# Data Mining SS25

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

Technology
Arts Sciences
**TH Köln**

# Agenda

- Company August Rüggeberg

- Topic of the semester project

- Data sources

- CRISP-DM

- General information

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

**Technology**
**Arts** **Sciences**

**TH Köln**

# Company

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics
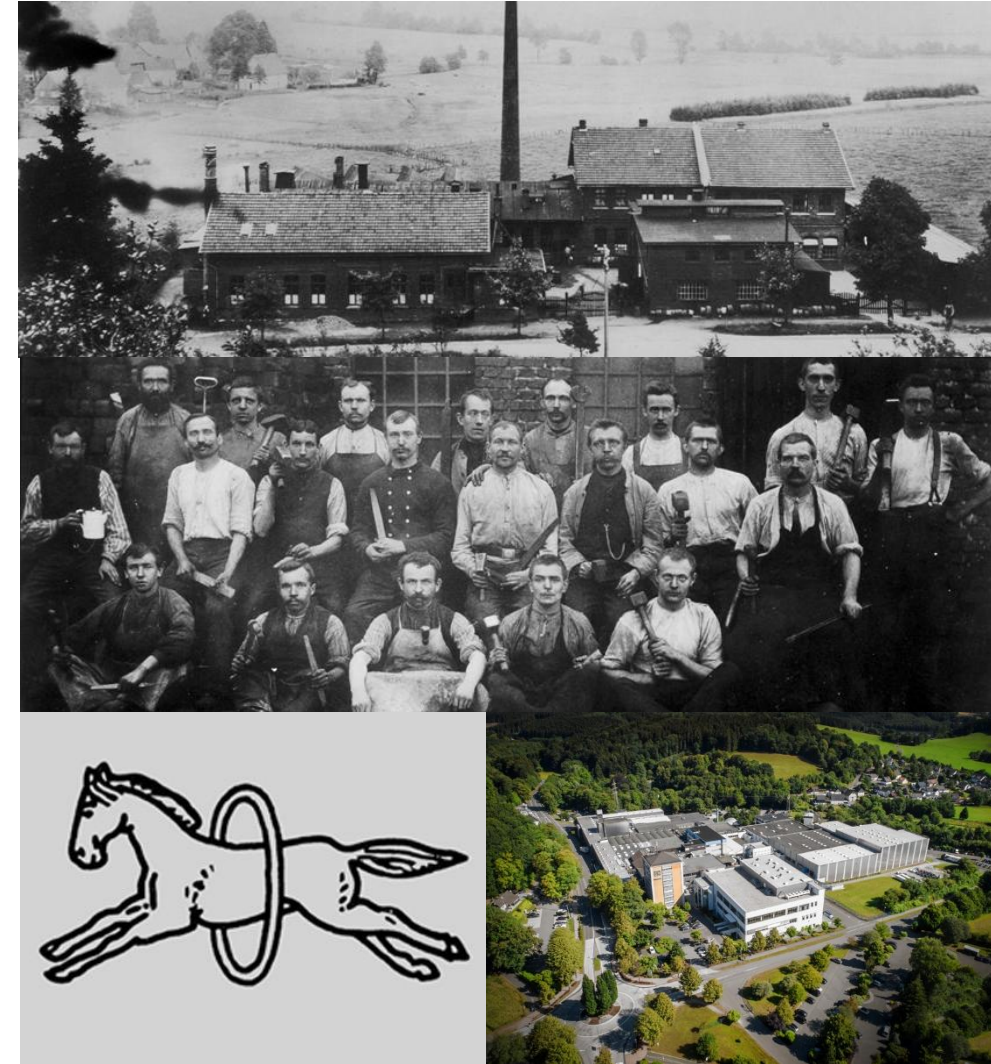
**Technology
Arts Sciences
TH Köln**

# August Rüggeberg

Company description

August Rüggeberg (PFERD) is leading in the

- development

- production

- support

- and distribution

of tool solutions for work on surfaces and material cutting.

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

**Technology**
**Arts Sciences**
**TH Köln**

# August Rüggeberg

Customers

- Metalworking is the primary focus among our top customers

- Key focus areas include aerospace, shipbuilding, and automotive, though the tools are also used in many other industrial sectors.



Foundry



Aerospace



Shipbuilding

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

Technology
Arts Sciences
TH Köln

# Topic of the semester project

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

**Technology**
**Arts Sciences**

**TH Köln**
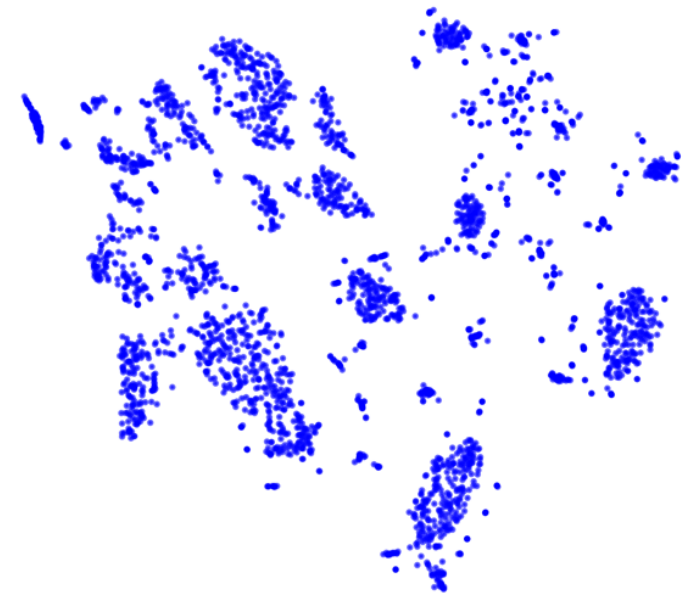
# Topic of the semester project

Topic area world economy

**Approach: Comprehensive Market Analysis** using clustering methods and tailored feature selection

**Indicators:**

- **Economic Indicators** (GDP, economic growth)

- **Demographic Indicators** (population size, population growth)

- **Trade Data** (import/export volumes of related products)

- **Business-Relevant Metrics** (Steel consumption (use-case reference))

- **Soft Factors** (Corruption Index, Economic Freedom Index)

**Goal:** Group countries into comparable market clusters to uncover high-potential markets.

Technology
Arts Sciences

TH Köln

# Data sources

15.04.2025

Lukas Bader

Data Mining SS25

Seite 9

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

**Technology
Arts Sciences
TH Köln**

# Data sources

The Global Ecomomy ([Link](#))

Platform providing reliable economic data for over 200 countries (since 1960)

- Includes 500+ indicators from:
    - Central banks,
    - National statistics offices,
    - International organizations

- Goal: Save time and effort through well-documented, downloadable data

- Free access for: Educators & students from low-income countries

- Founded in 2012 at Georgia State University (USA)

- Selection of approx. 800 annual key features are provided in a python df
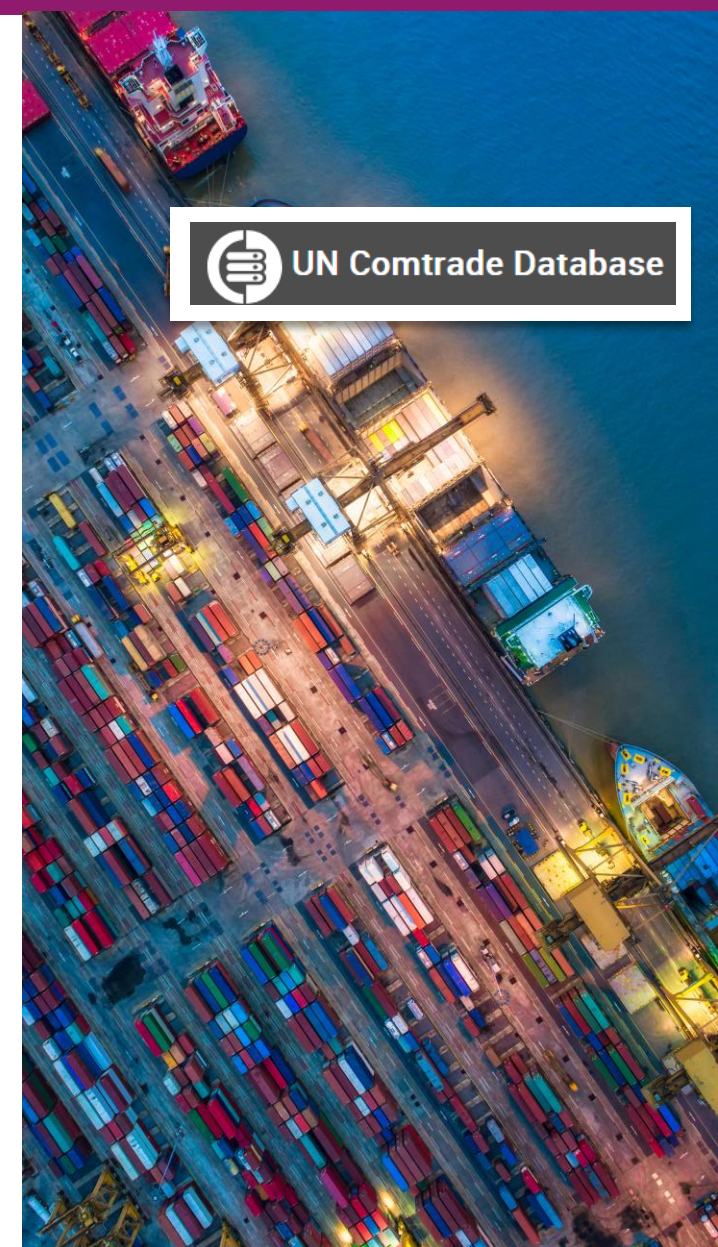
Technology
Arts Sciences

TH Köln

# Data sources

UN Comtrade ([Link](Link))

- Global trade database maintained by the United Nations

- Contains detailed import and export statistics from over 170 reporting countries

- Covers more than 5,000 product categories, classified by:

  - HS (Harmonized System), SITC, BEC, and other classification systems

- Data reported by national customs authorities and standardized for international comparison

- Available at annual and monthly frequency, depending on country

- Includes trade values and quantities, partner countries, and trade flows

- Data must be downloaded and pre-processed independently



UN Comtrade Database

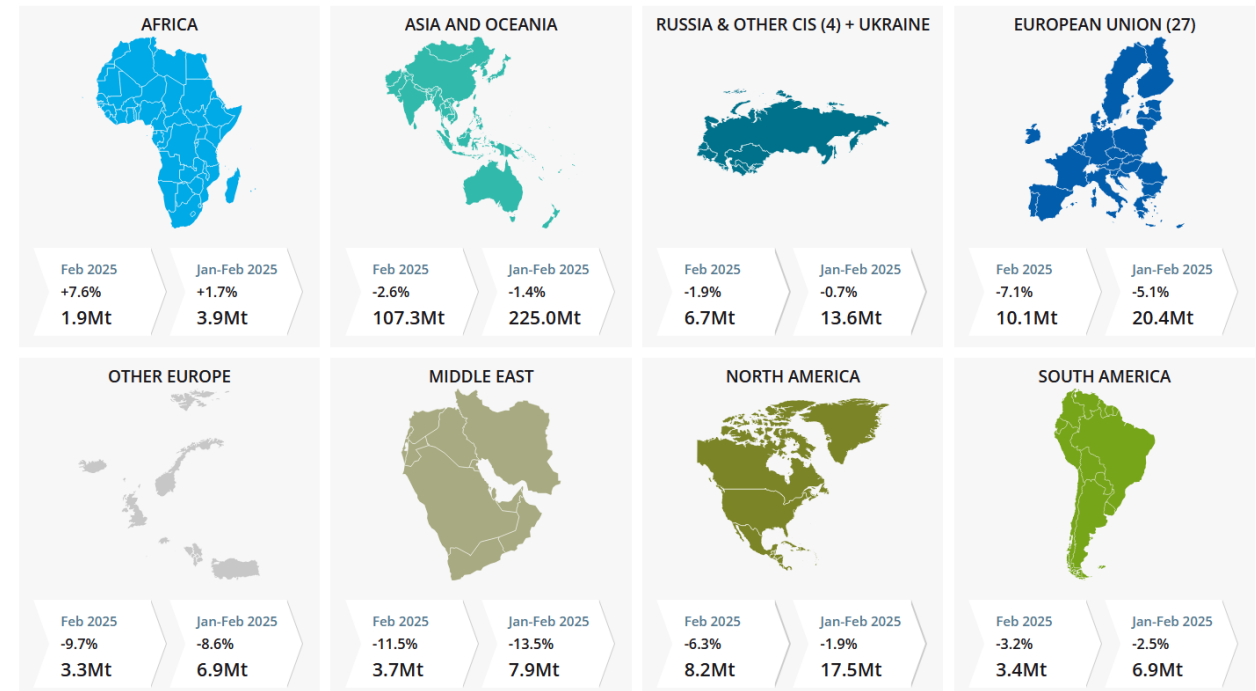**Technology**
**Arts Sciences**
**TH Köln**

# Data sources

World Steel Association ([Link](#))

- One of the largest industry associations in the world

- Represents steel producers, national/regional associations, and research institutes

- Covers around 85% of global steel production

- Provides comprehensive data and statistics on:
    - Crude steel production (monthly & annual)
    - Steel use by country and sector
    - Steel trade flows and demand forecasts



| AFRICA | | ASIA AND OCEANIA | | RUSSIA & OTHER CIS (4) + UKRAINE | | EUROPEAN UNION (27) | |
|---|---|---|---|---|---|---|---|
| Feb 2025 | Jan-Feb 2025 | Feb 2025 | Jan-Feb 2025 | Feb 2025 | Jan-Feb 2025 | Feb 2025 | Jan-Feb 2025 |
| +7.6% | +1.7% | -2.6% | -1.4% | -1.9% | -0.7% | -7.1% | -5.1% |
| 1.9Mt | 3.9Mt | 107.3Mt | 225.0Mt | 6.7Mt | 13.6Mt | 10.1Mt | 20.4Mt |

| OTHER EUROPE | | MIDDLE EAST | | NORTH AMERICA | | SOUTH AMERICA | |
|---|---|---|---|---|---|---|---|
| Feb 2025 | Jan-Feb 2025 | Feb 2025 | Jan-Feb 2025 | Feb 2025 | Jan-Feb 2025 | Feb 2025 | Jan-Feb 2025 |
| -9.7% | -8.6% | -11.5% | -13.5% | -6.3% | -1.9% | -3.2% | -2.5% |
| 3.3Mt | 6.9Mt | 3.7Mt | 7.9Mt | 8.2Mt | 17.5Mt | 3.4Mt | 6.9Mt |

Technology
Arts Sciences
TH Köln

# CRISP-DM cycle with regard to the project

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

Technology
Arts Sciences

TH Köln

# CRISP-DM cycle

Model planning 50%-70%

## Business Understanding

- What is the business use case?
- What is the aim of the project?
- What are the stakeholders?
- When is a market interesting for PFERD?
- Which features could be relevant?

## Data Understanding

- Have I chosen the right features?
- Are there redundant features?
- Is the data available for many markets?
- How should the various data sets be prepared?
- How do I combine the data sets into a df?
- How must the data be provided for clustering?

## Data Preparation

- Is the data suitable for modeling?
- Feature engineering?
- How is the data distributed?
- How to deal with outliers?
- How to deal with missing values?
- Does the data need to be scaled?
- Which method do I use for scaling?
- Must all columns of the df be scaled?

**Technology
Arts Sciences
TH Köln**

# CRISP-DM cycle

Model development 10%-20% & Model deployment / presentation 20%-40%

## Modeling

- Which clustering methods are suitable?
- How to find the best model parameters?
- How to save model parameters?
- Do the data need to be adjusted again?

## Evaluation

- Does clustering make sense in a visual inspection?
- Which metrics can I use for evaluation?
- What are the best clusters, and is that plausible?
- Are there differences between the methods?
- Which features have the greatest impact?
- Are the results relevant to the use case?

## Model (deployment) / presentation

- What are the most important results?
- How can modeling results be a strategic recommendation regarding the use case?
- How do I present my results?
- How do I create trust among stakeholders?

15.04.2025

Lukas Bader

Data Mining SS25

Seite 15

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

**Technology**
**Arts** Sciences
**TH Köln**

# General information

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

Technology
Arts Sciences

TH Köln

# General information

**Data:**

- Datasets will partly be published in a GitHub repo with further information

- Datasets are stored in the pickle file format, which can be easily loaded into a python dataframe

**Procedure:**

- After the group assignment, you will be assigned to a GitHub repo

- The project documentation takes place in the repo

- If you have any questions, you can contact me (lukas.bader@th-koeln.de)

**Interesting links:**

- Country similarity index (use of a distance matrix), Link

Lukas Bader

Data Mining SS25

Fakultät für Informatik und Ingenieurwissenschaften
Institut für Data Science, Engineering, and Analytics

Technology
Arts Sciences
TH Köln