

INTRODUCTION TO DIGITAL HUMANITIES

APPLICATIONS OF DH: LANGUAGE AND LITERATURE TOPIC MODELLING

Dr. Thomas Köntges,
Digital Humanities, University of Leipzig

10 November 2015

WHAT IS TOPIC-MODELLING

- Collective knowledge continues to grow; it becomes more difficult to find what one is looking for
- Topic Models more than search&link-approach
- Zooming in and out is possible with TMs
- Topic Model algorithms do not require prior annotation or labelling of the texts
- Topic Models discover the hidden thematic structure in large archives of documents

WHAT IS TOPIC-MODELLING

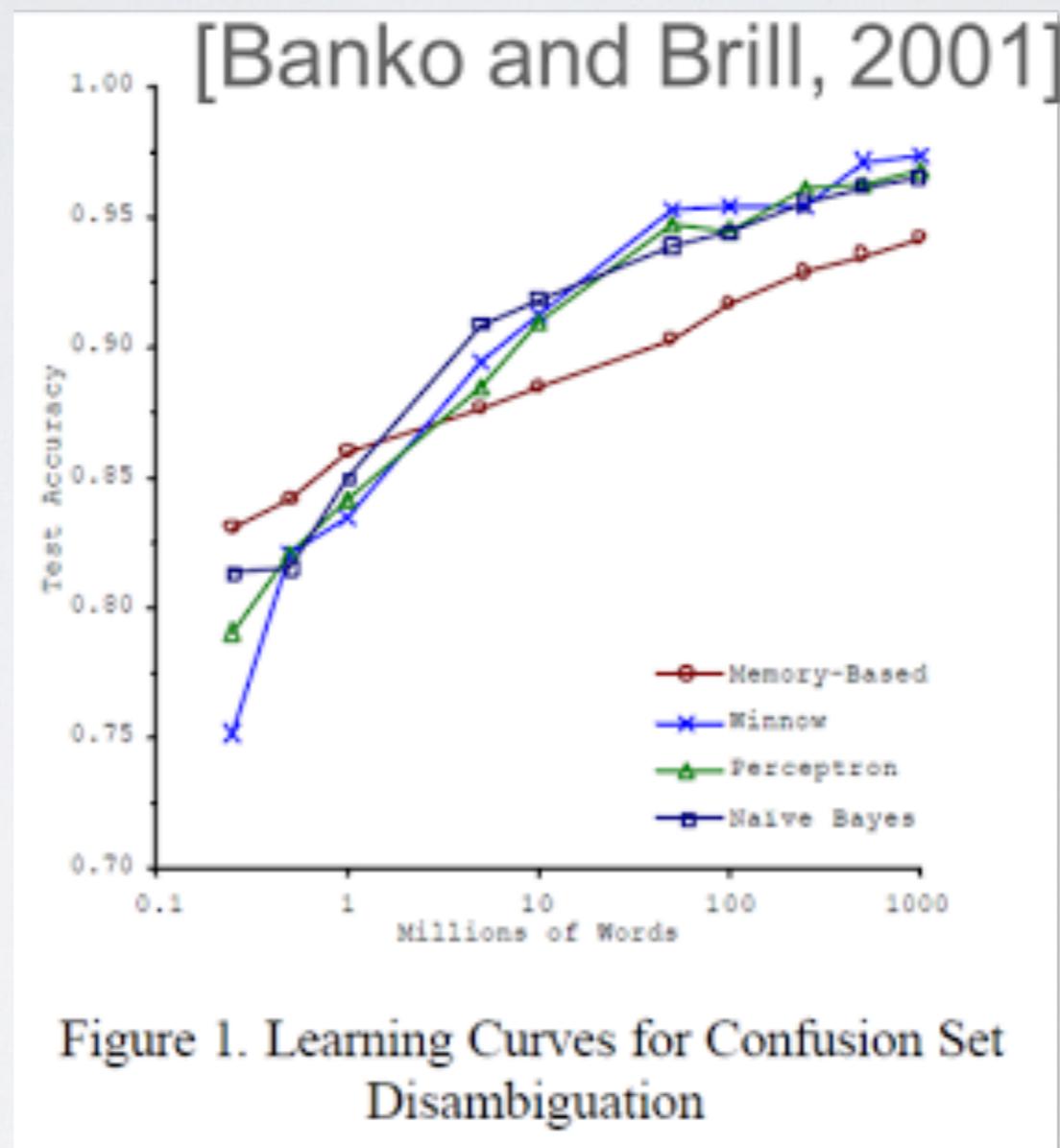
- A method to find clusters of words in large bodies of text
- Those clusters are called topics
- A topic is recurring pattern co-occurring words
- Topic models are probabilistic models that are often based on the number of topics in the corpus being assumed and fixed

WHAT IS TOPIC-MODELLING

[http://thomask81.github.io/INNZ_tm_vis/
index.html#topic=0&lambda=1&term=](http://thomask81.github.io/INNZ_tm_vis/index.html#topic=0&lambda=1&term=)

MORE (BETTER) DATA
AND
SIMPLE ALGORITHMS
OFTEN BEAT
COMPLEX ANALYTICS
MODELS

“We don’t have better algorithms. We just have more data ” - Peter Norvig (Google)



SO LET'S START WITH A
SIMPLE ALGORITHM

LATENT DIRICHLET ALLOCATION

STATISTICAL INFERENCE

**Deducing properties of
an underlying distribution
by analysis of data.**

LATENT DIRICHLET ALLOCATION (LDA)

**Assumes that one can find a
Dirichlet Distribution
by analysing the words in a corpus.**

DIRICHLET DISTRIBUTION

Distribution on a probability distribution

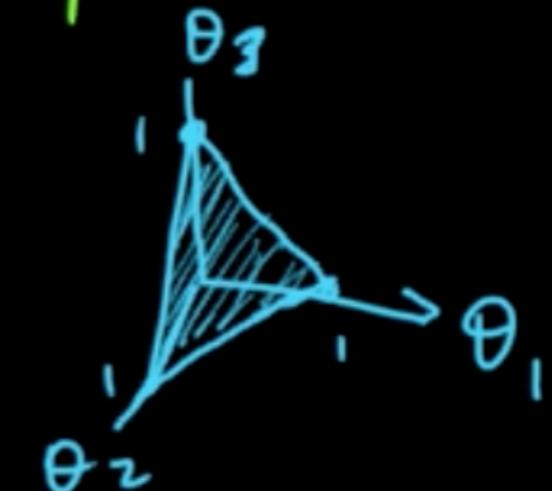
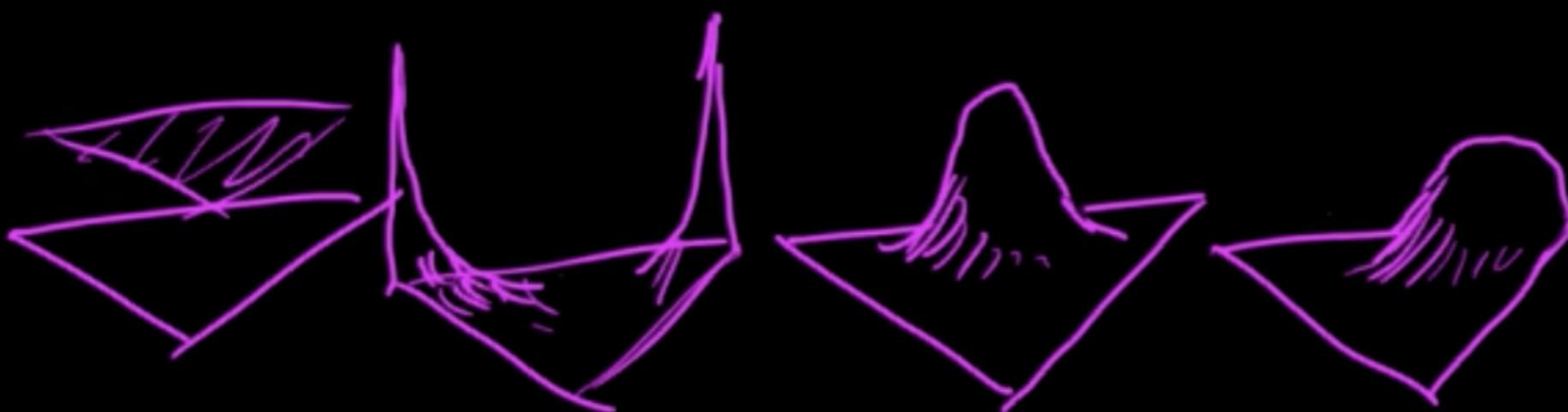
Dirichlet dist. $\theta = (\theta_1, \dots, \theta_n)$ $\alpha = (\alpha_1, \dots, \alpha_n), \alpha_i > 0$ $\alpha_0 = \sum_{i=1}^n \alpha_i$

$\theta \sim \text{Dir}(\alpha)$ $p(\theta) = \frac{1}{B(\alpha)} \prod_{i=1}^n \theta_i^{\alpha_i - 1} I(\theta \in S)$

$$\frac{1}{B(\alpha)} = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)}.$$

$$S = \left\{ x \in \mathbb{R}^n : x_i \geq 0, \sum_{i=1}^n x_i = 1 \right\}$$

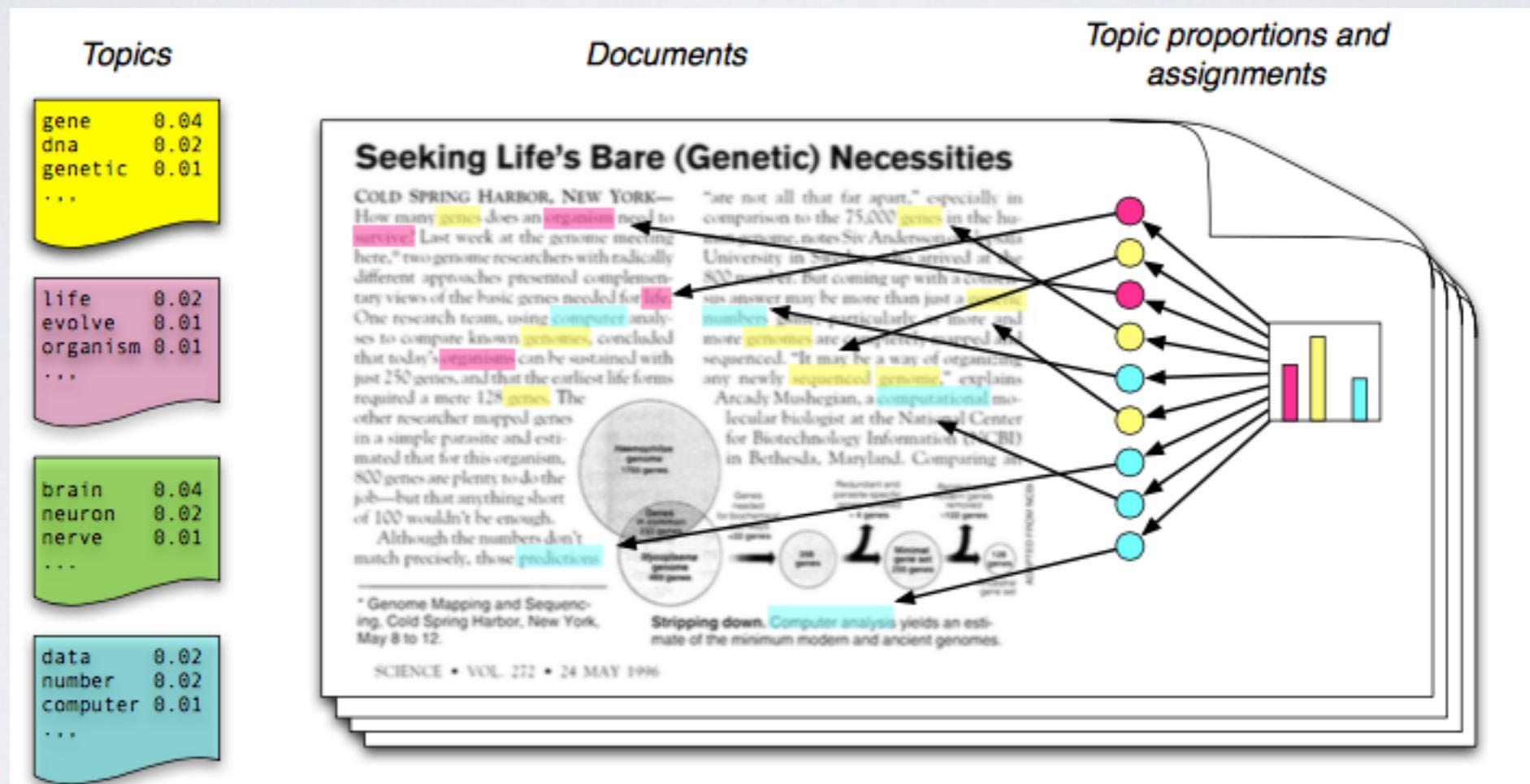
prob. simplex.



DIRICHLET DISTRIBUTION OF WORDS IN CORPUS

1. For each document draw a topic distribution
2. For each word in the document.
 - a. Draw a specific topic
 - b. Draw a word

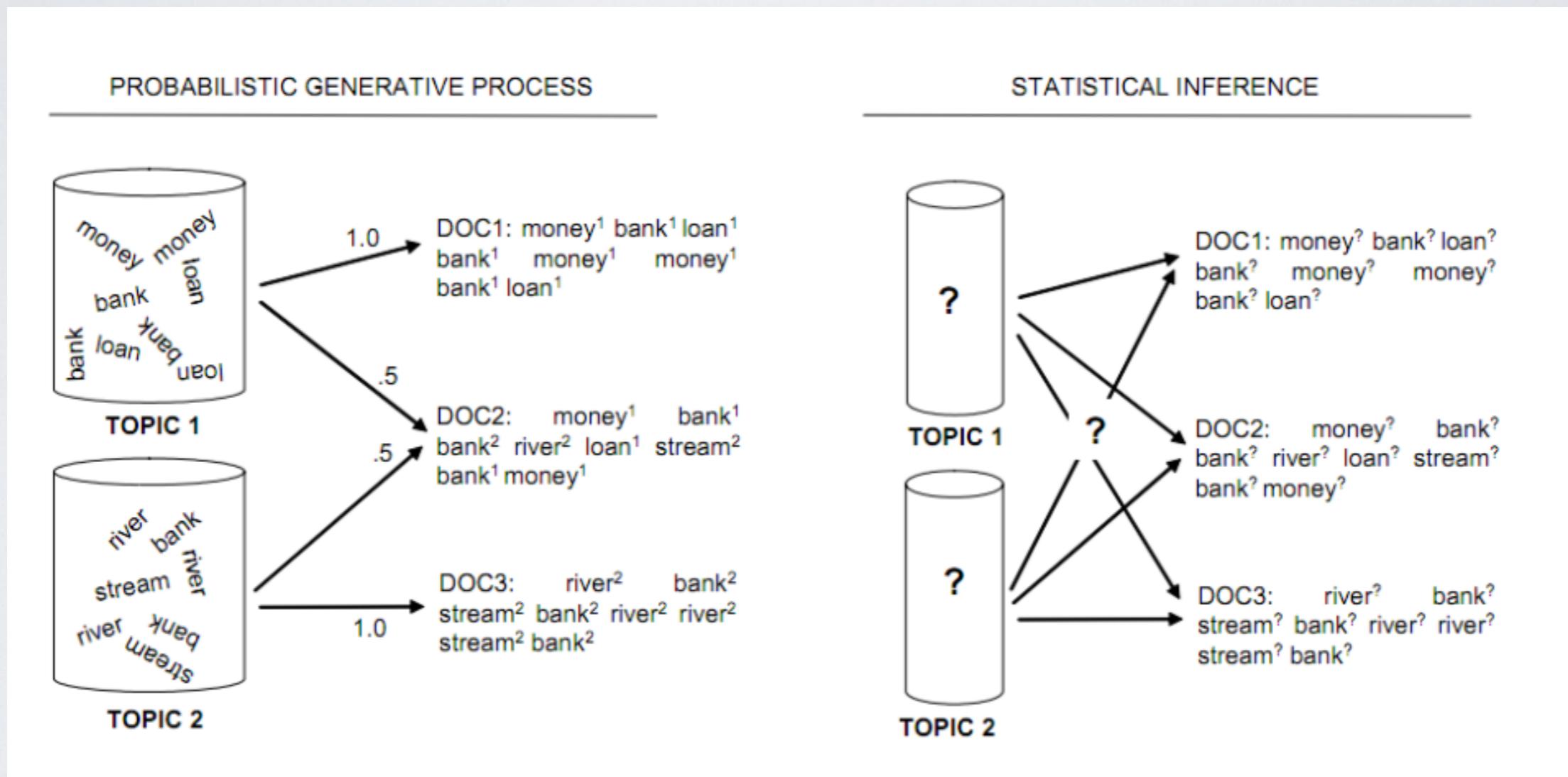
DIRICHLET DISTRIBUTION OF WORDS IN CORPUS



LATENT DIRICHLET ALLOCATION

- Simplification of how the documents in a dataset were created using Dirichlet Distribution
- Documents are Bag-of-words (order of the words in each document is irrelevant)
- Each corpus has words from a number of known topics
- We do not know which words belong to which topic
- “panel” in document A and “panel” in document B are not the same “word”

LATENT DIRICHLET ALLOCATION



LATENT DIRICHLET ALLOCATION

- The central goal of topic modelling is to automatically discover the topics from a collection of documents.
- The central inferential problem for LDA is determining the posterior distribution of the latent variables given the document

SUCCESS AND RESULTS OF LDA RELY ON APRIORI-SET VARIABLES

- Number of topics
- Number of iterations
- Normalisation of the data
- Stopwords

A PRACTICAL EXAMPLE

- <https://github.com/Thomask81/TopicModellingR>

WHAT DOES EACH VARIABLE DO?

- Search Text
- Collection
- Year
- Additional Stopwords
- Number of Topics
- Numbers of Terms Shown
- Iterations

ZOOMING IN AND OUT?

THE RISK OF NORMALISATION

- [http://thomask81.github.io/Greek_vis/
#topic=4&lambda=1&term=](http://thomask81.github.io/Greek_vis/#topic=4&lambda=1&term=)
- ArabicMorph
- ArabicTranslated
- ArabicTM

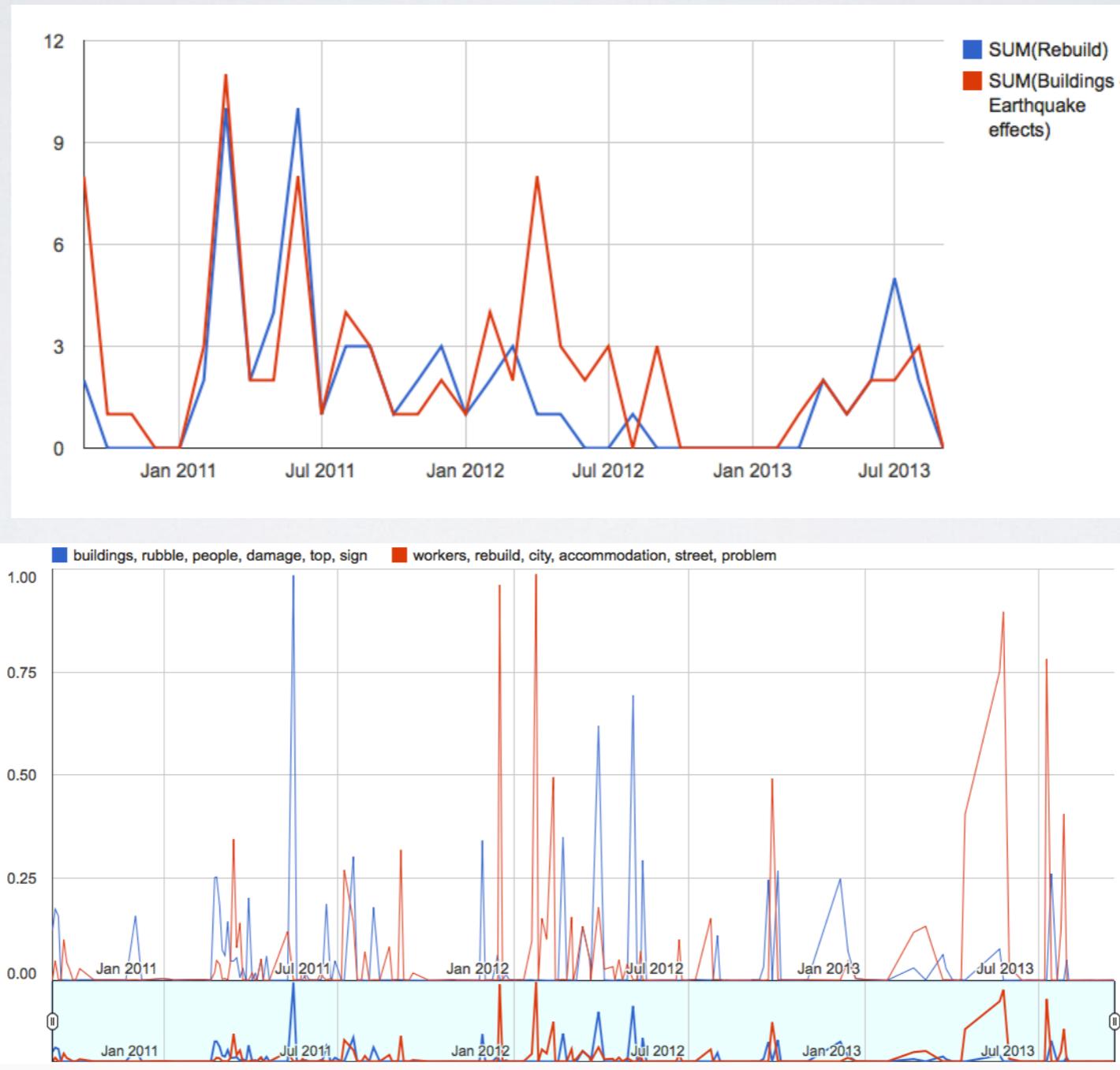
TOPIC MODELLING IS NOT AN END, BUT A MEANS TO AN END

- [https://www.google.com/fusiontables/DataSource?
docid=1HuEta4_g-
S0IX9LEow8F6llidt8s5u5_GKkur_so&pli=1#card:id=2](https://www.google.com/fusiontables/DataSource?docid=1HuEta4_g-S0IX9LEow8F6llidt8s5u5_GKkur_so&pli=1#card:id=2)
- [https://www.google.com/fusiontables/DataSource?
docid=1fppC5YlmhtxbtDWiEwbgq_N2FCQQKQKTtN8rPJV
9#card:id=2](https://www.google.com/fusiontables/DataSource?docid=1fppC5YlmhtxbtDWiEwbgq_N2FCQQKQKTtN8rPJV9#card:id=2)

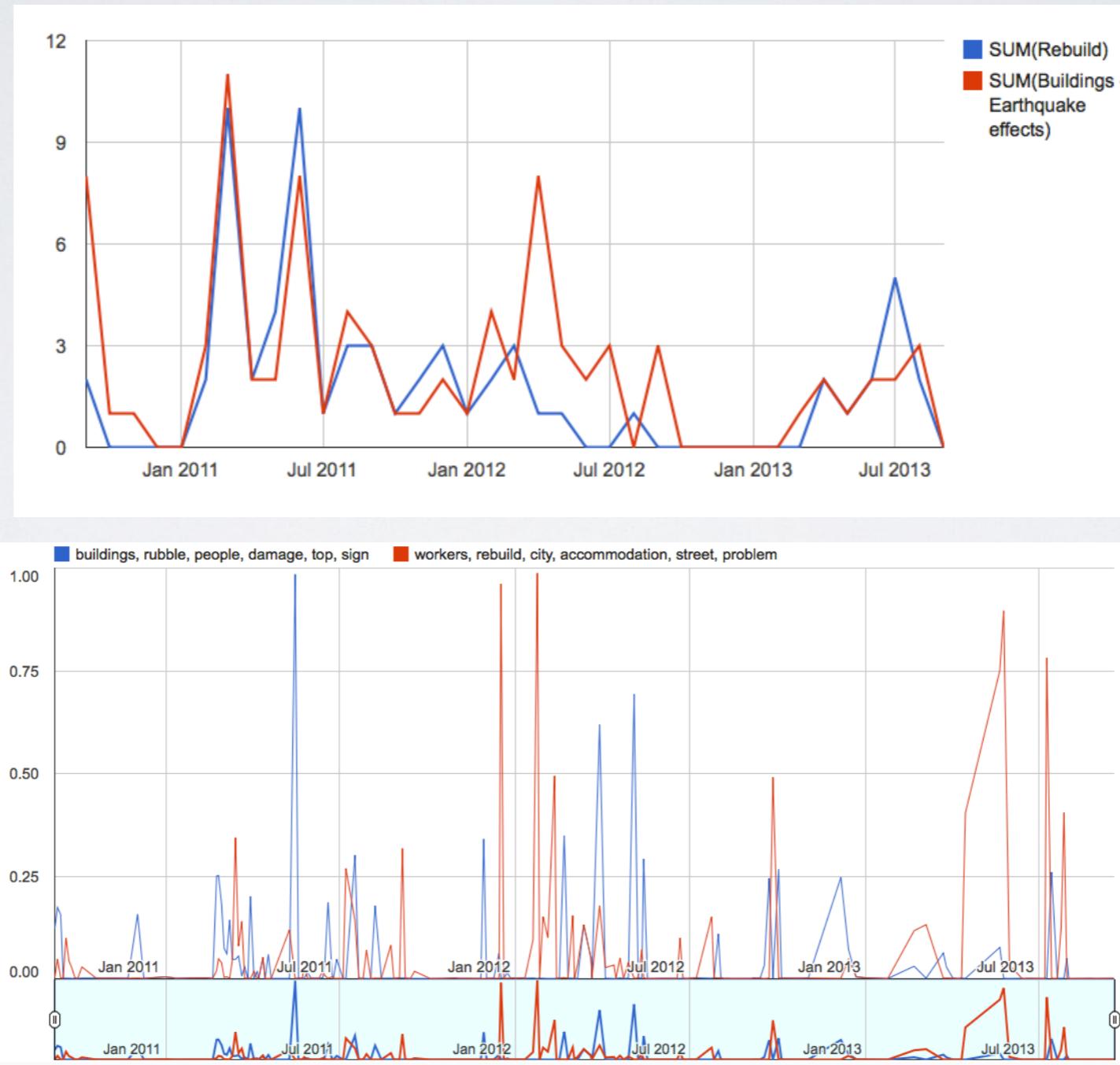
USE CASE: CARTOONS VERSUS ARTICLES

- Looking at the perception of the Canterbury EQ in published articles and published editorial cartoons
- Topic-Modelling is just one part of the overall research and is part of the quantitative analysis of the research data
- Quantitative analysis made it possible to reduce the sample from over 100,000 articles to around 2000 and reduce the number of cartoons from over 30,000 to under 1000.
- Combined with qualitative analysis (on the smaller sample)

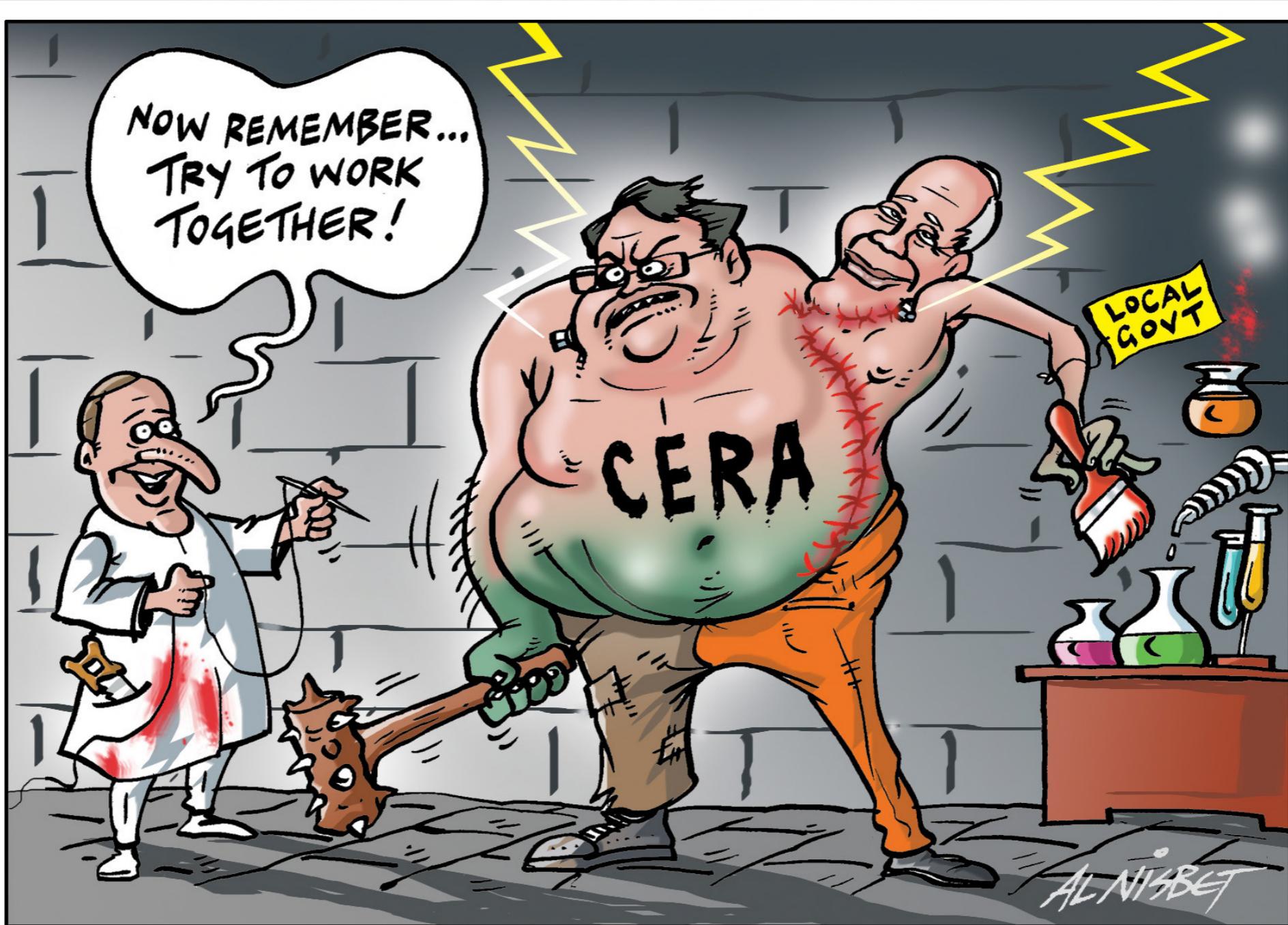
SUBJECT HEADINGS VS TM



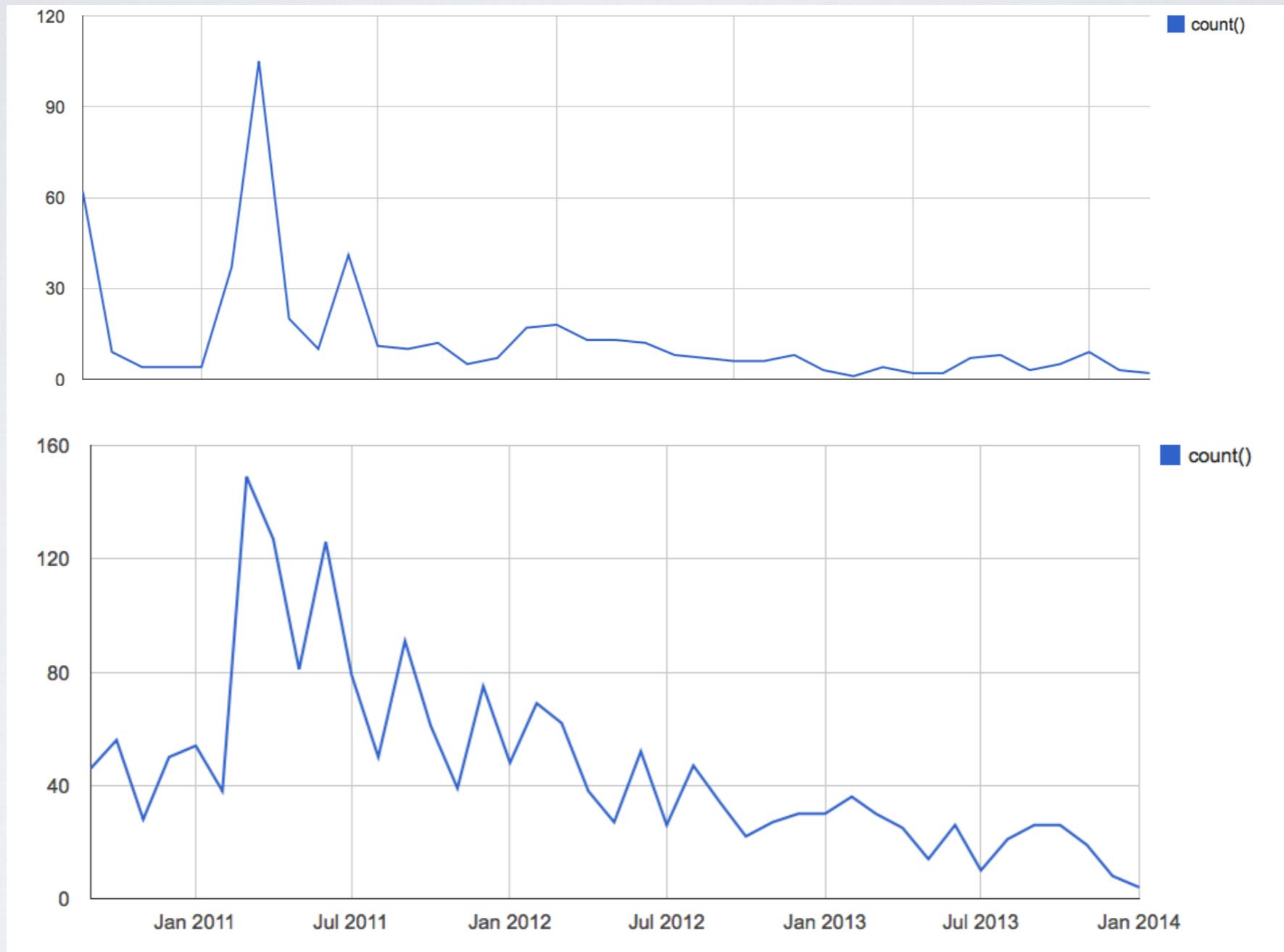
SUBJECT HEADINGS VS TM



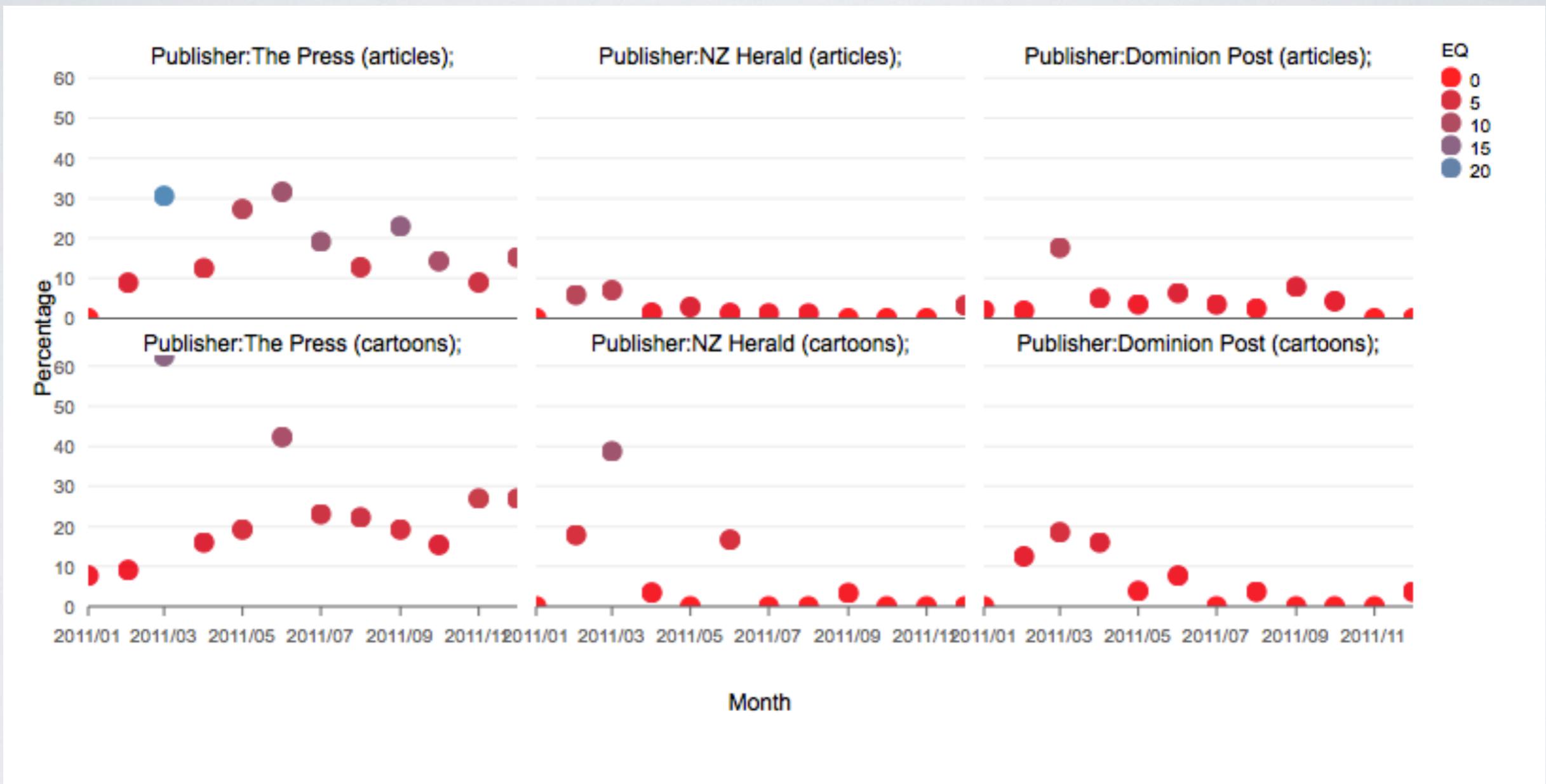
COMBINING DATA-SETS



COMBINING DATA-SETS



COMBINING DATA-SETS



QUESTIONS?

