

# **Deutsche Krebsgesellschaft**

Classification Usefulness

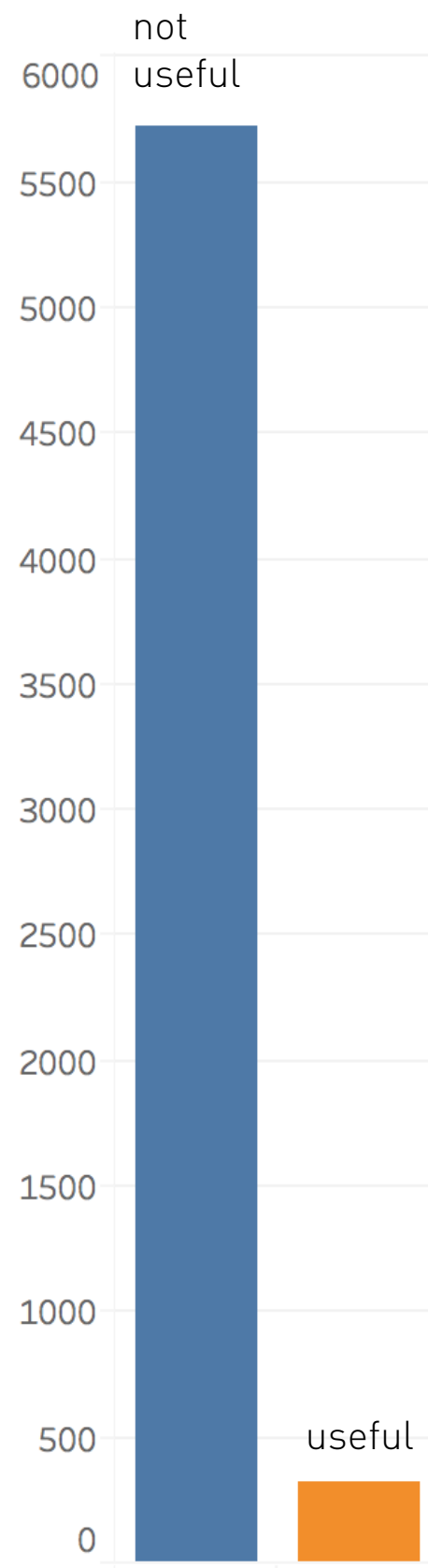
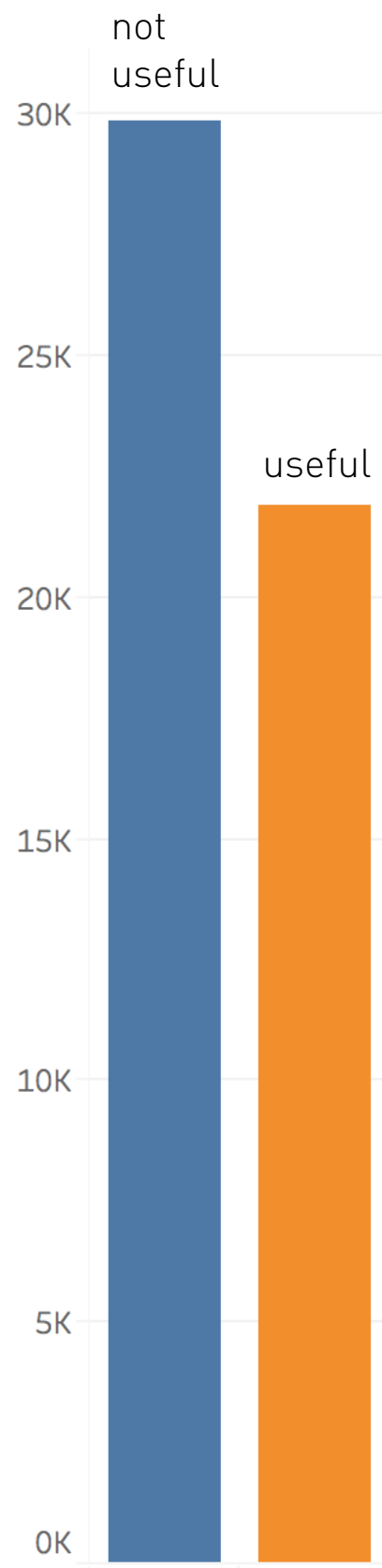
# Use case

Which article / PDF is useful for the website?

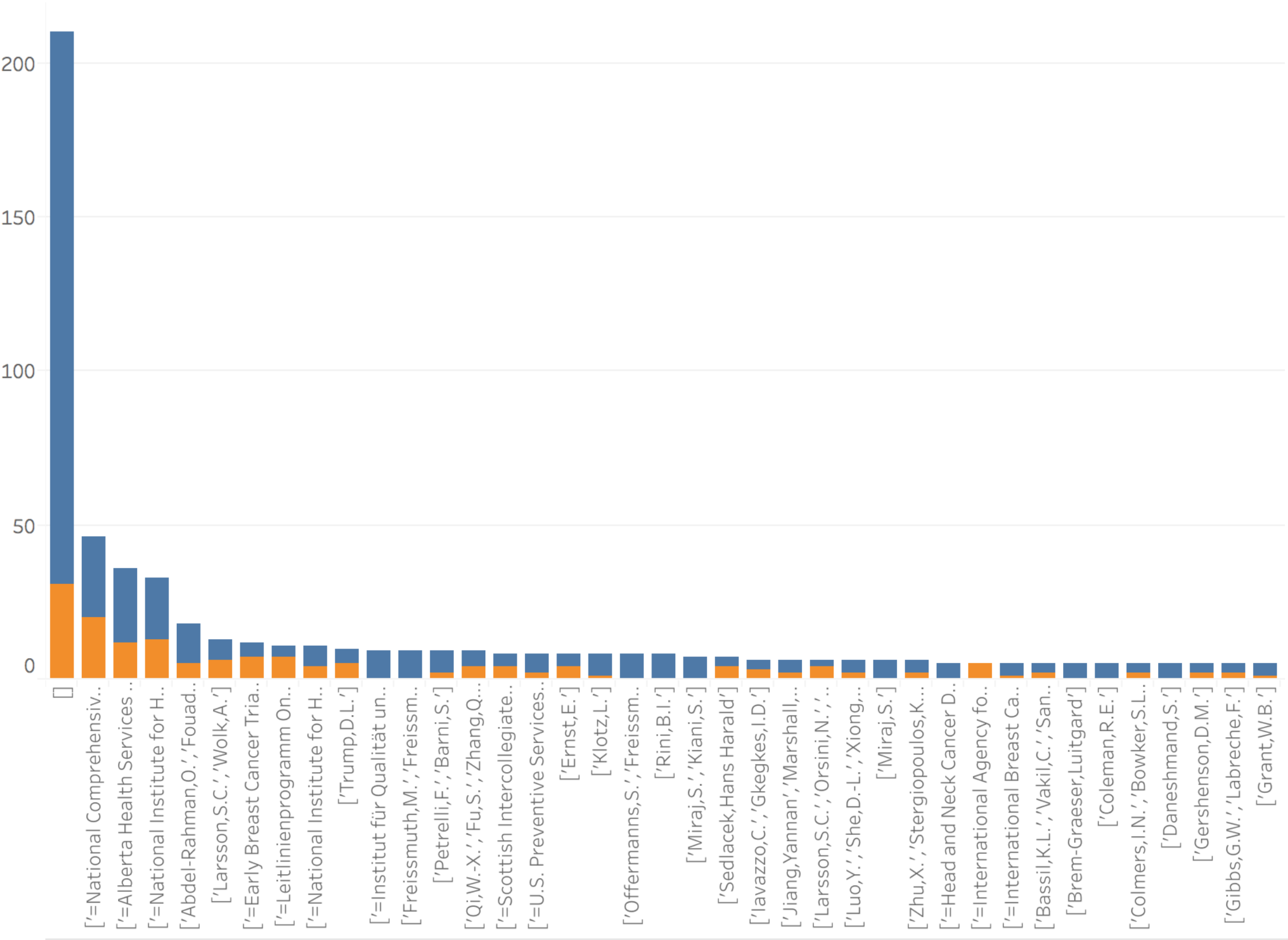
GOAL:

**classification of articles with yes/no**

# CLASSES



AUTHORS



# First Step

We singled out and analysed variables which could possibly influence usefulness: Author, Abstract ...

BUT we shouldn't get the ideal score and instead build something which is scalable and easy to use for DKG in the future

Thus...

# ElasticSearch

..we decided to classify on title, subtitle, author and article.

**1.** indexed the data in Elastic search continuously

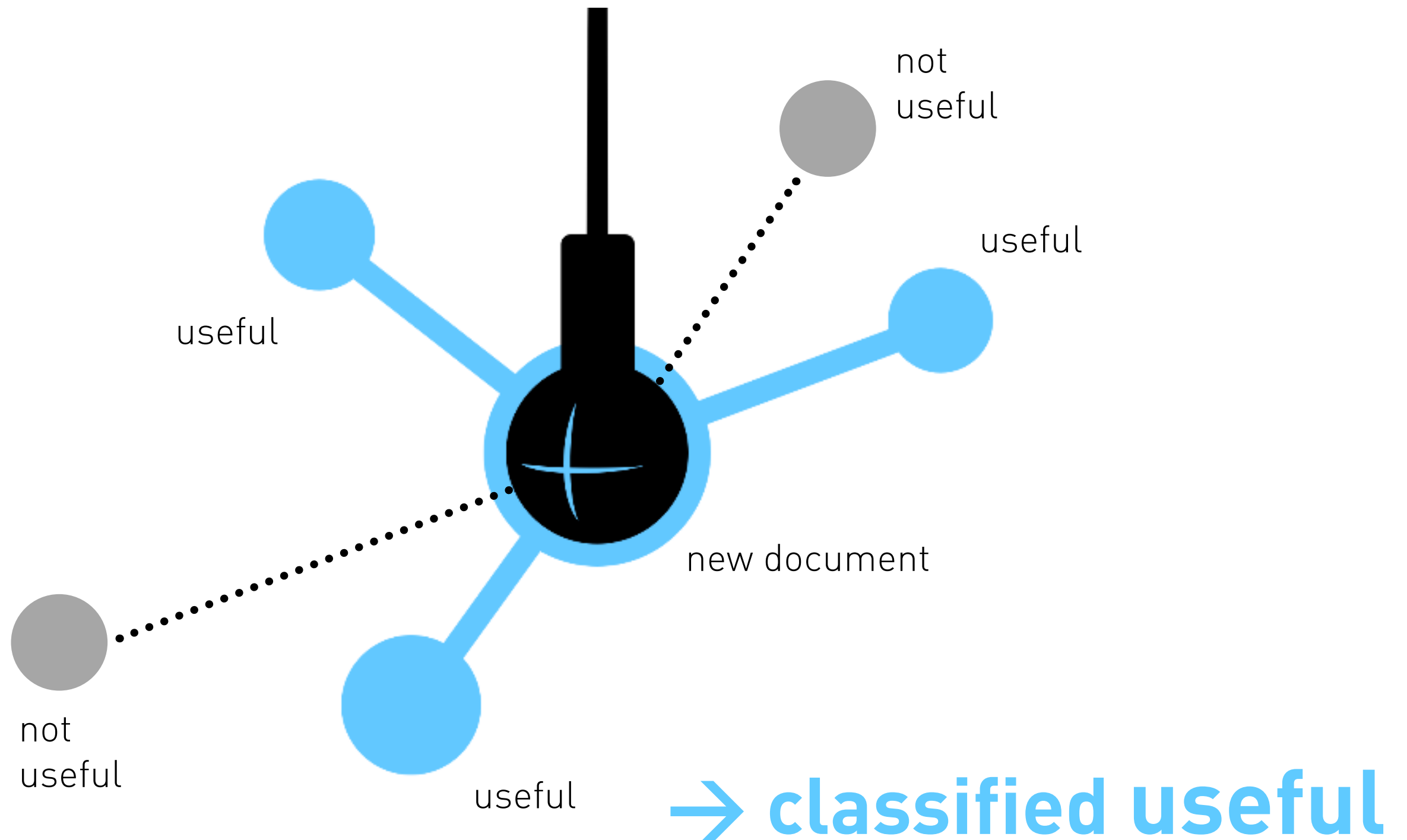
...

# ElasticSearch

**2.** used **k-nearest Neighbor** based on  
»more like this« query

→ under the hood »more like this« applies tf-idf on the articles

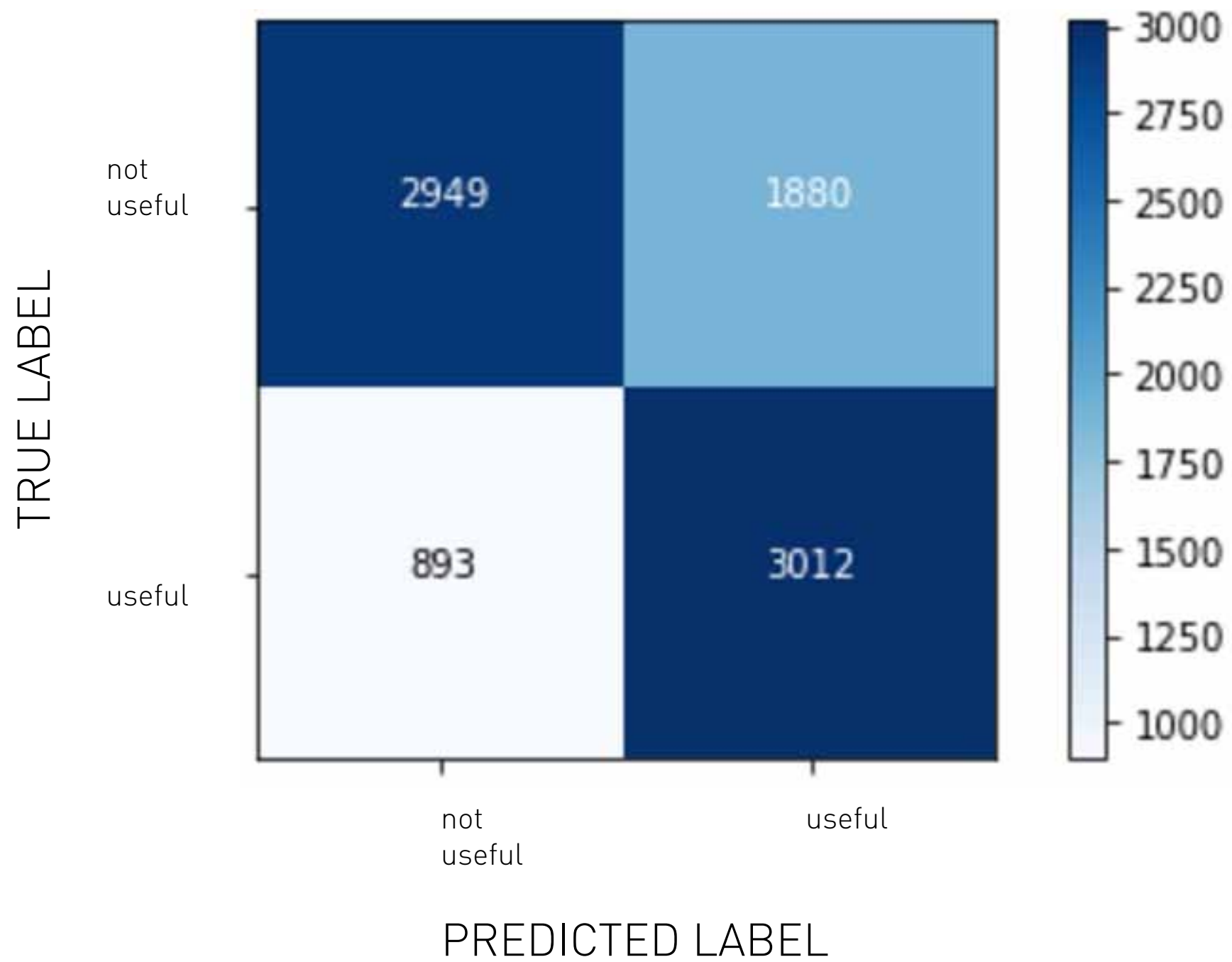
# k-nearest Neighbor





# RESULTS

confusion matrix, without normalisation



# Why Elasticsearch?

- no extra training time needed
- no additional infrastructure
- adapts to new data in real time
- updates runs continuously

# **Danke von unserem Team!**

Till Breuer

breuertill@gmail.com

Anke Nowottne

anke.nowottne@gmail.com

Oliver Rieger

o.rieger@posteo.de

Franziska Wittleder

franziska.wittleder@gmail.com