrubbing
project **AI4ART**
Data as a Service
Information Quality
Data Cleansing
Text Mining
Dependency Detection
Linked Open Data
CSV parsing
Knowledge Management for the Arts
Web Data
Distributed Computing
project **Janus**
Entity Recognition
Data Preparation
project **Metanome**
Change Exploration

Chart **4**

# Agenda

❏ Chair Introduction

❏ **Organizational Information**

❏ Data quality and AI

❏ Your Tasks

# What about you?

# Seminar Topic

❏ **Research Questions**

  ❏ How does (training/testing) data quality influence the performance of AI models?

  ❏ Is there a need for novel data quality dimensions? What could such dimensions be?

❏ **Deliverable**

  ❏ Collaborative paper-style technical report

  ❏ Code, models, and generated datasets

❏ **Teams**

  ❏ 3 teams of 2 students each (At most 6 participants)

*Seminar Webpage*

# Main Milestones

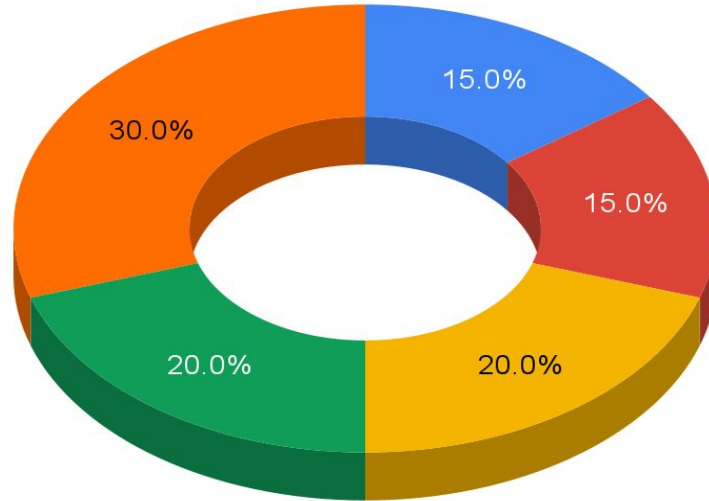| Group allocation | |
|---|---|
| Technical presentation of a paper | Implement your Polluters |
| Think about a beneficial new dimension | |
| Mid-term presentation | |
| Build a machine learning pipeline | Run experiments and write about your results |
| End-term presentation | |
| Final paper writing | |

Seminar Webpage

- ❏ 3 teams
- ❏ 3 ML Tasks
- ❏ 3 quality dimensions
- ❏ 3 new dimension !

# Grading



- Active participation in meetings and discussions
- Technical presentation of a scientific paper
- Mid- and End-term presentation
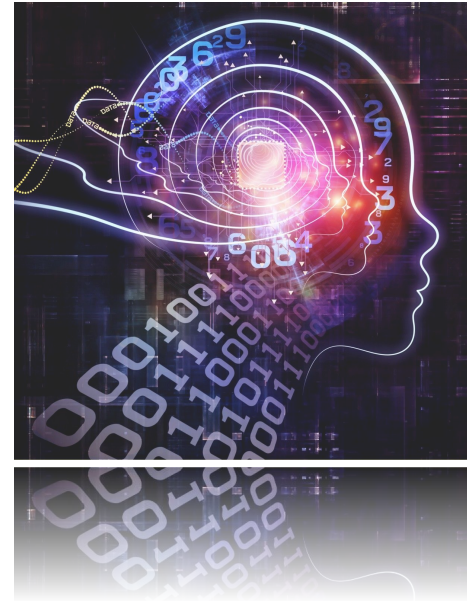- Quality of implementation and results
- Final paper-style submission

15.0%
15.0%
20.0%
20.0%
30.0%

Seminar Webpage

# Further Procedure

- ❏ To apply for this seminar (binding):
  - ❏ **Email** to hazar.harmouch@hpi.de
  - ❏ **Deadline**: Friday 29.10.2021 23:59
  - ❏ **Notification**: Monday 1.11.2021
  - ❏ Register with the Studienreferat
- ❏ In case of too many applications, we need to choose **randomly**.
- ❏ Group allocation deadline: **2.11.2021**
- ❏ Meeting next week: at **Campus II, Building F, Room 2.10 (instead of E.06)**.
- ❏ The course is **on-site**. However, we will switch to hybrid/online mode if the regulation changes.

*Seminar Webpage*

# Agenda

❏ Chair Introduction

❏ Organizational Information

❏ Data quality and AI

❏ Your Tasks

# AI is a Rock Star!

- **Prediction**
  - □ Weather, natural disaster, predictive maintenance, disease
- **Optimization**
  - □ Planning, traffic, logistics, machine efficiency, site selection
- **Individualization**
  - □ Digital health and personalized medicine, personalized learning, r
- **Comfort**
  - □ Sharing, smart home, authentication (face, gait)
  - □ Autonomous vehicles
- **Intelligence**
  - □ Fraud detection, translation, gaming
  - □ Robotics



https://unsplash.com/photos/JfolIjRnveY

# But…



If you have never failed you have never lived.

~ Abraham Lincoln

# AI Failure Example- Amazon's Recruiting Tool

❏ The tool automates the process of reviewing job applicants' resumes.

❏ It showed **bias** against women.

❏ There are many more types of bias.

# AI Failure Example- Microsoft Tay Chatbot



- ❏ Tay was built to learn from interactions to have better conversations in the future.

- ❏ Tay posted **racist** and **derogatory** offensive tweets.

# AI Failure Example- Uber Self Driving Car



❏ The incident on March 18th (2018) took place of the inability to classify an object as a pedestrian unless that object was near a crosswalk.

❏ It was trained on **unrepresentative** training data.

# AI Failure Example- Erroneous Labels

Helps me realize I am ok Not a big slob now I feel better!!!!!!! Yay Yay Ya! No more blues!

Amazon given label:
**Negative**

We guessed: **Positive**

ImageNet given label:
**dough**

We guessed: **pizza**

ImageNet given label:
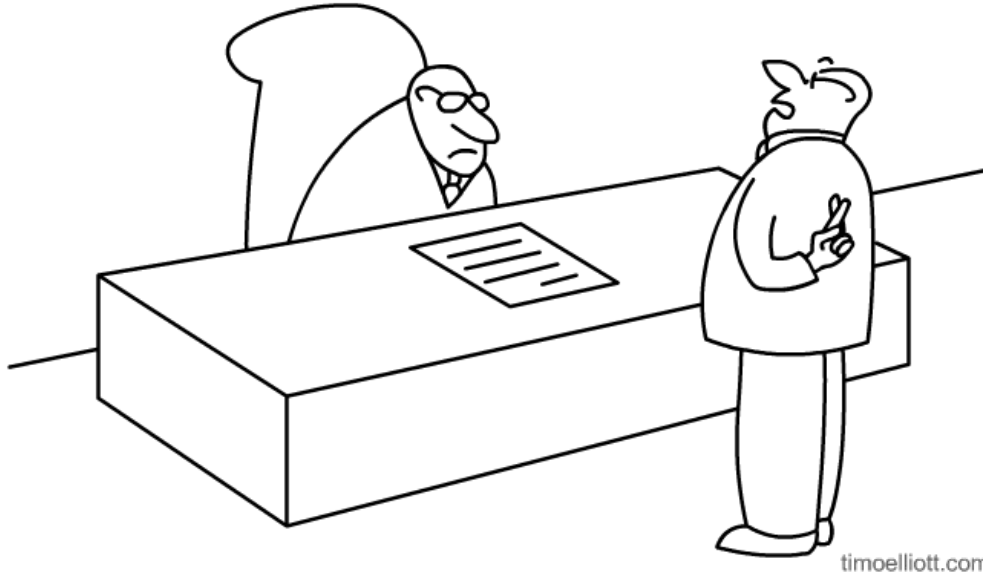**feather boa**

We guessed: **Chihuahua**

Caltech-256 given label:
**yo-yo**

We guessed: **golf-ball**

# Lesson Learned!



"Yes sir, you can absolutely trust those numbers"

❏ AI performance is heavily influenced by the underlying data.

❏ It is important to understand this correlation!

# Real-world data is raw and dirty

"Garbage in, garbage out"

Your analysis is as good as your data.

# Real-world data is raw and dirty

| | | | | | |
|---|---|---|---|---|---|
| 488941 britney spears | 29 britent spears | 9 brinttany spears | 5 brney spears | 3 britiy spears | 2 brirreny spears |
| 40134 brittany spears | 29 brittnany spears | 9 britanay spears | 5 broitney spears | 3 britmeny spears | 2 brirtany spears |
| 36315 brittney spears | 29 britttany spears | 9 britinany spears | 5 brotny spears | 3 britneeey spears | 2 brirttany spears |
| 24342 britany spears | 29 btiney spears | 9 britn spears | 5 bruteny spears | 3 britnehy spears | 2 brirttney spears |
| 7331 britny spears | 26 birttney spears | 9 britnew spears | 5 btiyney spears | 3 britnely spears | 2 britain spears |
| 6633 briteny spears | 26 breitney spears | 9 britneyn spears | 5 b | 3 | 2 |
| 2696 britteny spears | 26 brinity spears | 9 britrney spears | 5 g | | |
| 1807 briney spears | 26 britenay spears | 9 brtiny spears | 5 s | | |
| 1635 brittny spears | 26 britneyt spears | 9 brtittney spears | 4 b | | |
| 1479 brintey spears | 26 brittan spears | 9 brtny spears | 4 b | | |
| 1479 britanny spears | 26 brittne spears | 9 brytny spears | 4 b | | |
| 1338 britiny spears | 26 btittany spears | 9 rbitney spears | 4 b | | |
| 1211 britnet spears | 24 beitney spears | 8 birtiny spears | 4 b | | |
| 1096 britiney spears | 24 birteny spears | 8 bithney spears | 4 b | | |
| 991 britaney spears | 24 brightney spears | 8 brattany spears | 4 b | | |
| 991 britnay spears | 24 brintiny spears | 8 breitny spears | 4 b | | |
| 811 brithney spears | 24 britanty spears | 8 breteny spears | 4 b | | |
| 811 brtiney spears | 24 britenny spears | 8 brighty spears | 4 b | | |
| 664 birtney spears | 24 britini spears | 8 brintay spears | 4 b | | |
| 664 brintney spears | 24 britnwy spears | 8 brinttey spears | 4 b | | |
| 664 briteney spears | 24 brittni spears | 8 briotney spears | 4 b | | |
| 601 bitney spears | 24 brittnie spears | 8 britanys spears | 4 b | | |
| 601 brinty spears | 21 biritney spears | 8 britley spears | 4 b | | |
| 544 brittaney spears | 21 birtany spears | 8 britneyb spears | 4 b | | |
| 544 brittnay spears | 21 biteny spears | 8 britnrey spears | 4 b | | |
| 364 britey spears | 21 bratney spears | 8 britnty spears | 4 b | | |
| 364 brittiny spears | 21 britani spears | 8 brittner spears | 4 b | | |
| 329 brtney spears | 21 britanie spears | 8 brottany spears | 4 b | | |
| 269 bretney spears | 21 briteany spears | 7 baritney spears | 4 b | | |
| 269 britneys spears | 21 brittay spears | 7 birntey spears | 4 b | | |
| 244 britne spears | 21 brittinay spears | 7 biteney spears | 4 b | | |
| 244 brytney spears | 21 brtany spears | 7 bitiny spears | 4 b | | |
| 220 breatney spears | 21 brtiany spears | 7 breateny spears | 4 britneuy spears | 2 barittany spears | 2 britneyh spears |
| 220 britiany spears | 19 birney spears | 7 brianty spears | 4 britnewy spears | 2 bbbritney spears | 2 britneym spears |

**LIVE BBC NEWS CHANNEL**

Page last updated at 11:45 GMT, Thursday, 19 February 2009

E-mail this to a friend    Printable version

## The mystery of Ireland's worst driver

Details of how police in the Irish Republic finally caught up with the country's most reckless driver have emerged, the Irish Times reports.

He had been wanted from counties Cork to Cavan after racking up scores of speeding tickets and parking fines.

However, each time the serial offender was stopped he managed to evade justice by giving a different address.

But then his cover was blown.

It was discovered that the man every member of the Irish police's

Poles are Ireland's largest immigrant population

PERMIS DE CONDUIRE
PRAWO JAZDY
1.
2. JULIUSZ
3. KRAKÓW
4a. 26.02.03   4c. PREZY
MIASTA V
4b.           4d.

SEE A
Cou
03 Fe

RELAT
Irish

The BB
internet

TOP N
Oma
Sinn
City

# FIFA registration form (2010)

# Hidden Values / Hidden Value

| Datenelement | Feld Name1 | Name2 | Name3 | City | District | Street | Sum |
|---|---|---|---|---|---|---|---|
| Handy-Nummer | 41 | 501 | 10 | 0 | **2677** | 297 | 3526 |
| Festnetznummer | 15 | 98 | 6 | 0 | **221** | 9579 | 9919 |
| Kostenstelle | 283 | 1112 | 73 | 2 | **87** | 16 | 1573 |
| Registriernummer | 11 | 583 | 1 | 1 | **0** | 3 | 599 |
| Lieferungsnummer | 55 | 390 | 9 | 0 | **212** | 15 | 681 |
| Abteilung | 3711 | **9997** | 115 | 60 | **439** | 175 | 14497 |
| Sperrkennzeichen | 129 | 143 | 2 | 0 | **66** | 9 | 349 |
| Löschkennzeichen | 1028 | 442 | 5 | 36 | **113** | 10 | 1634 |
| Rechtsform | **131700** | 66136 | 187 | 6 | **64** | 57 | 198150 |
| Kreditoreninfo | 0 | 100 | 11 | 0 | **18** | 0 | 129 |
| Kommissionsinfo | 216 | 352 | 1 | 2 | **36** | 10 | 617 |
| Baustelle | 2013 | 3452 | 42 | 5 | **124** | 222 | 5858 |
| Abladestelle | 2923 | 3808 | 94 | 1503 | **958** | 3065 | 12351 |
| Behörde | 13410 | 12461 | 172 | 19 | **295** | 7075 | 33432 |
| **Summe** | 155535 | 99575 | 728 | 1634 | **5310** | 20533 | |

Source: Joachim Schmid, FUZZY! Informatik AG

# From Data Errors (aka. Data Quality) to Data Problems (aka. Information Quality)

- Incorrect data:                   Accuracy
- Missing data:                   Completeness
- Poor formatting:                   Representational consistency

- Old data:                   Timeliness
- Unknown data source:                   Trustworthiness

- Hard to reach data:                   Accessibility
- Slow connection:                   Latency

- And many more information quality dimensions
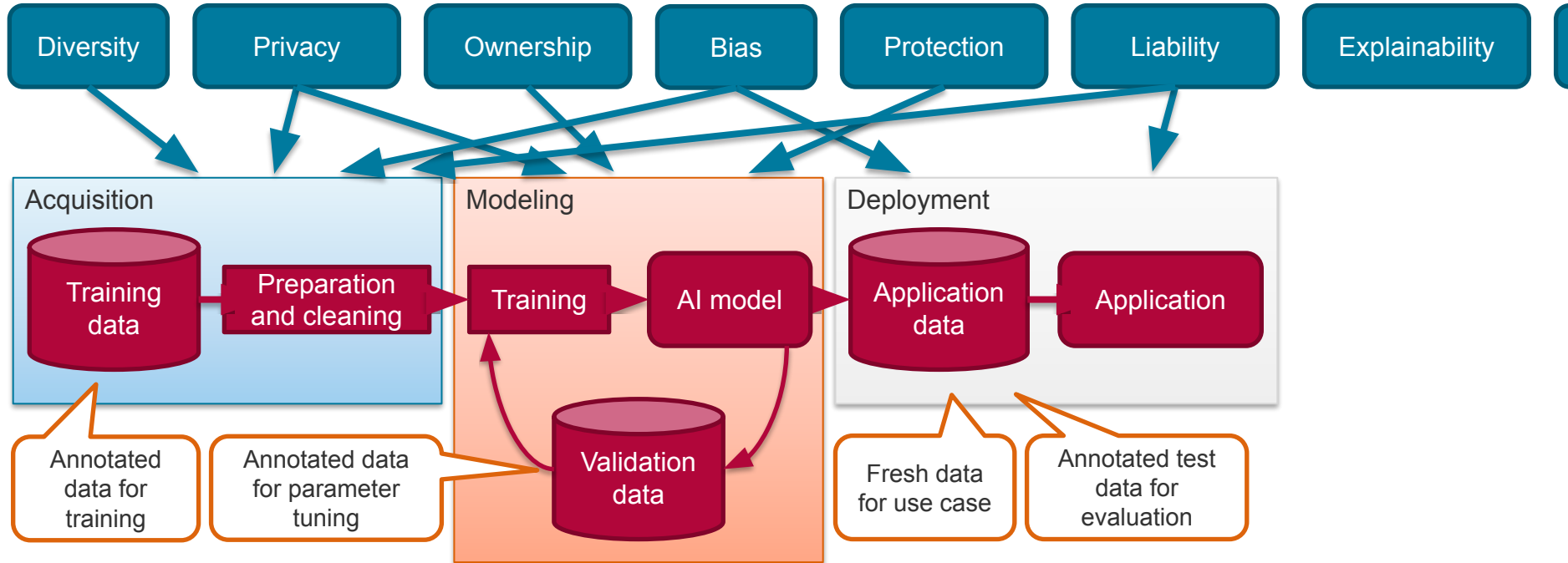
# IQ Classification of Wang and Strong

- Intrinsic IQ
  - □ Believability, Accuracy, Objectivity, Reputation
- Contextual IQ
  - □ Value-added, Relevancy, Timeliness, Completeness, Amount
- Representational IQ
  - □ Interpretability, Understandability, Repr. Consistency, Repr. conciseness
- Accessibility IQ
  - □ Accessibility, Security

- And more
  - □ Customer support, documentation, reliability, latency, price, response time, verifiability
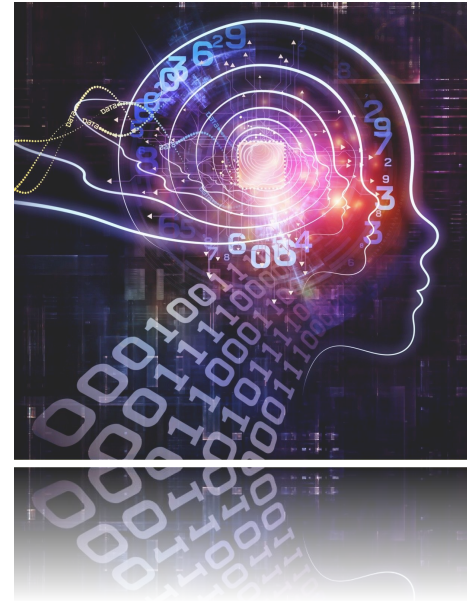
Wang & Strong
Beyond Accuracy:
What data quality
means to data
consumers
*Management of
Information Systems,*
1996*, 12(4)*, 5-34

# New AI-specific Data Quality Dimensions and Where to Find Them

# Agenda



❏    Chair Introduction

❏    Organizational Information

❏    Data quality and AI

❏    Your Tasks

# Team Tasks

- ❏ Form a team and choose a ML task: Classification, Regression, or Clustering
  - ❏ A specific data quality dimension will be assigned to each team.

- ❏ Each team will have the following tasks:
  - ❏ Find 3 datasets used for your task.
  - ❏ Implement a polluter in terms of your assigned data quality dimension.
  - ❏ Think of a novel data quality dimension and implement a respective polluter.
  - ❏ Implement a ML pipeline with 5 different algorithms that belong to your ML task.
  - ❏ Write about the results in the technical report.
  - ❏ In between: present a related work paper