DATA-DRIVEN SEARCH FOR TRAFFIC DRIVERS
WITH DATA PROVIDED BY

EFAHRER.com

# We are an interdisciplinary team with diverse background in business, science, music and IT

**Renzo Torrecuso**

Neuroscience-Engineer-Violinist

**Thomas Brandstätter**

Analytics Engineer, Data Products

**Ekaterina Burakova**

Physical chemist, Dr. rer. nat.

**Clara Thümecke**

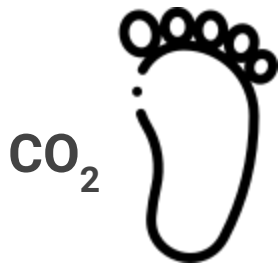Business Development

**Preethi Karumathil**

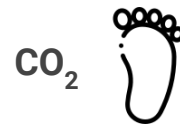IT Engineer

# We'll show how to create impactful content

data driven and backed by machine learning

# We help EFAHRER in empowering their users to contribute to carbon reduction

**EFAHRER.com is a media portal** which strives to influence users to take actions that support $CO_2$ reduction.

$CO_2$

$CO_2$

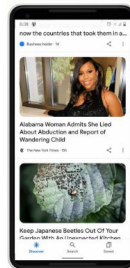We want to **provide valuable insights** for the editorial team

We want to perform a **prognosis of article success**

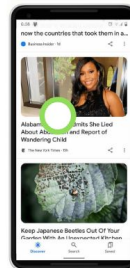# For that we analyzed one of EFAHRER's biggest traffic sources for news articles and enriched the data

Editorial data
**6.899** unique articles

Feed
Impression

Feed
Click

**Billions**

**Millions**

# Data quality and completeness of raw data led to extensive preprocessing and analysis for modelling

## Challenges

- Missing values
- Lack of article versions
- 3 different aggregation levels for selected metrics
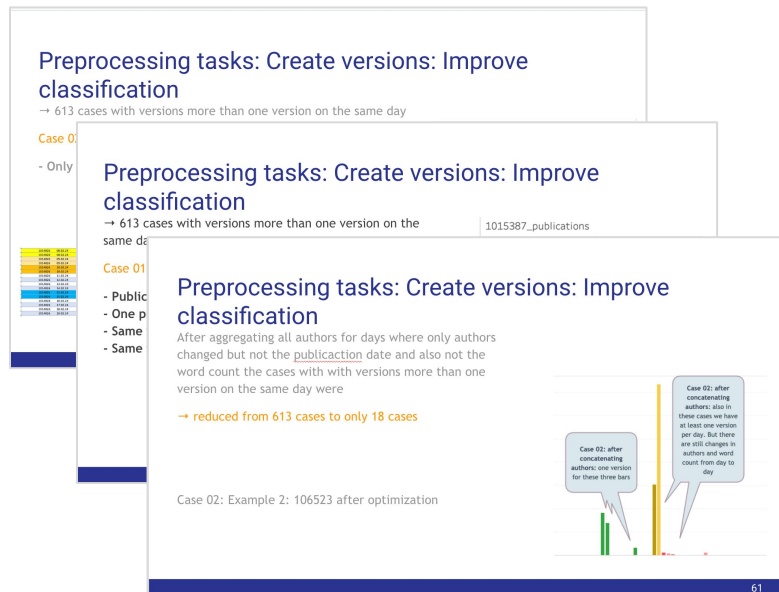
## Strategies:


Delete


Impute


Scrape

# Scraping of 6 000 articles increased the data quality and added new features

Respecting ethics in web scraping we managed to add

**+ 6** visible features

**+ 5** meta features (invisible)

# By importing related search terms for 17 product categories we added a trends score to our data

Using a NLP classifier we matched each article with a related search label and trend score

**+ 2** features

Search trend for "E-Bike"

Related queries

# We identified the following relevant features

⇒ Article genre and topic

⇒ Type of the first media: video or image

⇒ Word count and lengths of the metadata

Live demo

🔄 Features in the positive feedback loop with the target were **ignored** or **normalized**:
- Number of likes, dislikes, video views
- Number of URLs → *URL update frequency*

# We verified our hypotheses and created prediction tool

**Baseline**

**Feature Engineering**

**Final Model**

# With the baseline model we created an advanced starting point for our modeling

**Baseline**

Performance (R² = 0.35)
Input:
- Text fields size
- Product classification
- URL update frequency

Feature importance

# We engineered additional features based on the existing data



**Feature Engineering**

NLP

**Raw data**
- Authors
- Article length

**Meta data**
- Meta title
- Meta description
- Media type

**Calculated data**
- Number of versions
- Urls per day

**Clickbait**
- Clickbait label
- Title has colon

**Sentiment**
- Sentiment title
- Sentiment abstract

**Trends Class.**
- Trends label
- Trends score

# Our final model is a stable starting point for predicting article performance



**Final Model**

- Modeling with AutoML Tables by Google Vertex AI
- Best performance (R² = 0.49) for simpler model without full text features

Live demo

# Updating of articles enhances outcomes, the media plays relevant role

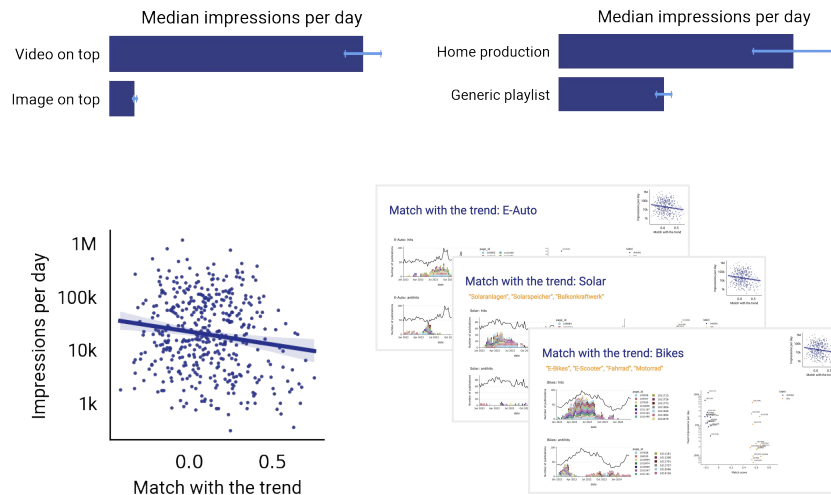**1** Change in URL has **tangible** impact on impressions

**2** The algorithm prefers articles with videos over images as the first media on page

**3** Video production pays off!

**4** Optimize publication timing alongside trends

**5** Algorithm does not punish clickbait behavior



Median impressions per day

Video on top

Image on top

Median impressions per day

Home production

Generic playlist



Impressions per day

1M
100k
10k
1k

0.0    0.5

Match with the trend

Match with the trend: E-Auto

Match with the trend: Solar
"Solaranlage", "Solarspeicher", "Balkonkraftwerk"

Match with the trend: Bikes
"E-Bikes", "E-Scooter", "Fahrrad", "Motorrad"

# Further improvements promise a reliable prediction of page impressions

➡️ Try out different semantic segmentation and model each segment individually (e.g. News)

➡️ Dive deeper into the video and image content and formats

➡️ Refine the trend-related features (e.g. different keyword & time matching, trend sources)

➡️ Improve the evaluation of "clickbaitness"

➡️ Fine tune sentiment analysis

➡️ The full article history would provide new valuable features

# Thank you for your attention

Special thanks to:

 **neue fische** coach team who made this possible

- Nico Steffen
- Aljosha Wilhelm
- Lina Willing
- Jin-Ho Lee
  … and others

 **efahrer.com** team who kindly supported this project

- Markus Höllmüller
- Analytics team