



"Using a simple tool to solve a complex problem does not result in a simple solution." Larry Wall

Classifying unstructured text

Deterministic and machine learning approaches

Stephanie Fischer
Dr. Christian Winkler

Hamburg München Berlin Köln Leipzig

Apache Big Data Sevilla , 15th November 2016

- 01 About us
- 02 Text statistics
- 03 Categories
- 04 Text classification
- 05 Conclusion and outlook



Stephanie Fischer

Product Owner Text Analytics
Big Data, Agile & Change
mgm consulting partners



Dr. Christian Winkler

Enterprise Architect
Big Data, Data Science
mgm technology partners

Agenda

Speaker

01

About us

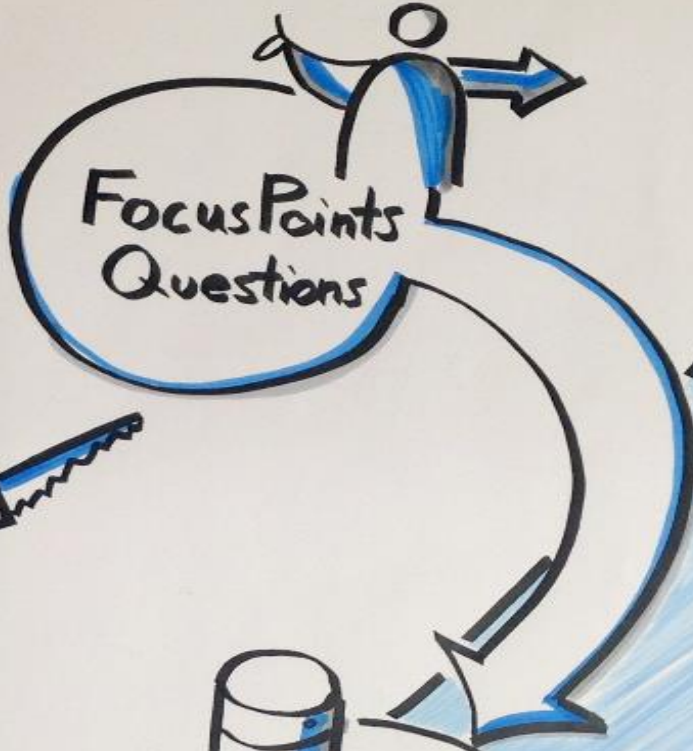
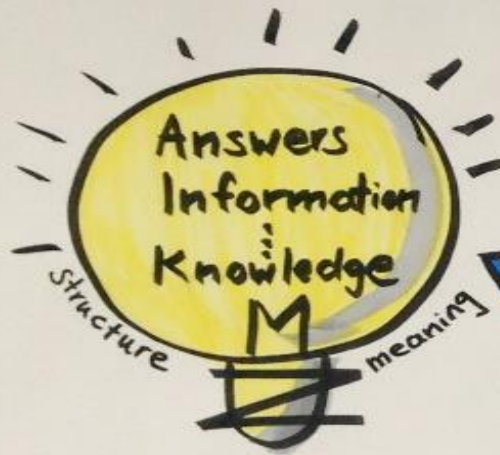
Stephanie and Christian according to their browser history



02

Text statistics

Why text mining



↑
40.000
exabytes



Comparing word frequency of news from Reuters, Telegraph, Aljazeera



Reuters World News

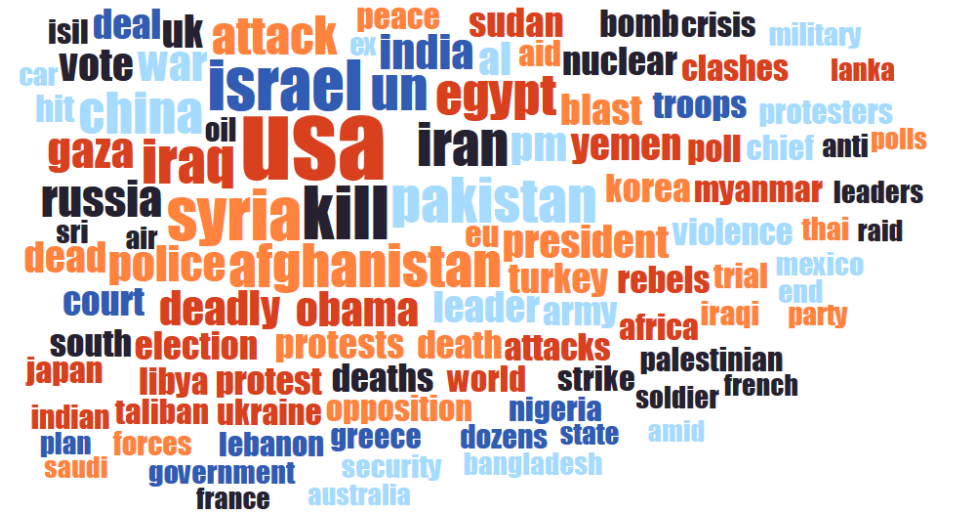
163,919 headlines

🕒 9 years

Telegraph

958,996 headlines

🕒 9.5 years



Aljazeera

94,309 headlines

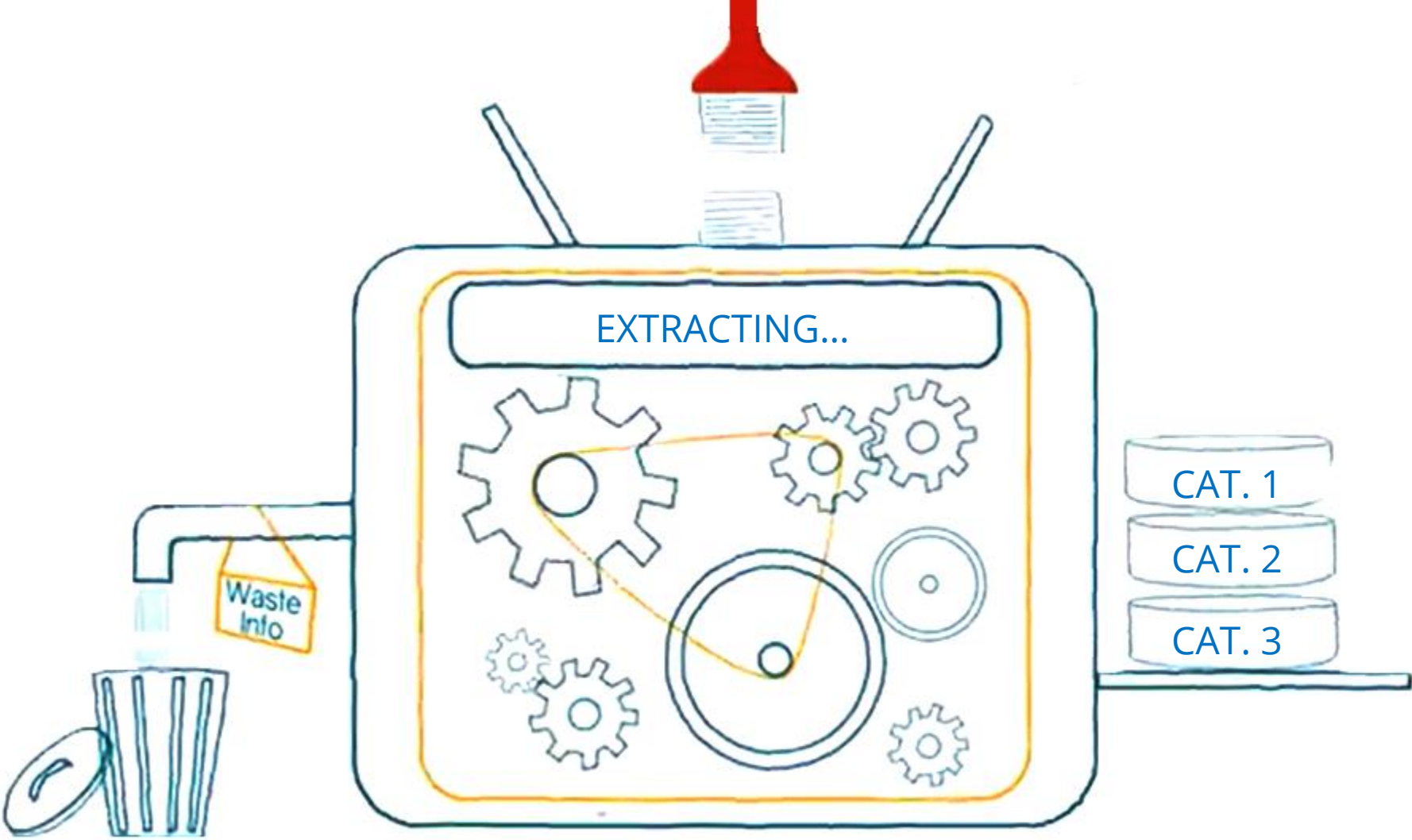
🕒 8.5 years



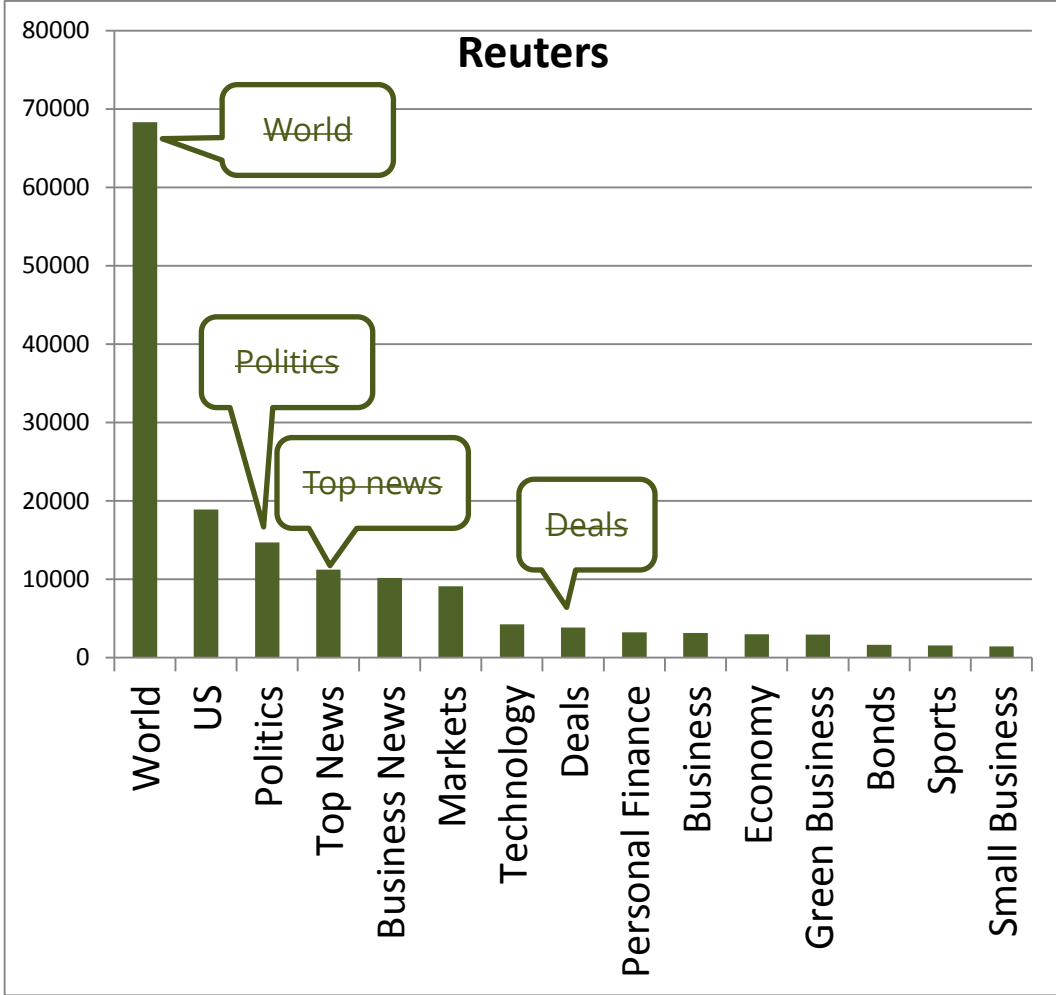
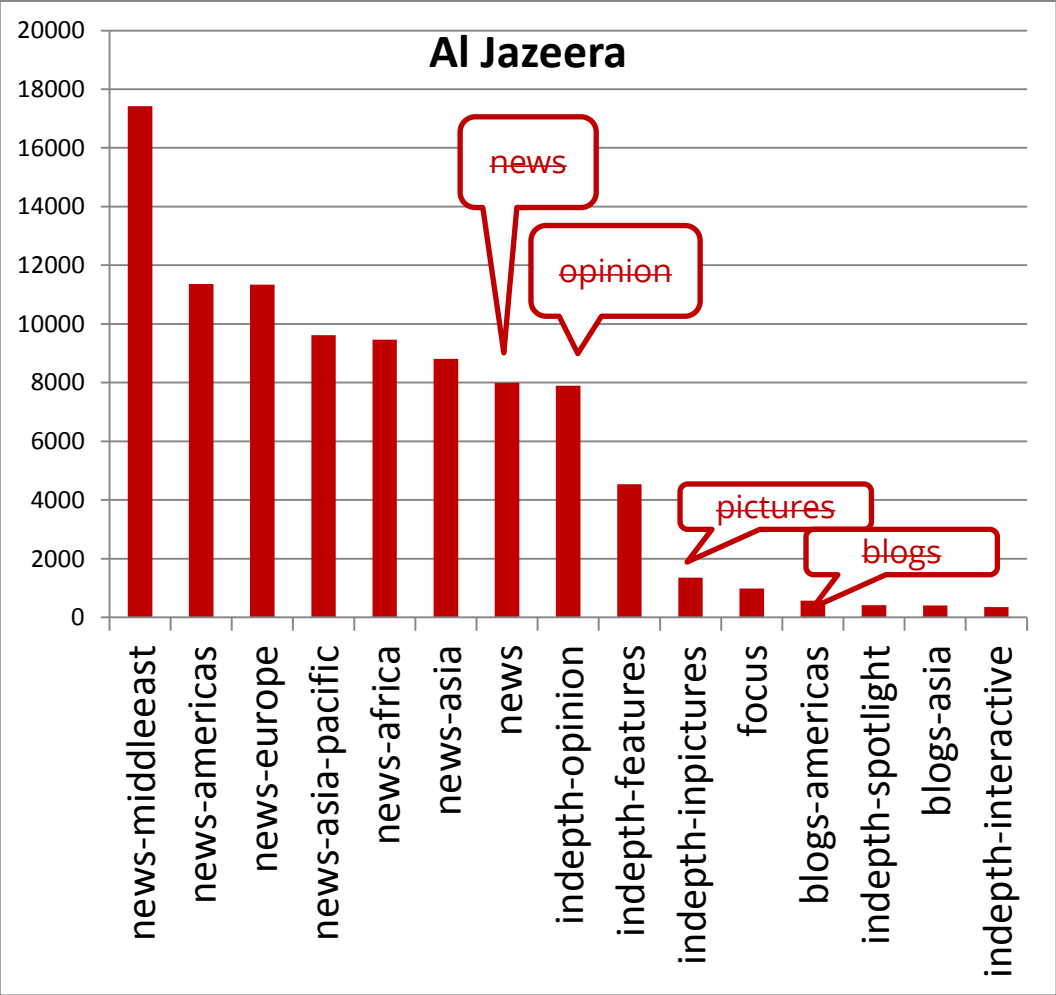
03

Categories

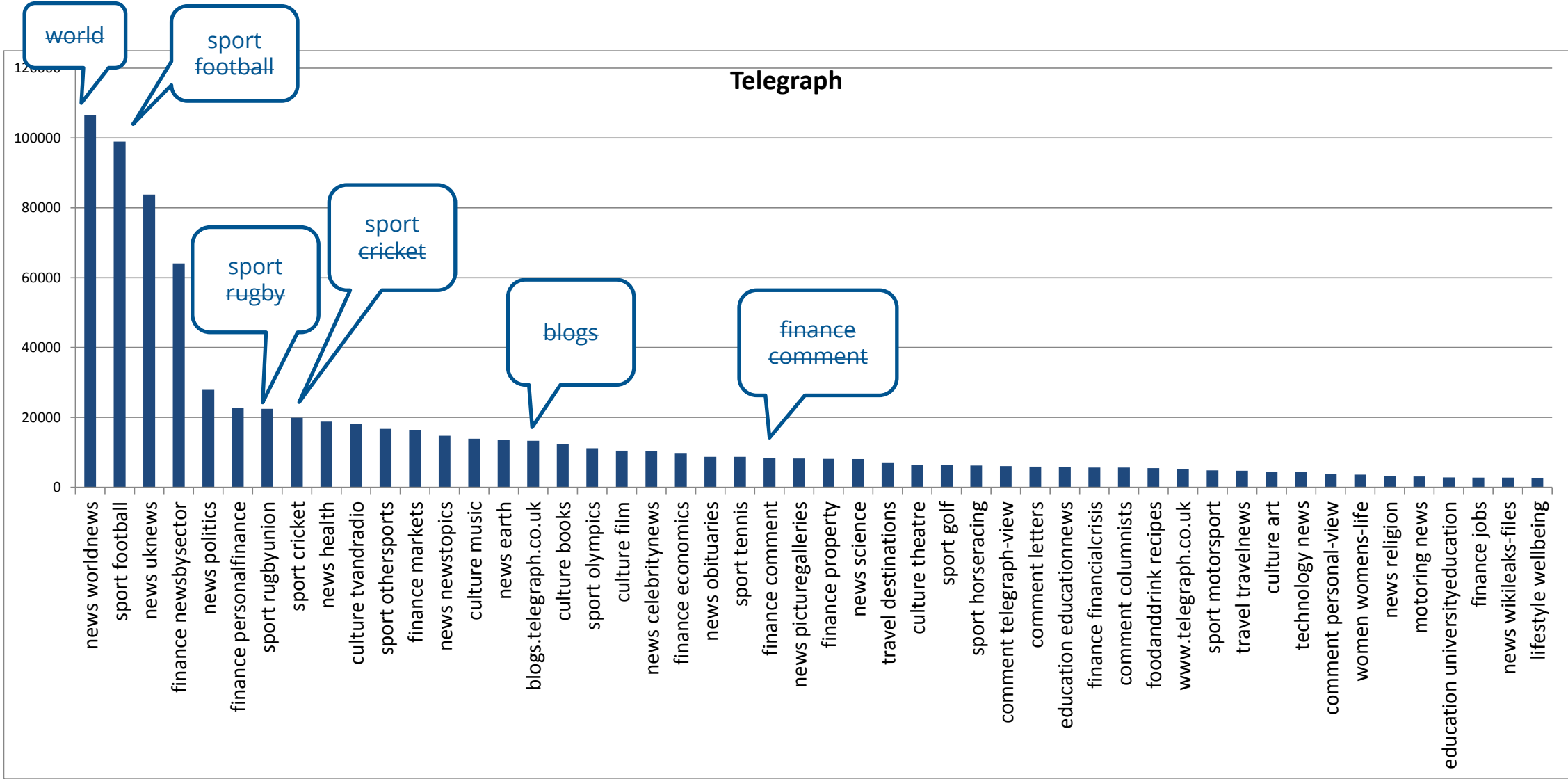
Finding meaningful categories. Each text is different. Challenge accepted!



Comparing pre-defined categories of Al Jazeera, Reuters...



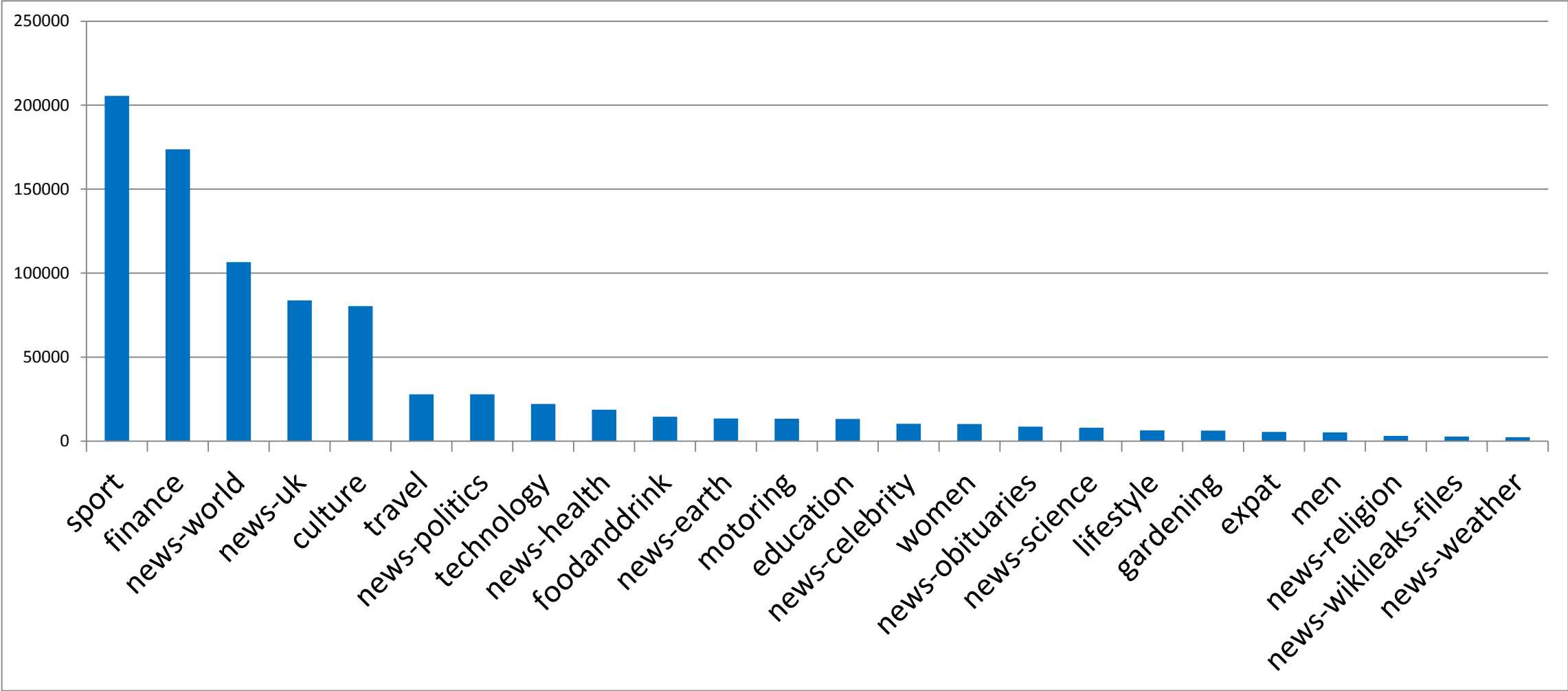
... and the Telegraph categories



It's **not**
so
easy.



Our selection: Functionally relevant, mutually exclusive categories derived from Telegraph categories



Finding **meaningful categories** for the Telegraph News was fun!

Lets go on and do a whole **text classification experiment**. Our **aim** is to **classify 1 million Telegraph News documents** with an ML algorithm.

While doing this we want to **find out...**

... if a **ML algorithm** will be able to classify the Telegraph news documents

... what are the **steps** we need to work out in order to make the ML algorithm work?

Handy for us: We will be able to **train** the **ML algorithm with the pre-classified data set** of the Telegraph News!

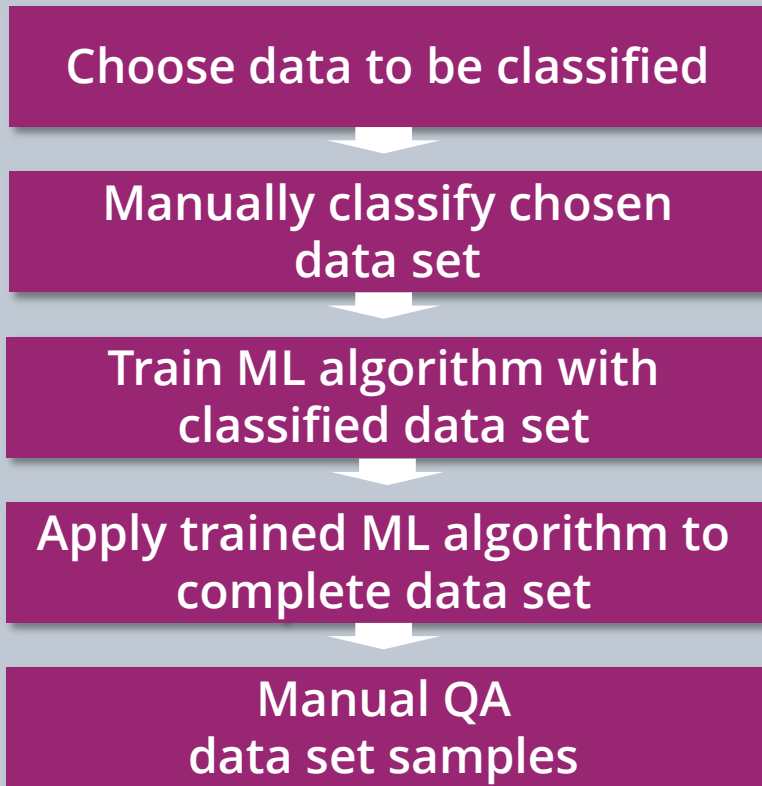


04

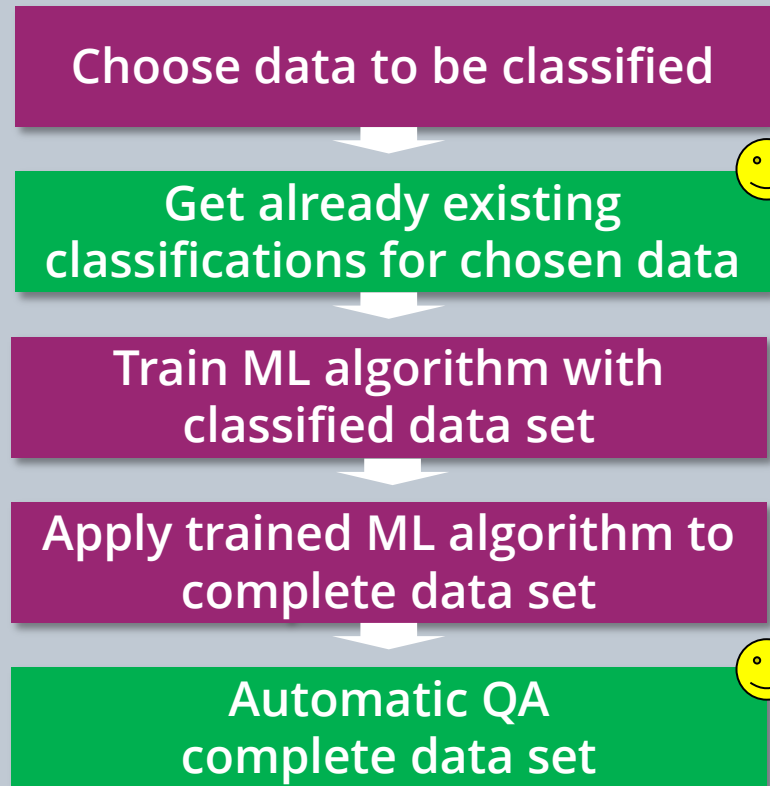
Text classification

Typical text classification projects and our experiment set-up

Typical set-up:
no classification scheme, no classified data



Our Telegraph experiment with
pre-classified documents

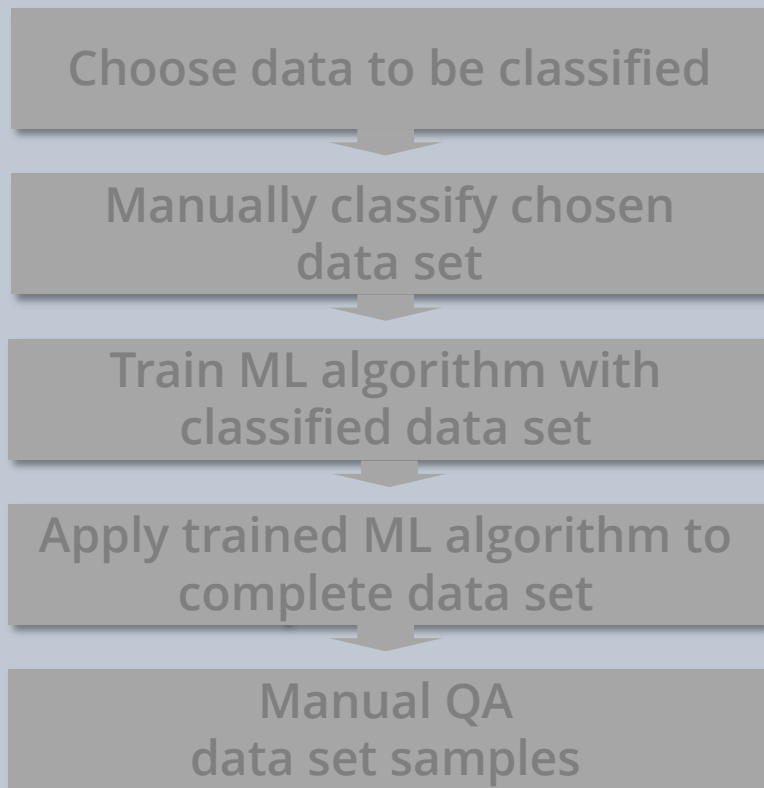


Advantages for us:

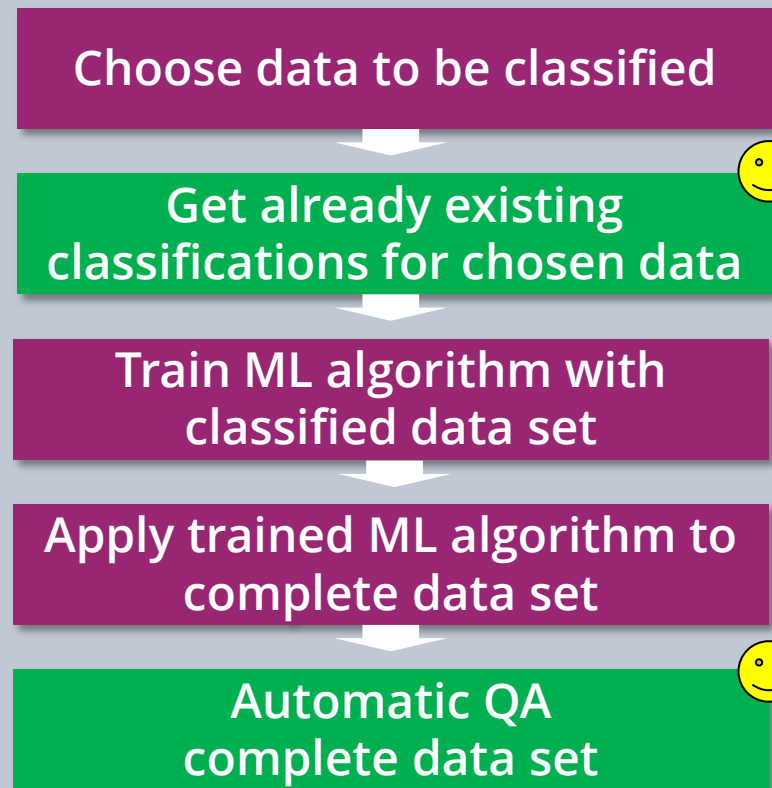
- ✓ No manual classification & QA necessary
- ✓ Existing classification scheme
- ✓ Playground easily set up
- ✓ Free to choose both manual data set & categories

Our experiment for the next 30 minutes

Typical set-up:
no classification scheme, no classified data



Our Telegraph experiment with
pre-classified documents



Our aims in the next 30 minutes:

- ✓ Train & apply the ML algorithm to 1 of Telegraph News
- ✓ See how well ML performs

This process sounds easy and very structured. The people in the audience who have already done text classification projects probably now that in reality, data can become **pretty challenging**.

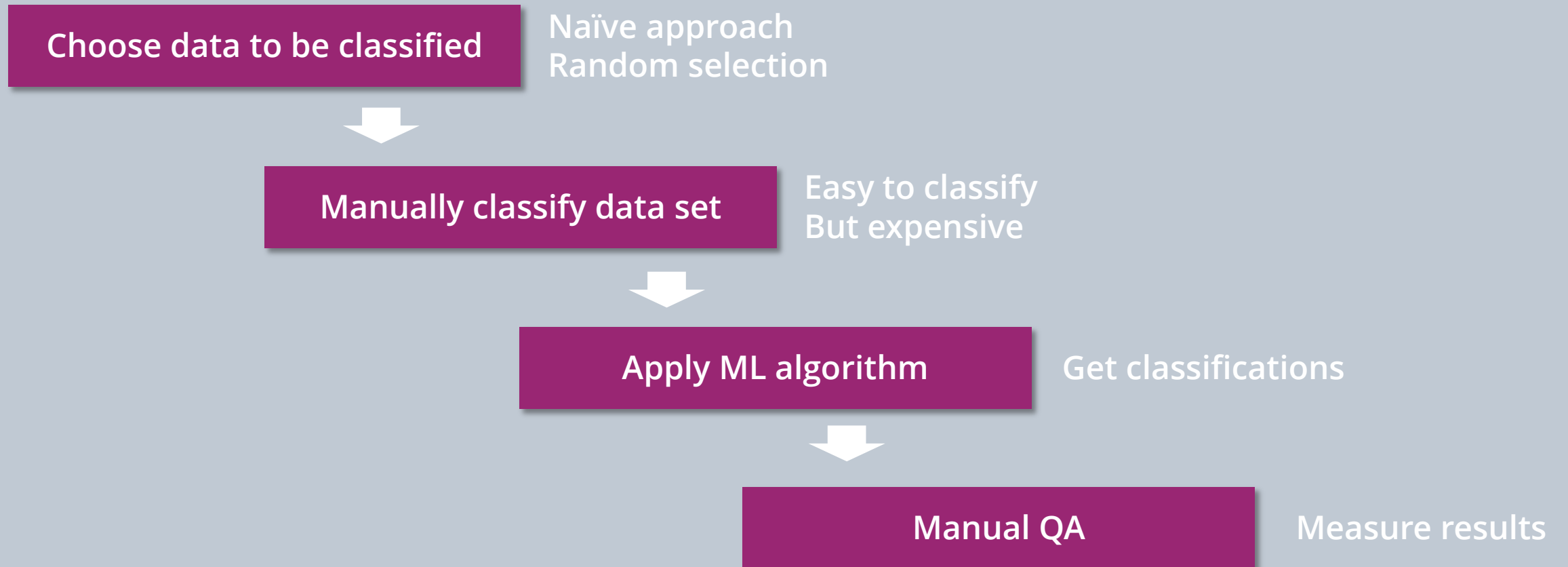
The next slides show you the process of **how** we classified 1 Million Telegraph news.

What is the **reality** we deal with?

And what are **good practices/our learnings?**

The devil is in the data

Getting started: Preparing data for and executing ML



The result is **BAD!**

WHY?

Lets take a step back and find out:

How does **ML WORK?**

How can I **MEASURE** its results?



ML algorithm explained – Support Vector Machine (SVM)

Machine learning is linear algebra

- Need to discretize first

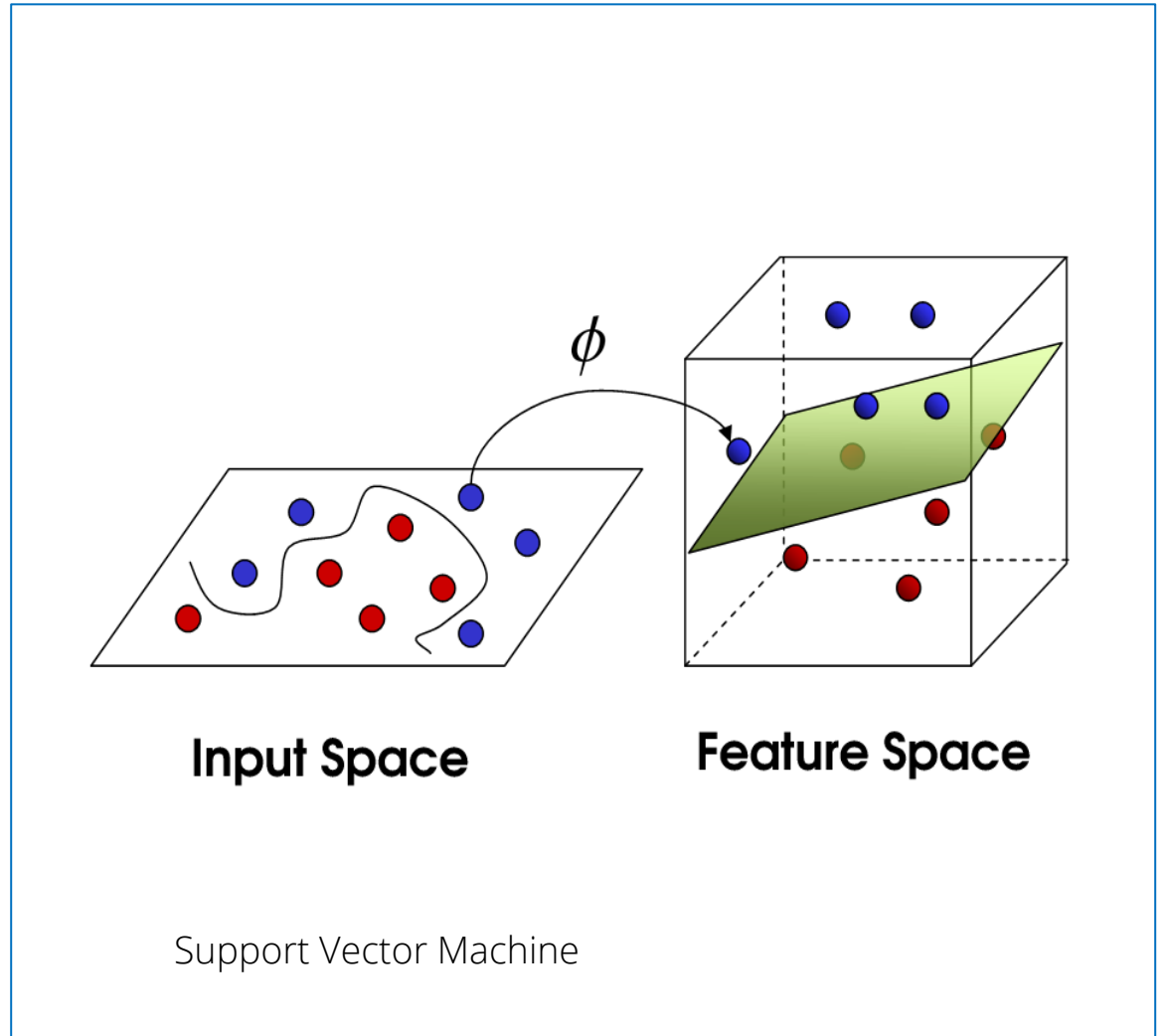
Categories are already discrete

More complicated for text

- Bag of words = detect words
- TF/IDF matrix = use document and total frequency

Many different possible learning models

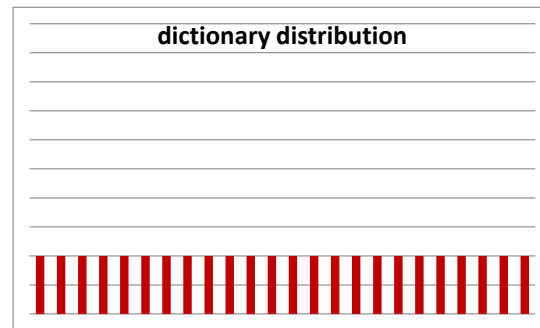
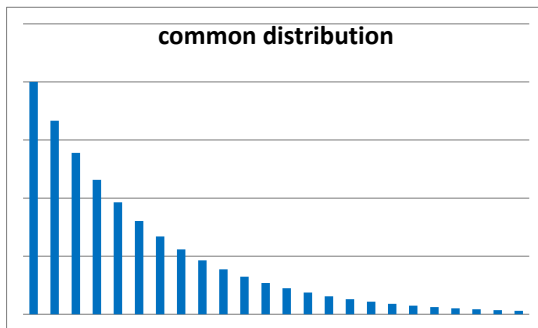
- Support Vector Machines (most popular)
- Neural Network
- Random forest
- Decision tree



Preparation of manually classified set

Choosing set for manual classification

- Select documents with highest word variability
 - Metric:
Word heterogeneity
= Number of words in all documents
(→ stopwords)
 - Even distribution
 - Long tail distribution
(→ many, many words use infrequently)
- Complicated: knapsack-like problem
- Use an approximate approach (like genetic algorithm)
- Crucial for all following tasks



1. Good situation:

The manually classified data set contains all the words of the complete data set.

Word heterogeneity in manual set			Word heterogeneity complete data set		
w01	w02	w03	w01	w02	w03
w04	w05	w06	w04	w05	w06
w07	w08	w09	w07	w08	w09
w10	w11	w12	w10	w11	w12
w13	w14	w15	w13	w14	w15
w16	w17	w18	w16	w17	w18

2. Not so good situation:

The manually classified data contains only a fraction of all the words in the complete data set

w19	w20	w21
w22	w23	w24
w25	w26	w27
w28	w29	w30
w31	w32	w33
...	...	w99

Complete set

Intelligently choose data set to be classified manually



Final data set available



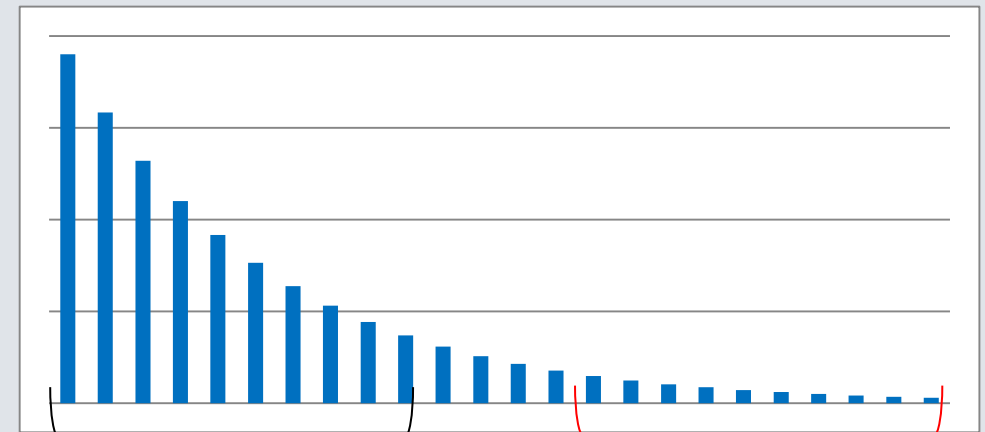
- Choose training data set in a way to create maximal word overlap with complete data set
- $W_M = \{ \text{words in training set} \}$
 $W_C = \{ \text{words in complete set} \}$
find maximum for $| W_C \cap W_M | = | W_M |$
- Improved approach: choose training set to minimize headlines with unknown words in complete data set
- Find minimum for $| C \cap \overline{W_M} |$
- More complicated, but worth it



Final data set not available



- Optimize for high variability and high usage



Select this

Don't select that

Measure classification quality: precision and recall

Precision

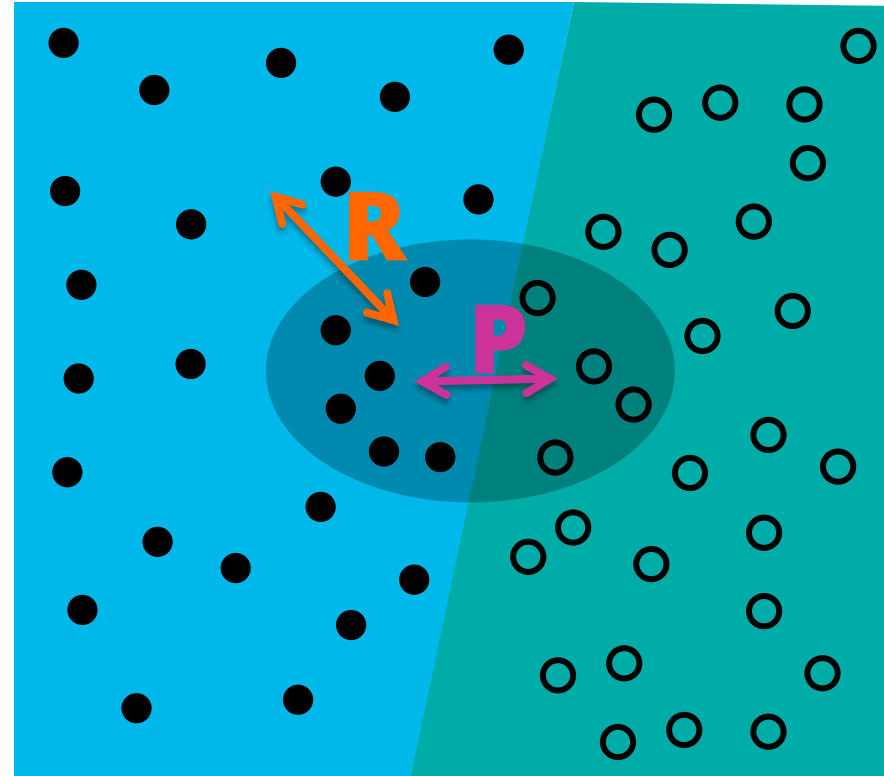
- „positive predictive value“
- Precision is the probability that a (randomly selected) retrieved document is classified correctly

Recall

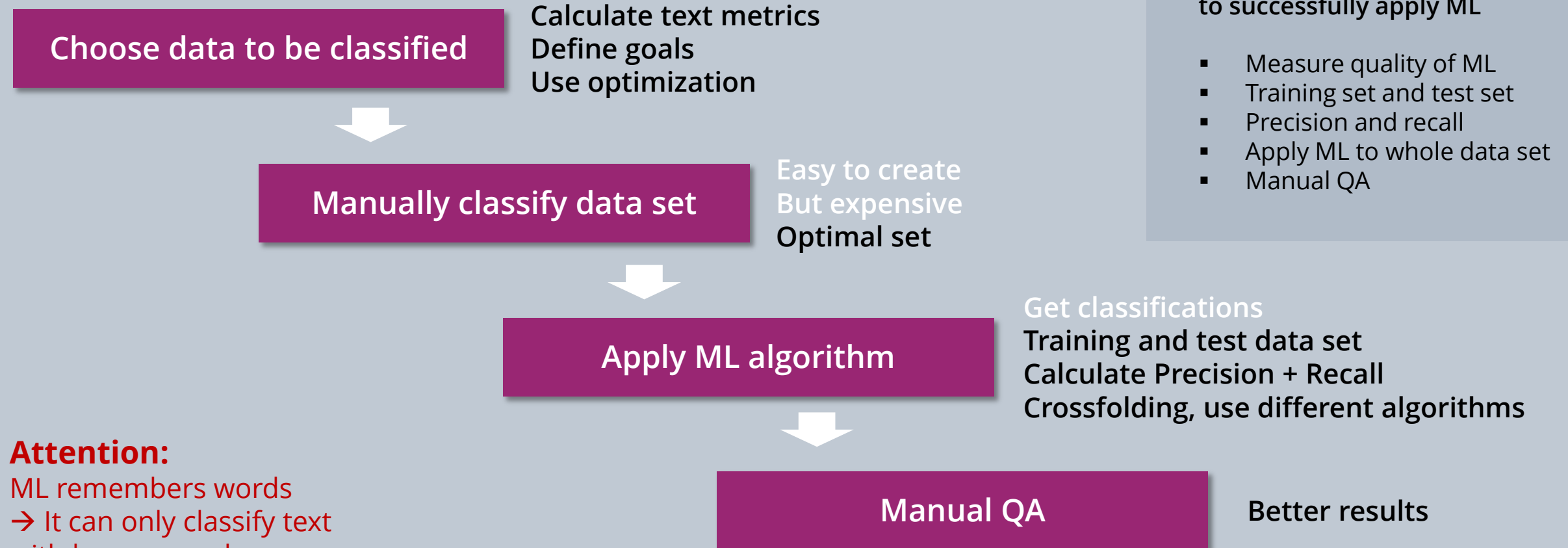
- Sensitivity or „true positive rate“
- Recall is the probability that a (randomly selected) classified document is found

Example

- Africa has very high precision for category „Africa“, but bad sensitivity (recall)



Now we know why the naive approach of preparing data for and executing ML is not enough. Lets try the following instead...



Attention:

ML remembers words
→ It can only classify text with known words

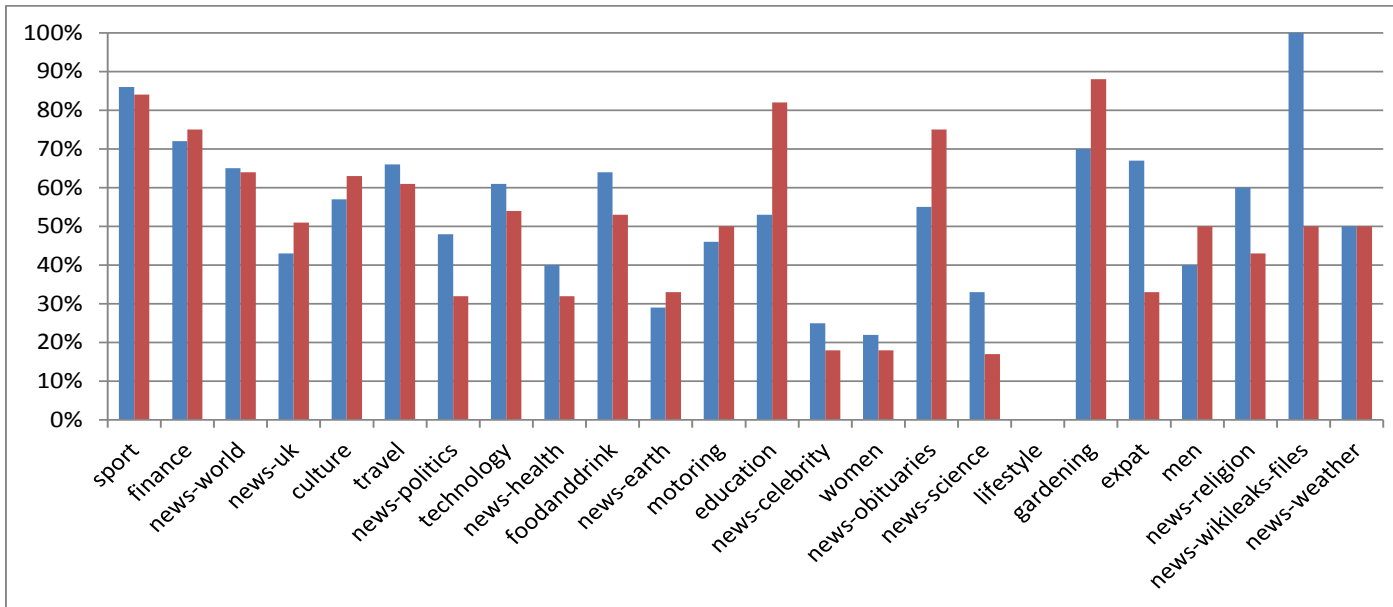
What we have done & achieved so far

1. Data cleaning and preparation: Docs with same headline but different classification removed
2. Category definition: Mutually exclusive, functionally relevant
3. Precision & Recall per categories and different training/test-sets

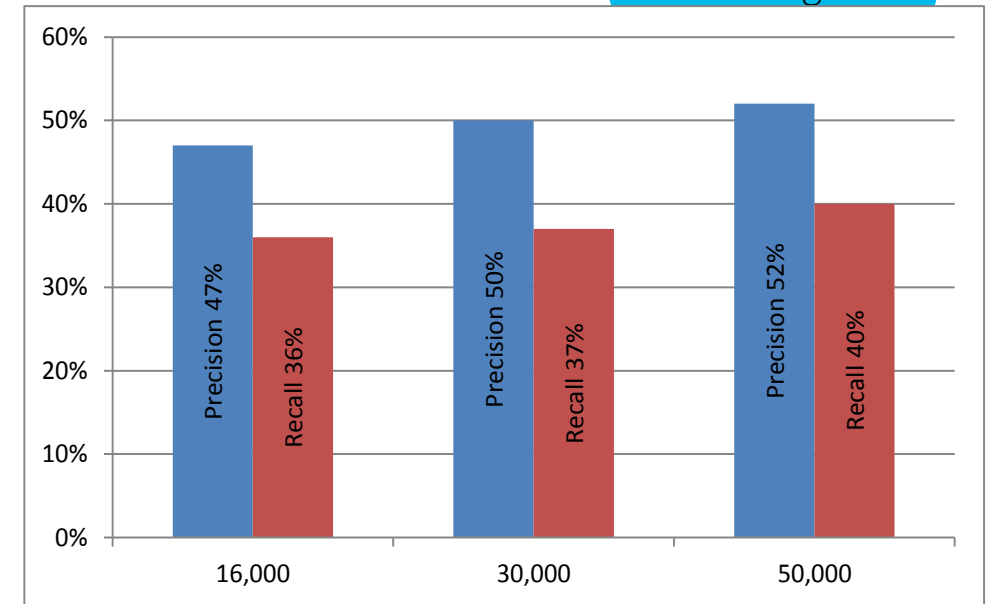
All eyes on eclipse of the Moon	news-UK
All eyes on eclipse of the Moon	news-science

4. Eliminating Longtail

abdurahim kerimbakiev
 abbot placid spearritt
 abdullahi sudi arale
 abduh alhamiri
 abib sarajuddin
 acer nethercott
 abdulli feghouel



Precision/recall per category



Precision/recall for different training/test- sizes

The result is **BETTER!**

Now...

... what **options** do you have if you don't have a pre-categorized data set **to train your ML?**

... or your **manually classified data set is too small?**



What to do about the things that can still go wrong

Manually classified data set is too small for training

- Data set is too heterogenous
- ML cannot detect patterns
- Bad precision and recall

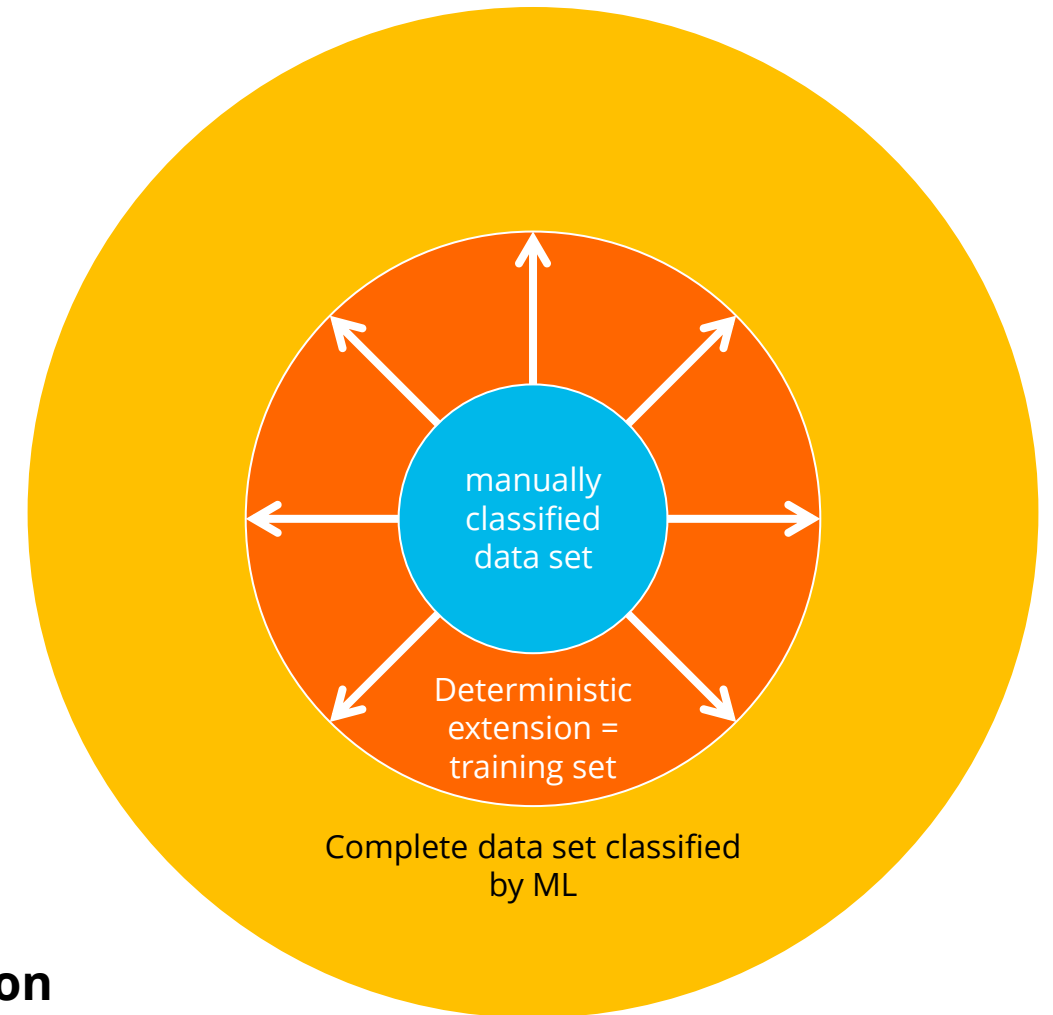
Extend data set

- Requires manual classification
- Too expensive

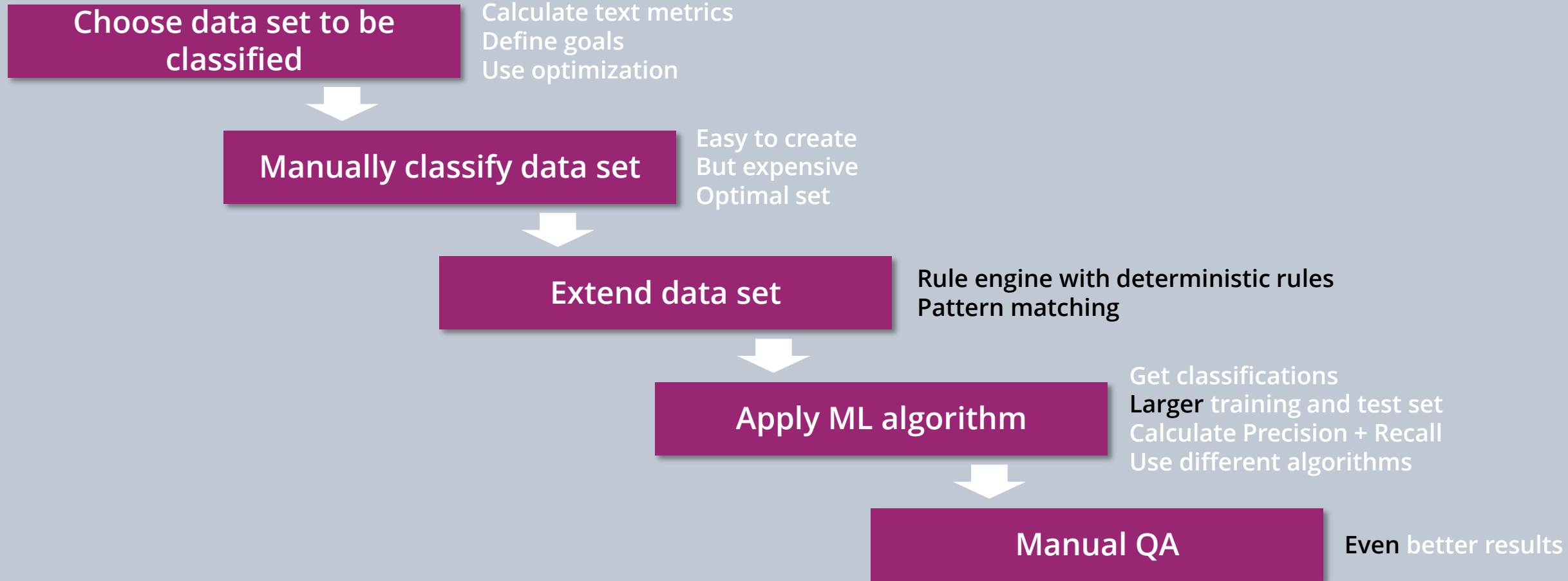
Try to understand structure of manual classification

- Find category-specific keywords
- Find patterns
- Use NLP etc.

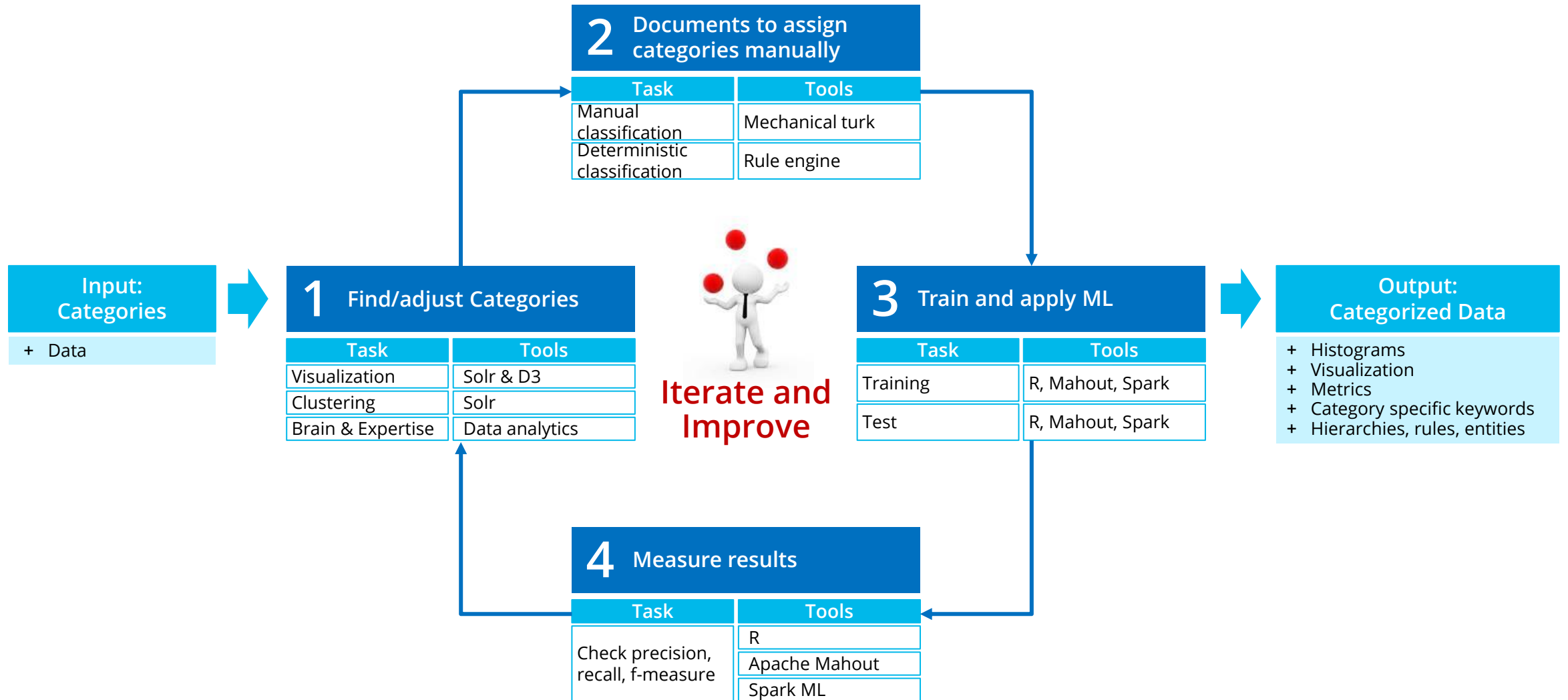
→ Extension of training set by deterministic classification



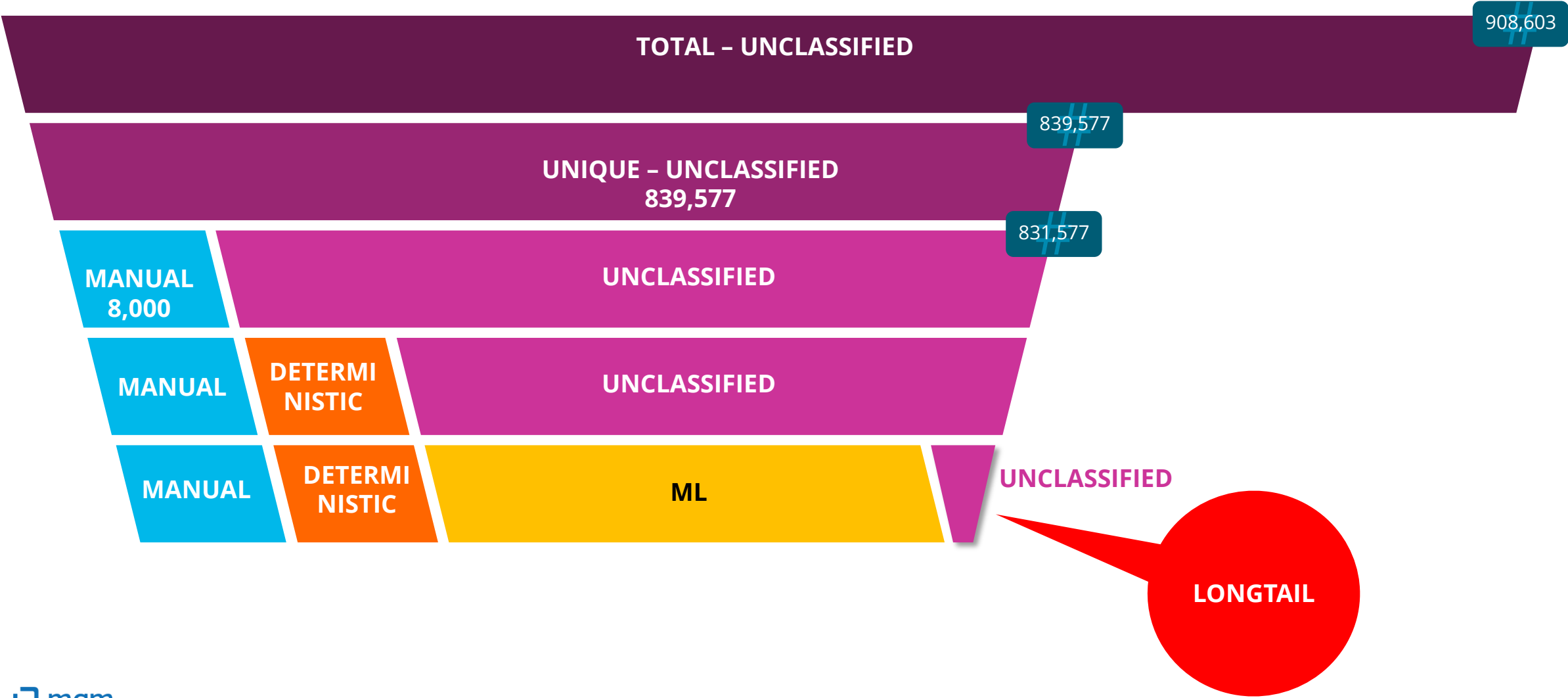
Improved approach with deterministic extension



Be prepared for a long journey: Often results get better incrementally



Classification cascade



Talking about longtail: Variability

Reasons for longtail

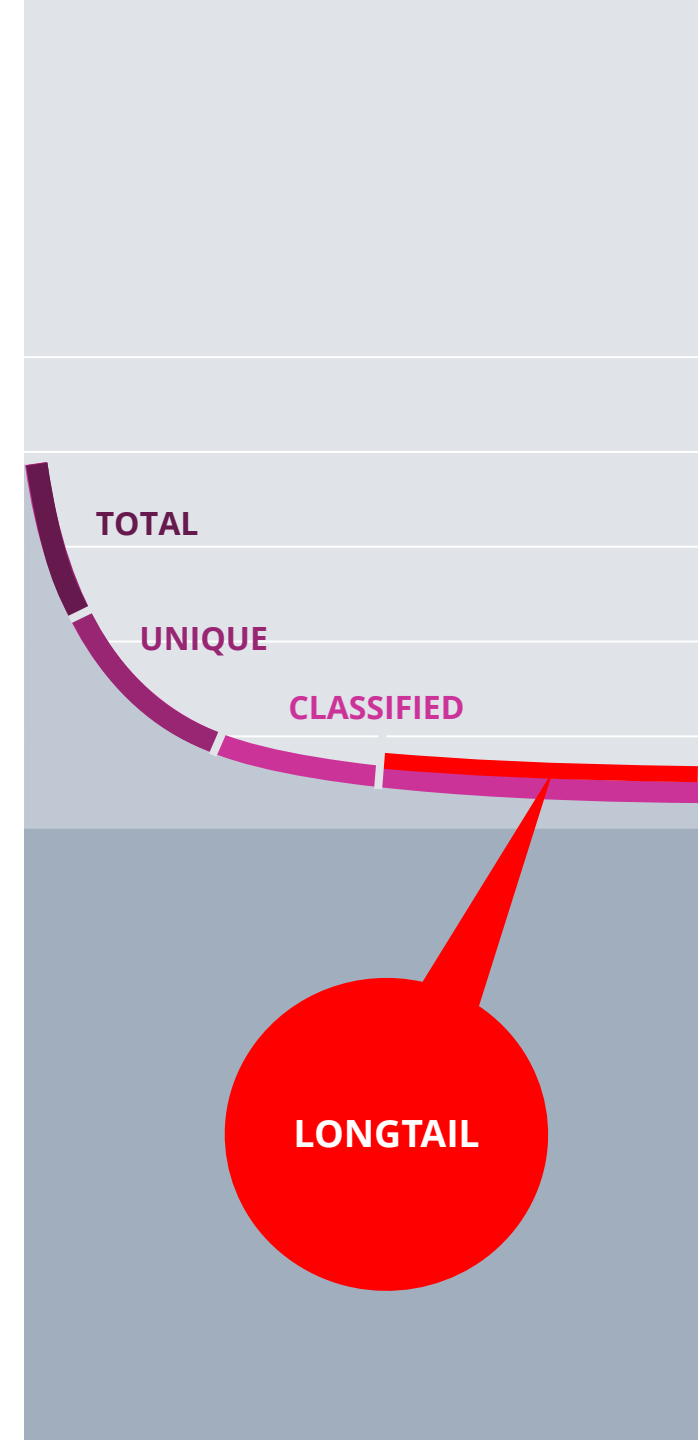
- Flat longtail via dictionary-type texts
- Decreasing longtail from domain specific language

Analyze the longtail

- Count words
- Measure heterogeneity

Elimination strategies

- Foreign language detection
- Eliminate typos (n-grams)
- Manual classification if not too many documents
- Put into separate category (aka „miscellaneous“)



How to make the decision in your data analytics project data-driven?
... measurable?

Metrics help making objective decisions during the project



The project is finished when your cost/benefit ratio (or its prediction) of classifying the longtail becomes negative.



05

Conclusion and outlook

10 Lessons learned

Really naïve

Run it on your notebook

Complex data structure & complicated classification scheme

Trying to understand ML

Thinking the functional specification is finished before the project is finished

Design manually classified data set very simple so ML will reach a high Precision/Recall

Sounds clever

Sounds naïve

Increase the ML test & training set manually and deterministically

Check data heterogeneity immediately, then choose technic

Get creative to find useful pre-categorized data

Mutually exclusive categories

Understand data qualitatively & quantitatively

Really clever

Outlook

Getting more pre-categorized data by

- Categories from other sources
- Semantic extraction
- NLP
- Meaning

Not yet analyzed text is everywhere

- Discretization helps in understanding
- Toolbox with ML. Deterministic rules helpful

Big potential: use already classified data to classify new data

Innovation Implemented.



München



Bamberg



Berlin



Boswil



Đà Nẵng



Dresden



Grenoble



Hamburg



Köln



Leipzig



Nürnberg



Prag

mgm technology partners GmbH

Frankfurter Ring 105a
80807 München
Tel.: +49 (89) 35 86 80-0
Fax: +49 (89) 35 86 80-288
<http://www.mgm-tp.com>

mgm consulting partners GmbH

Holländischer Brook 2
20457 Hamburg
Tel.: +49 (0) 40 / 80 81 28 20-0
Fax: +49 (0) 40 / 80 81 28 20-388
<http://www.mgm-cp.com>

