

ANALYSIS OF UNSTRUCTURED TEXT DATA WITH TOPIC MODELS

WK RECH 2019, Frankfurt School of Finance & Management



PROF. DR. OLIVER MÜLLER

LS Wirtschaftsinformatik, insb. Data Analytics

- Since October 2018: Professor at the Department of Management Information Systems, Paderborn University
- 2016-2018: Associate Professor at the Business IT Department, IT University of Copenhagen
- 2011-2016: Assistant Professor at the Institute of Information Systems, University of Liechtenstein
- 2007-2011: PhD at European Research Center for Information Systems (ERCIS), Westfälische Wilhelms-Universität Münster



Research Interests

1. Using big data and machine learning to solve relevant business and societal problems
2. Analysis of unstructured data (e.g., text, images)
3. Acceptance and value of big data analytics

- I. Big Text Data
- II. Fundamentals of Topic Modeling
- III. Topic Modeling Walkthrough

BIG TEXT DATA



A Very Short History Of Big Data

+50k views in the last 24 hours



Expect 'Pokémon GO' To Make More Halloween-Like Events After Huge 133% Revenue Jump

+16k views in the last 24 hours



WWE Smackdown Results: Winners, Analysis, Reaction and Highlights From November 1

Apple iOS 10.1.1: Should You Upgrade?



IBM Voice: How Blockchain Could Help To Make The Food We Eat Safer... Around The World

+12k views in the last 24 hours

Tech



MAY 9, 2013 @ 09:45 AM 168,556 VIEWS

A Very Short History Of Big Data



Gil Press, CONTRIBUTOR

I write about technology, entrepreneurs and innovation. [FULL BIO](#)

Opinions expressed by Forbes Contributors are their own.

The story of how data became big starts many years before the current buzz around big data. Already seventy years ago we encounter the first attempts to quantify the growth rate in the *volume of data* or what has popularly been known as the “information explosion” (a term first used in 1941, according to the *Oxford English Dictionary*). The following are the major milestones in the history of sizing data volumes plus other “firsts” in the evolution of the idea of “big data” and the “information explosion.”

Last Update: December 21, 2013

1944 Fremont Rider, Wesleyan University Librarian, publishes *The Scholar and the Future of the Research Library*. He estimates that American university libraries were doubling in size every sixteen years. Given this growth



Source: <http://www.forbes.com/sites/gilpress/2013/05/09/a-very-short-history-of-big-data/#41b6f62055da>

Forbes

LOG IN

YOUR READING LIST



A Very Short History Of Big Data

+50k views in the last 24 hours



Expect 'Pokémon GO' To Make More Halloween-Like Events After Huge 133% Revenue Jump

+16



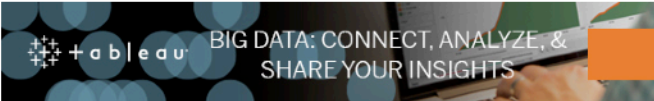
+22



Around The World

+12k views in the last 24 hours

Tech




BIG DATA: CONNECT, ANALYZE, & SHARE YOUR INSIGHTS

MAY 9, 2013 @ 09:45 AM 168,556 VIEWS

A Very Short History Of Big Data





Gil Press, CONTRIBUTOR

I write about technology, entrepreneurs and innovation. [FULL BIO](#)

Opinions expressed by Forbes Contributors are their own.

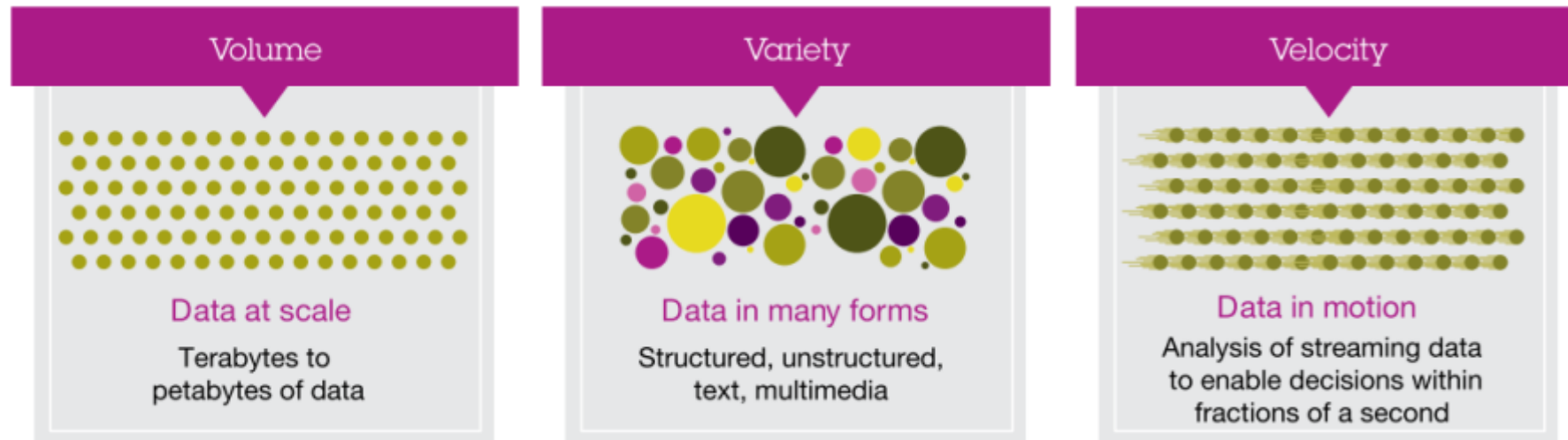
story of how data became big starts many years before the current buzz and big data. Already seventy years ago we encounter the first attempts to identify the growth rate in the *volume of data* or what has popularly been known as the “information explosion” (a term first used in 1941, according to the *Oxford English Dictionary*). The following article is a brief history of sizing data volumes plus other “firsts” in the evolution of the idea of “data” and observations pertaining to data or information explosion.

Update: December 21, 2013

4 Fremont Rider, Wesleyan University Librarian, publishes *The Scholar and the Future of the Research Library*. He estimates that American university libraries were doubling in size every sixteen years. Given this growth

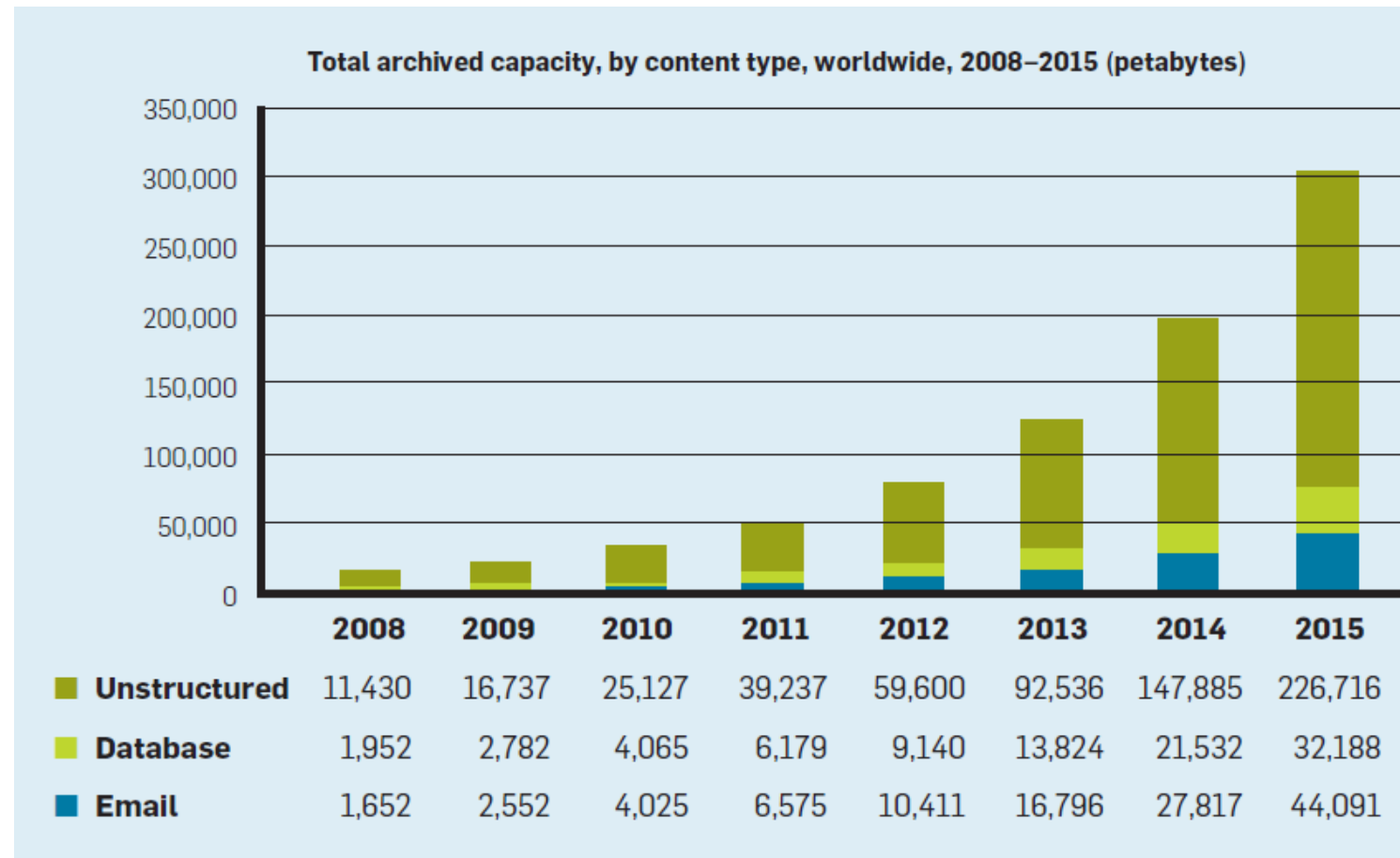
The 3 "V"s of Big Data (Laney, 2001)

Source: <http://www.forbes.com/sites/gilpress/2013/05/09/a-very-short-history-of-big-data/#41b6f62055da>



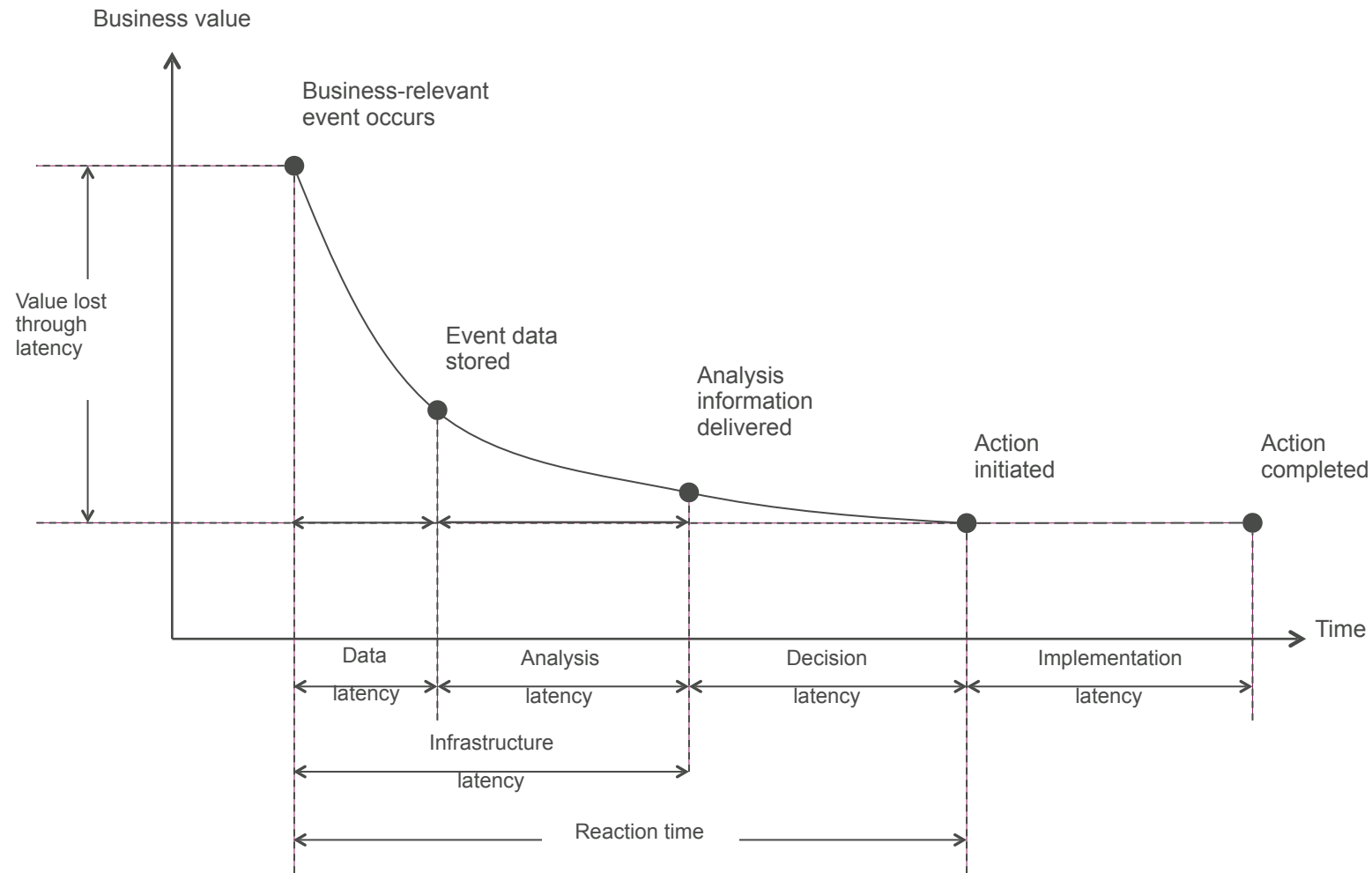
Source: Laney (2001), IBM (2012)

Volume and Variety



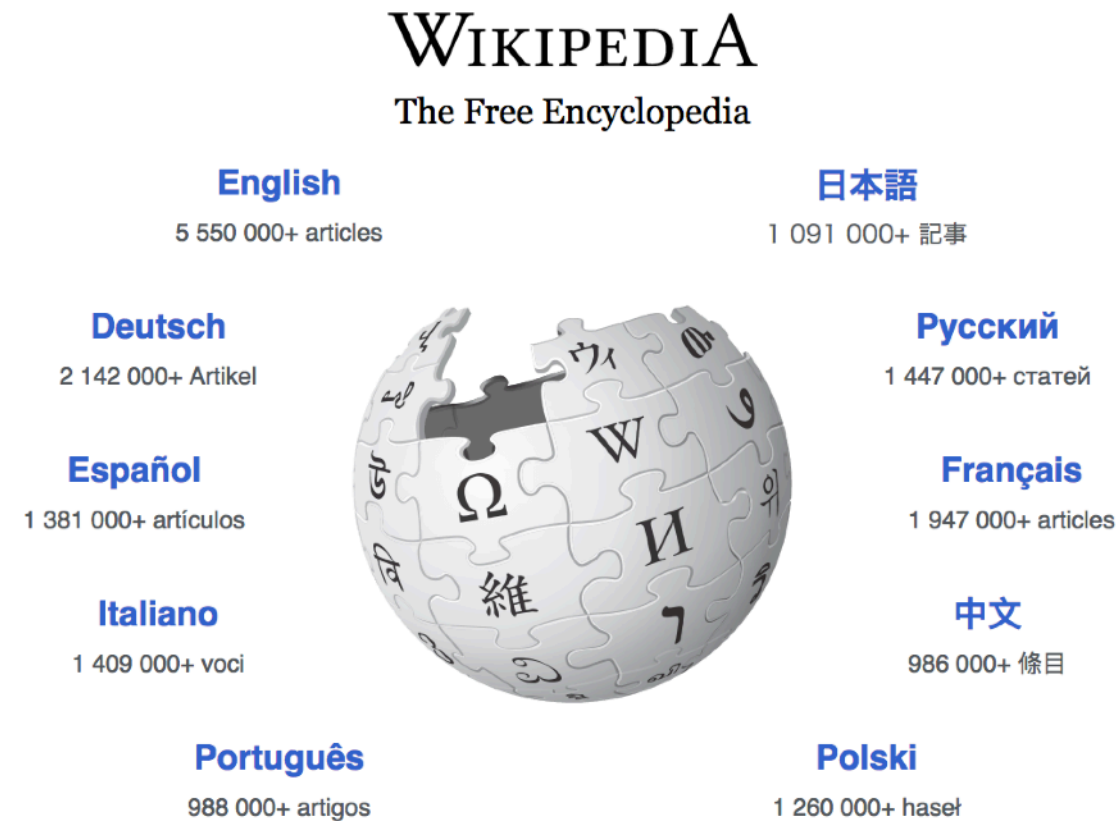
Source: Dhar (2013)

Velocity



Source: Zur Mühlen & Shapiro (2010)

Big or not?



Big or not?



Why is Text Analytics Difficult?

- **Text is Messy**
 - Cannot easily be represented in rows and columns of tables
 - Has complex linguistics structures that differ across languages
- **Text is Dirty**
 - Lots of words that are in no dictionary (e.g., spelling mistakes, slang, abbreviations, technical terms)
- **Text is Ambiguous**
 - Meaning of words depends on context

Text is Messy

country	year	cases	population
Afghanistan	1999	31737	1508071
Afghanistan	2000	2566	20595360
Brazil	1999	31737	17206362
Brazil	2000	80488	174604898
China	1999	212258	127091272
China	2000	210766	1280423583

variables

country	year	cases	population
Afghanistan	1999	31737	1508071
Afghanistan	2000	2566	20595360
Brazil	1999	31737	17206362
Brazil	2000	80488	174604898
China	1999	212258	127091272
China	2000	210766	1280423583

observations

country	year	cases	population
Afghanistan	1999	31737	1508071
Afghanistan	2000	2566	20595360
Brazil	1999	31737	17206362
Brazil	2000	80488	174604898
China	1999	212258	127091272
China	2000	210766	1280423583

values

Lorem Ipsum

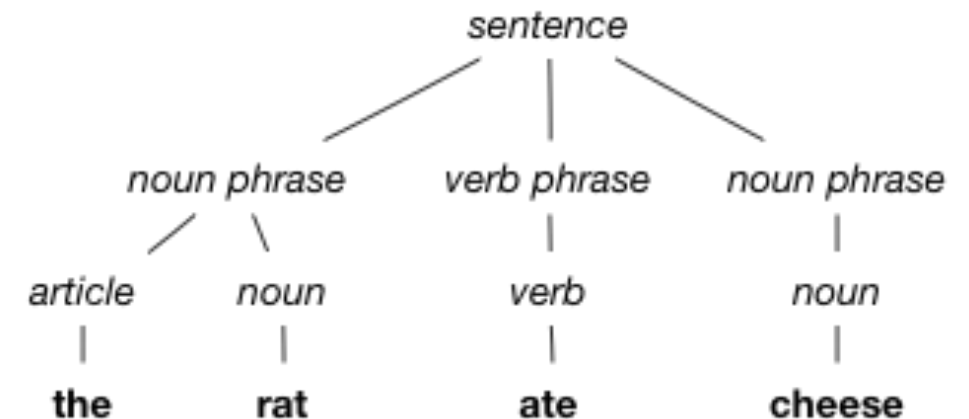
"Neque porro quisquam est qui dolorem ipsum quia dolor sit amet, consectetur, adipisci velit..."
 "There is no one who loves pain itself, who seeks after it and wants to have it, simply because it is pain..."

What is Lorem Ipsum?

Lorem Ipsum is simply dummy text of the printing and typesetting industry. Lorem Ipsum has been the industry's standard dummy text ever since the 1500s, when an unknown printer took a galley of type and scrambled it to make a type specimen book. It has survived not only five centuries, but also the leap into electronic typesetting, remaining essentially unchanged. It was popularised in the 1960s with the release of Letraset sheets containing Lorem Ipsum passages, and more recently with desktop publishing software like Aldus PageMaker including versions of Lorem Ipsum.

Why do we use it?

It is a long established fact that a reader will be distracted by the readable content of a page when looking at its layout. The point of using Lorem Ipsum is that it has a more-or-less normal distribution of letters, as opposed to using 'Content here, content here', making it look like readable English. Many desktop publishing packages and web page editors now use Lorem Ipsum as their default model text, and a search for 'lorem ipsum' will uncover many web sites still in their infancy. Various versions have evolved over the years, sometimes by accident, sometimes on purpose (injected humour and the like).

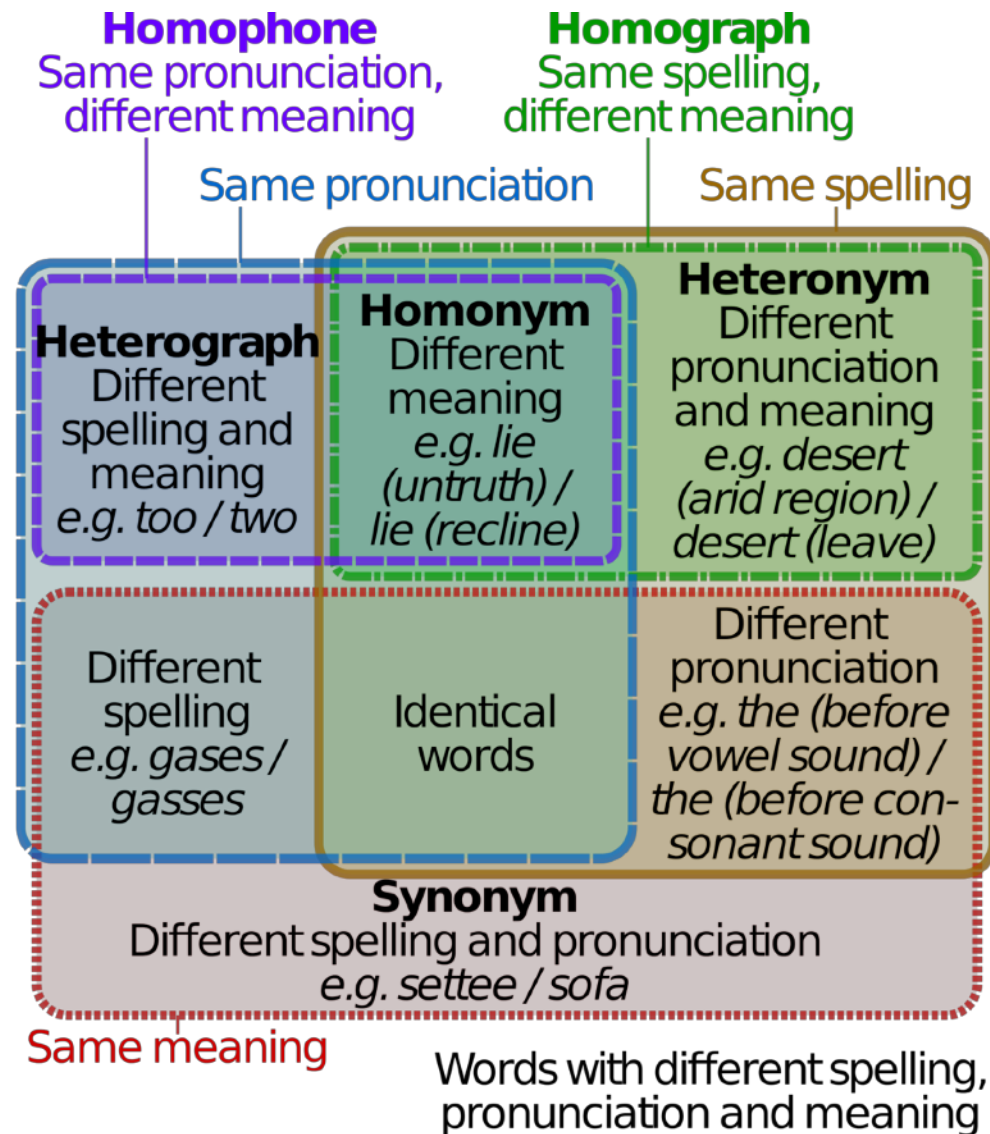


A scatter plot showing the relationship between Rank (X-axis) and 2013 Q4 Relative Frequency Per Million Words (Y-axis). The X-axis ranges from 0 to 60, and the Y-axis ranges from 0 to 50,000. The plot shows a steep decline in frequency as rank increases, with the word 'i' being the most frequent.

Rank	2013 Q4 Rel Freq Per Million Words	Word
1	47000	i
2	23000	the
3	20000	you
4	17000	a
5	16000	my
6	12000	and
7	11000	me
8	10000	that
9	9000	for
10	8000	one
11	7000	year
12	6000	at
13	5000	had
14	4000	in
15	3000	to
16	2000	of
17	1500	and
18	1000	the
19	800	you
20	600	are
21	500	is
22	400	on
23	300	at
24	200	the
25	150	of
26	100	in
27	80	to
28	60	and
29	40	the
30	30	you
31	20	are
32	15	is
33	10	on
34	8	at
35	6	the
36	4	of
37	3	in
38	2	to
39	1	and
40	1	the
41	1	you
42	1	are
43	1	is
44	1	on
45	1	at
46	1	the
47	1	of
48	1	in
49	1	to
50	1	and
51	1	the
52	1	you
53	1	are
54	1	is
55	1	on
56	1	at
57	1	the
58	1	of
59	1	in
60	1	to

AFAIK	As Far as I Know	MMB	Message Me Back
AFK	Away from Keyboard	msg	Message
ASL	Age/Sex/Location?	MYOB	Mind Your Own Business
ATM	At The Moment	N/A	Not Available
b/c	Because	NC	No Comment
b/w	Between	ne1	Anyone
b4	Before	NM	Not much
BBIAB	Be Back in a bit	noob	Newbie
BBL	Be back later	NP	No Problem
BFF	Best Friends Forever	NTN	No Thanks Needed
BRB	Be Right Back	OMG	Oh My Gosh
BTW	By The Way	OMW	On My Way
CTN	Can't Talk Now	OT	Off Topic
CYE	Check Your E-mail	PC	Personal Computer
dl	Download	pls	Please
ETA	Estimated Time of Arrival	POS	Parent Over Shoulder
FWIW	For What It's Worth	ppl	People
FYI	For Your Information	qt	Cutie
GG	Good Game	re	Regarding
GJ	Good Job	SMH	Shaking my head
GL	Good Luck	Sry	Sorry
gr8	Great	TBA	To Be Announced
GTG	Got To Go	TBC	To Be Continued
GMV	Got My Vote	TC	Take Care
HTH	Hope this helps	thx	Thanks
hw	Homework	TIA	Thanks In Advance
IAC	In Any Case	TLC	Tender Loving Care
IC	I See	TMI	Too Much Information
IDK	I Don't Know	TTFN	Ta-ta For Now
IIRC	If I Remember Correctly	TYTL	Talk To You Later
IKR	I Know, Right?	txt	Text
IM	Instant Message	TY	Thank You
IMO	In My Opinion	w/e	Whatever
IMHO	In My Humble Opinion	w/o	Without
IRL	In Real Life	W8	Wait
J/K	Just kidding	XOXO	Hugs and kisses
K	OK	Y	Why
L8	Late	YNt	Why Not
I8r	Later	YOLO	You Only Live Once
LMK	Let Me Know	YW	You're Welcome
LOL	Laughing Out Loud	ZZZ	Sleeping

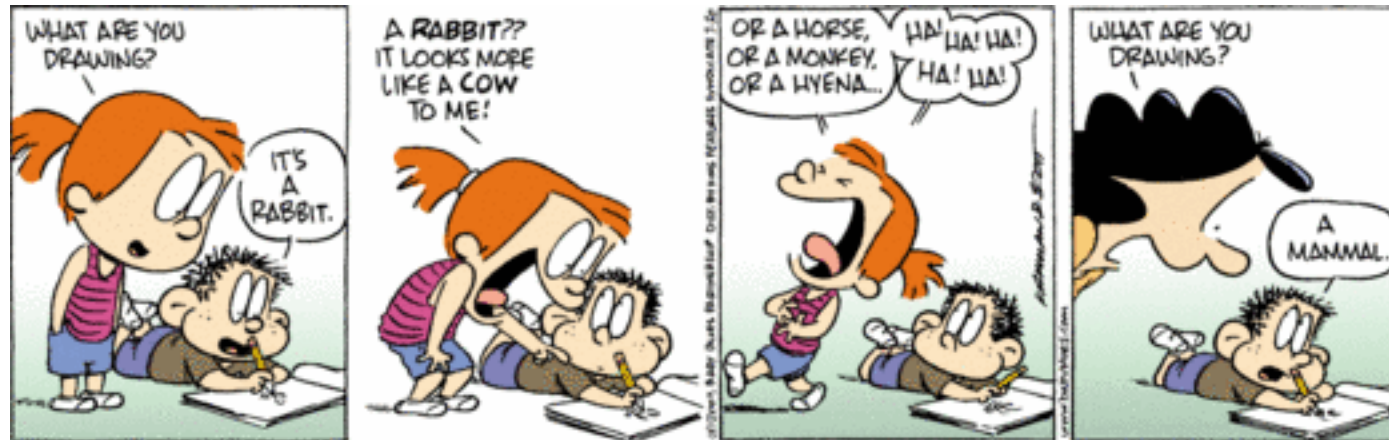
Text is Ambiguous



Source: Wikipedia

FUNDAMENTALS OF TOPIC MODELING

From Counting to Categorizing



<http://www.cse.buffalo.edu/~rapaport/575/categories.html>

Text Categorization

	Manual coding (bottom up)	Manual coding (top down)	Dictionaries	Supervised machine learning	Unsupervised machine learning
Assumptions					
Categories are predefined	No	Yes	Yes	Yes	No
Relevant text features are known	Yes	Yes	Yes	Yes	Yes
Mapping between text features and categories is known	No	No	Yes	No	No
Costs					
<i>Pre-analysis costs</i>					
Person-hours spent conceptualizing	Low	High	High	High	Low
Level of substantive knowledge	Low	High	High	High	Low
<i>Analysis costs</i>					
Person-hours spent per text	High	High	Low	Low	Low
Level of substantive knowledge	Moderate	Moderate	Low	Low	Low
<i>Post-analysis costs</i>					
Person hours spent interpreting	Moderate	Low	Low	Low	Moderate
Level of substantive knowledge	High	High	High	High	High

Source: Debortoli et al. (2016)

Text Categorization

	Manual coding (bottom up)	Manual coding (top down)	Dictionaries	Supervised machine learning	Unsupervised machine learning
Assumptions					
Categories are predefined	No	Yes	Yes	Yes	No
Relevant text features are known	Yes	Yes	Yes	Yes	Yes
Mapping between text features and categories is known	No	No	Yes	No	No
Costs					
<i>Pre-analysis costs</i>					
Person-hours spent conceptualizing	Low	High	High	High	Low
Level of substantive knowledge	Low	High	High	High	Low
<i>Analysis costs</i>					
Person-hours spent per text	High	High	Low	Low	Low
Level of substantive knowledge	Moderate	Moderate	Low	Low	Low
<i>Post-analysis costs</i>					
Person hours spent interpreting	Moderate	Low	Low	Low	Moderate
Level of substantive knowledge	High	High	High	High	High

Source: Debortoli et al. (2016)

Text Categorization

	Manual coding (bottom up)	Manual coding (top down)	Dictionaries	Supervised machine learning	Unsupervised machine learning
Assumptions					
Categories are predefined	No	Yes	Yes	Yes	No
Relevant text features are known	Yes	Yes	Yes	Yes	Yes
Mapping between text features and categories is known	No	No	Yes	No	No
Costs					
<i>Pre-analysis costs</i>					
Person-hours spent conceptualizing	Low	High	High	High	Low
Level of substantive knowledge	Low	High	High	High	Low
<i>Analysis costs</i>					
Person-hours spent per text	High	High	Low	Low	Low
Level of substantive knowledge	Moderate	Moderate	Low	Low	Low
<i>Post-analysis costs</i>					
Person hours spent interpreting	Moderate	Low	Low	Low	Moderate
Level of substantive knowledge	High	High	High	High	High

Source: Debortoli et al. (2016)

Text Categorization

	Manual coding (bottom up)	Manual coding (top down)	Dictionaries	Supervised machine learning	Unsupervised machine learning
Assumptions					
Categories are predefined	No	Yes	Yes	Yes	No
Relevant text features are known	Yes	Yes	Yes	Yes	Yes
Mapping between text features and categories is known	No	No	Yes	No	No
Costs					
<i>Pre-analysis costs</i>					
Person-hours spent conceptualizing	Low	High	High	High	Low
Level of substantive knowledge	Low	High	High	High	Low
<i>Analysis costs</i>					
Person-hours spent per text	High	High	Low	Low	Low
Level of substantive knowledge	Moderate	Moderate	Low	Low	Low
<i>Post-analysis costs</i>					
Person hours spent interpreting	Moderate	Low	Low	Low	Moderate
Level of substantive knowledge	High	High	High	High	High

Source: Debortoli et al. (2016)

Text Categorization

	Manual coding (bottom up)	Manual coding (top down)	Dictionaries	Supervised machine learning	Unsupervised machine learning
Assumptions					
Categories are predefined	No	Yes	Yes	Yes	No
Relevant text features are known	Yes	Yes	Yes	Yes	Yes
Mapping between text features and categories is known	No	No	Yes	No	No
Costs					
<i>Pre-analysis costs</i>					
Person-hours spent conceptualizing	Low	High	High	High	Low
Level of substantive knowledge	Low	High	High	High	Low
<i>Analysis costs</i>					
Person-hours spent per text	High	High	Low	Low	Low
Level of substantive knowledge	Moderate	Moderate	Low	Low	Low
<i>Post-analysis costs</i>					
Person hours spent interpreting	Moderate	Low	Low	Low	Moderate
Level of substantive knowledge	High	High	High	High	High

Source: Debortoli et al. (2016)

Text Categorization

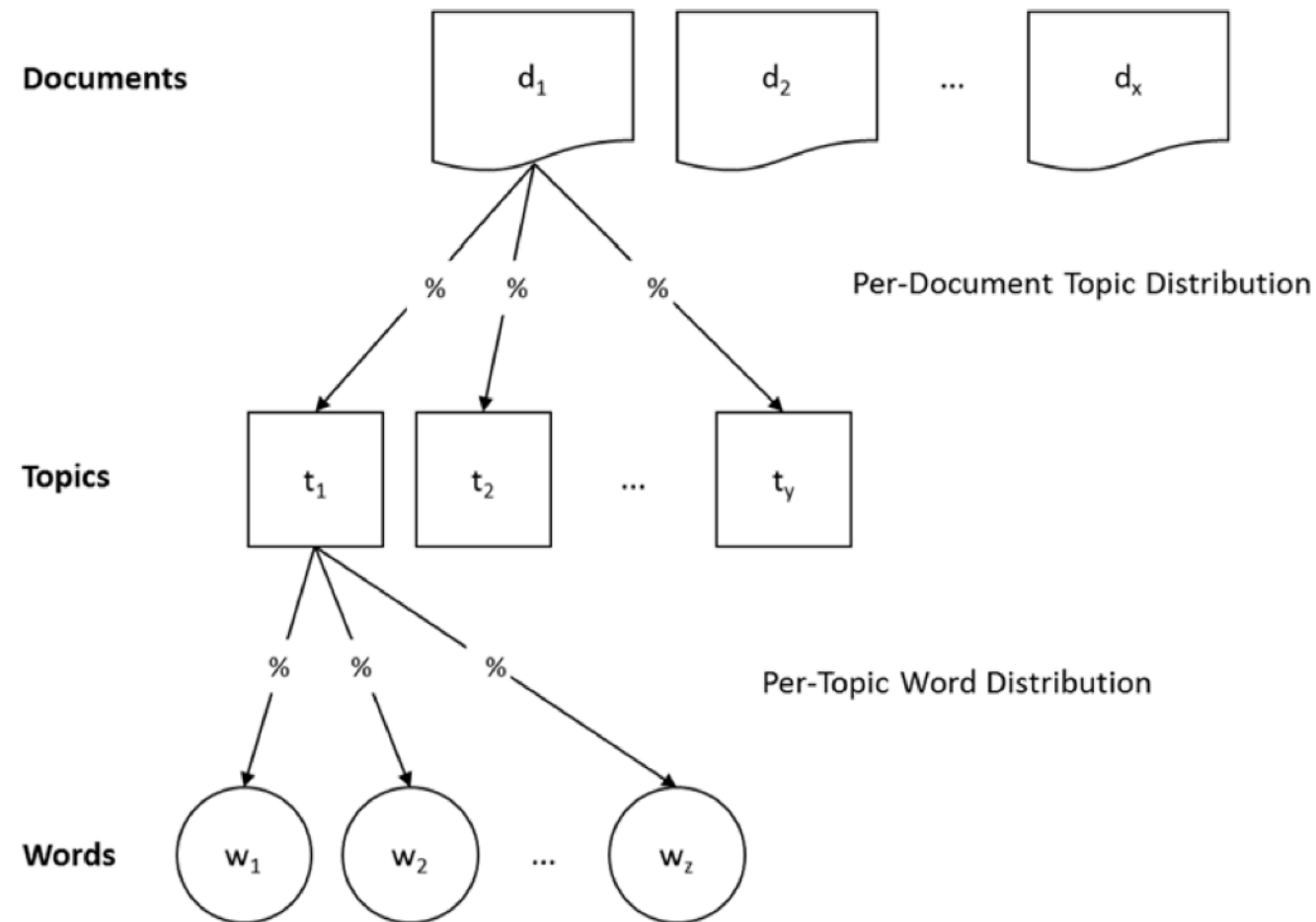
	Manual coding (bottom up)	Manual coding (top down)	Dictionaries	Supervised machine learning	Unsupervised machine learning
Assumptions					
Categories are predefined	No	Yes	Yes	Yes	No
Relevant text features are known	Yes	Yes	Yes	Yes	Yes
Mapping between text features and categories is known	No	No	Yes	No	No
Costs					
<i>Pre-analysis costs</i>					
Person-hours spent conceptualizing	Low	High	High	High	Low
Level of substantive knowledge	Low	High	High	High	Low
<i>Analysis costs</i>					
Person-hours spent per text	High	High	Low	Low	Low
Level of substantive knowledge	Moderate	Moderate	Low	Low	Low
<i>Post-analysis costs</i>					
Person hours spent interpreting	Moderate	Low	Low	Low	Moderate
Level of substantive knowledge	High	High	High	High	High

Source: Debortoli et al. (2016)

What are Topic Models?

- **Unsupervised** machine learning methods for text mining (e.g., Latent Semantic Analysis, Latent Dirichlet Allocation)
- Theoretical grounding: **Distributional hypothesis** of linguistics
 - Words that co-occur together in similar contexts (e.g., ball, goal, offside) tend to have similar meanings
 - Co-occurrence patterns can be interpreted as topics (e.g., football) and used to cluster documents

Schematic Overview of Probabilistic Topic Modeling with LDA



Source: Debortoli et al. (2016)

Illustrative Example of Probabilistic Topic Modeling with LDA

Exemplary Customer Review about a **Fitbit Flex**

I bought this for my 14 year old daughter as a gift. She received it in July. It works great - she lost 6 pounds in 2 weeks. The Fitbit makes staying in shape easy. The iPhone app works fine.



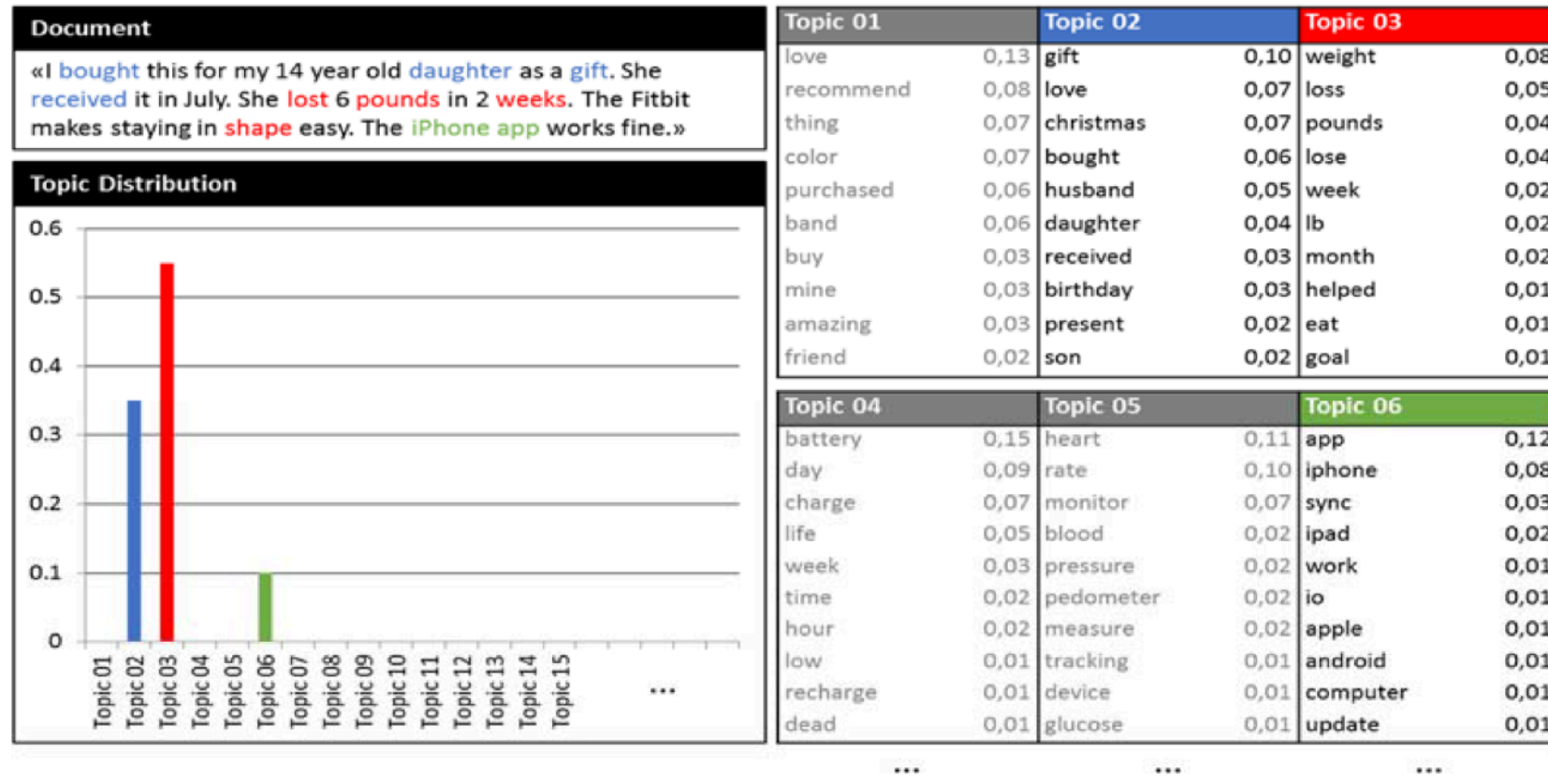
Topics

Birthday present

Loosing weight

Mobile app

Illustrative Example of Probabilistic Topic Modeling with LDA

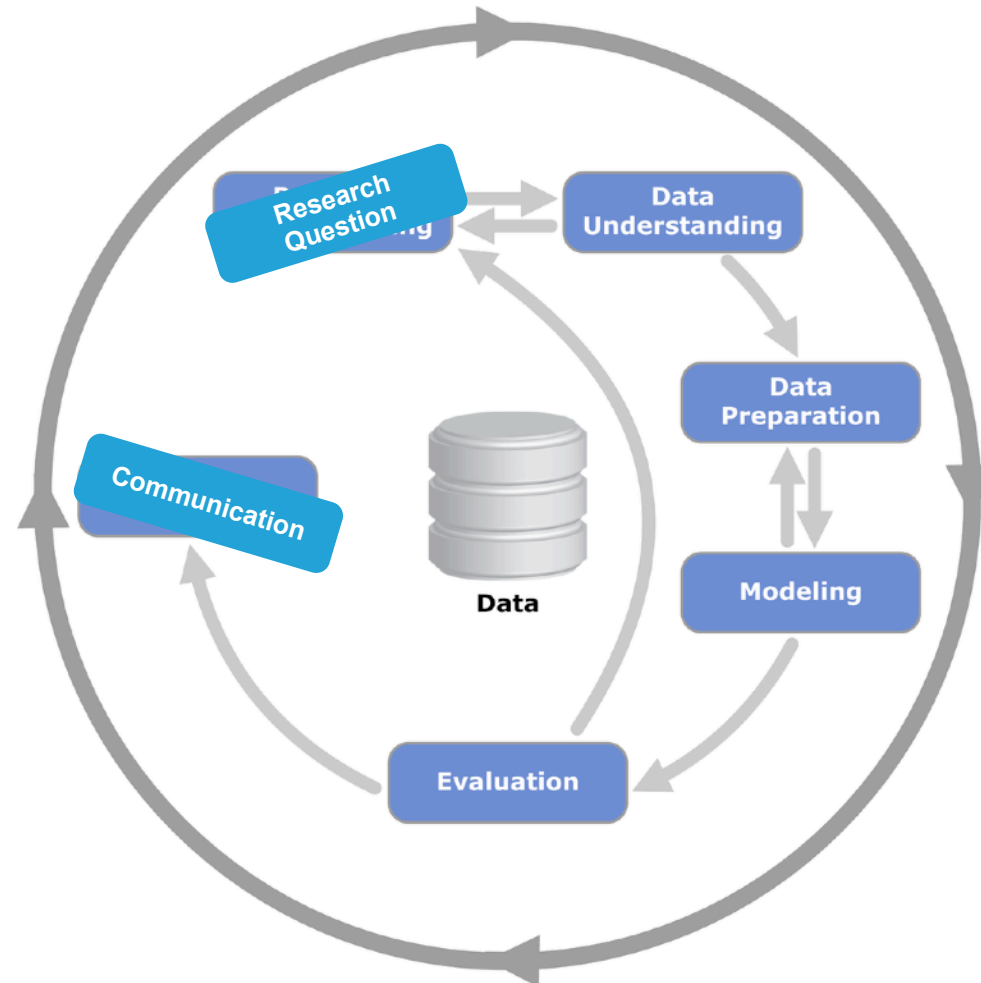


Source: Debortoli et al. (2016)

TOPIC MODELING WALKTHROUGH



Cross Industry Standard Process for Data Mining



Source: Shearer et al. (2000)

Research Question

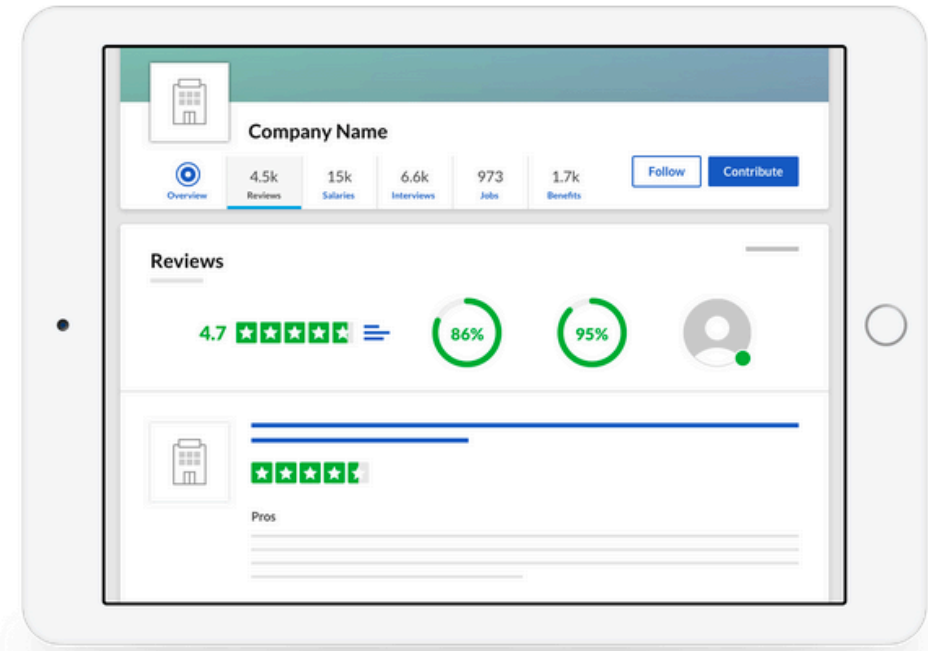
“What factors drive employees’ company ratings?”



Company reviews and ratings. Get the whole story.

Search ratings and reviews of over 600,000 companies worldwide. Get the inside scoop and find out what it's really like from people who've actually worked there.

Write a Review

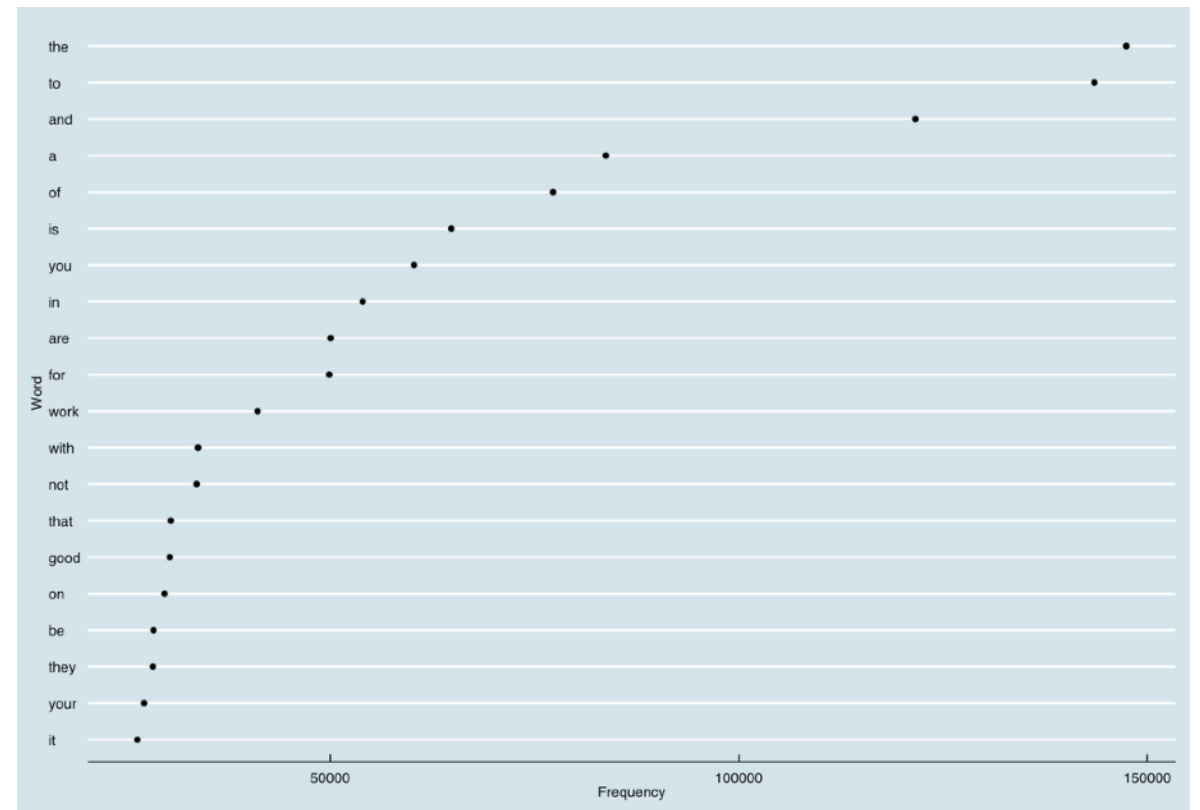


Data Understanding

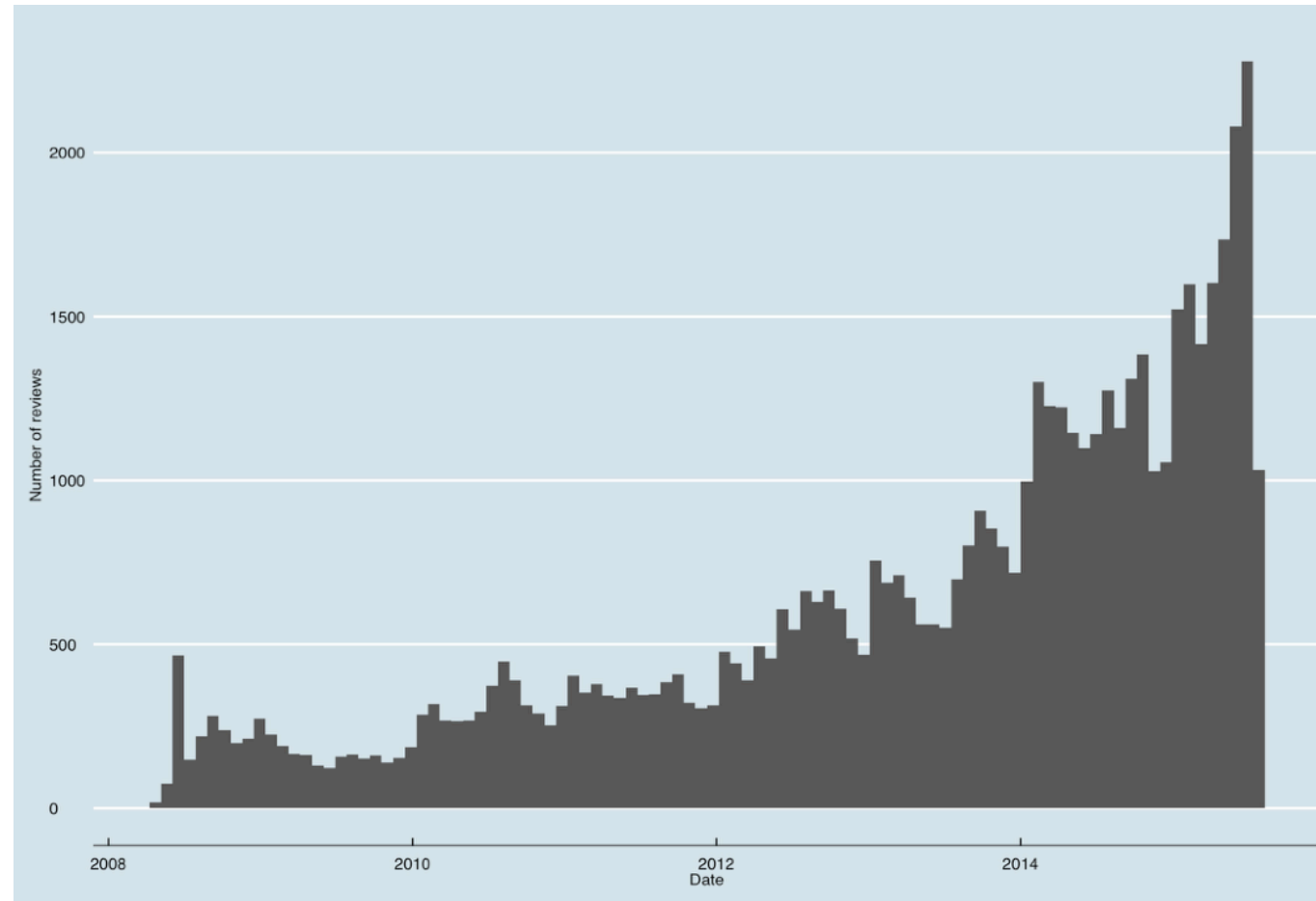
```
{  
  "shortName": "Bank of America",  
  "date": "2010-10-03",  
  "stars": 3,  
  "text": "Great benefits for associates, Paid maternity/paternity leave, most associates receive 3  
    weeks of vacation leave per year (SSS, PB, AM and four weeks for BCM). Micro-management,  
    poor leadership, lack of recognition, extremely under staffed. Do not forget the human aspect.  
    Micro-management is not the answer to every situation. Put more people in the branches."  
}
```

Data Understanding

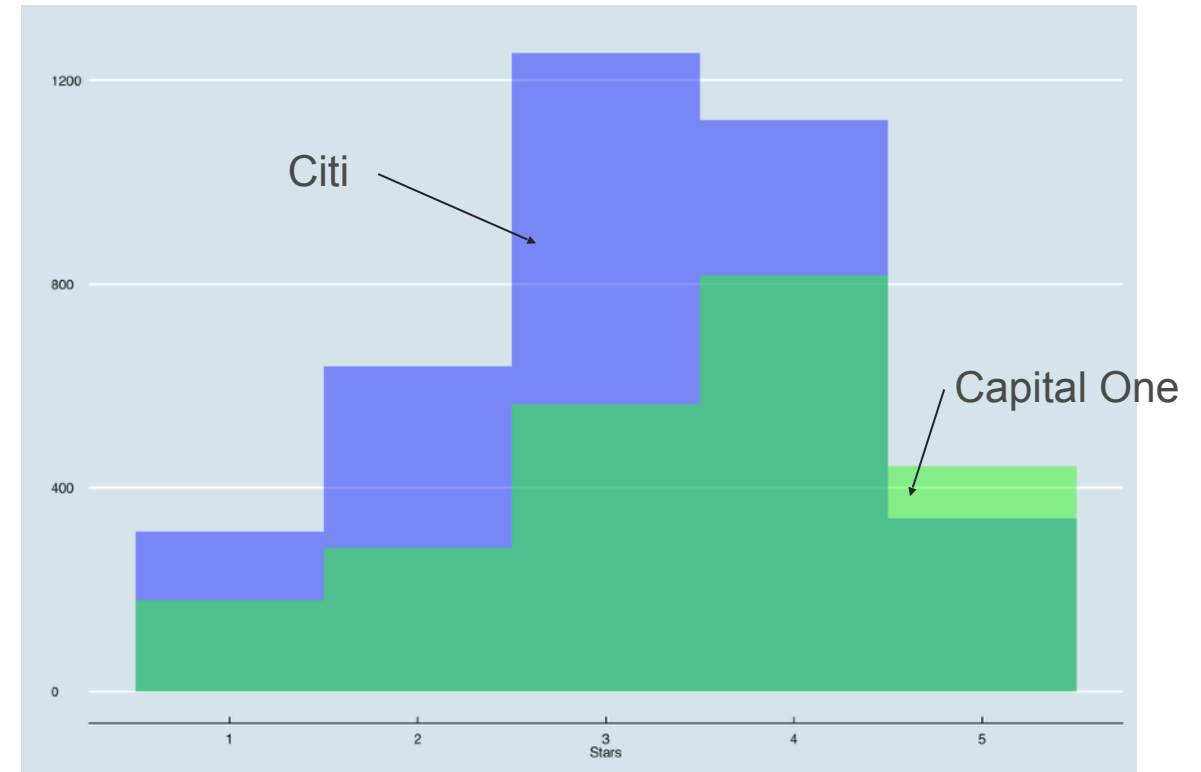
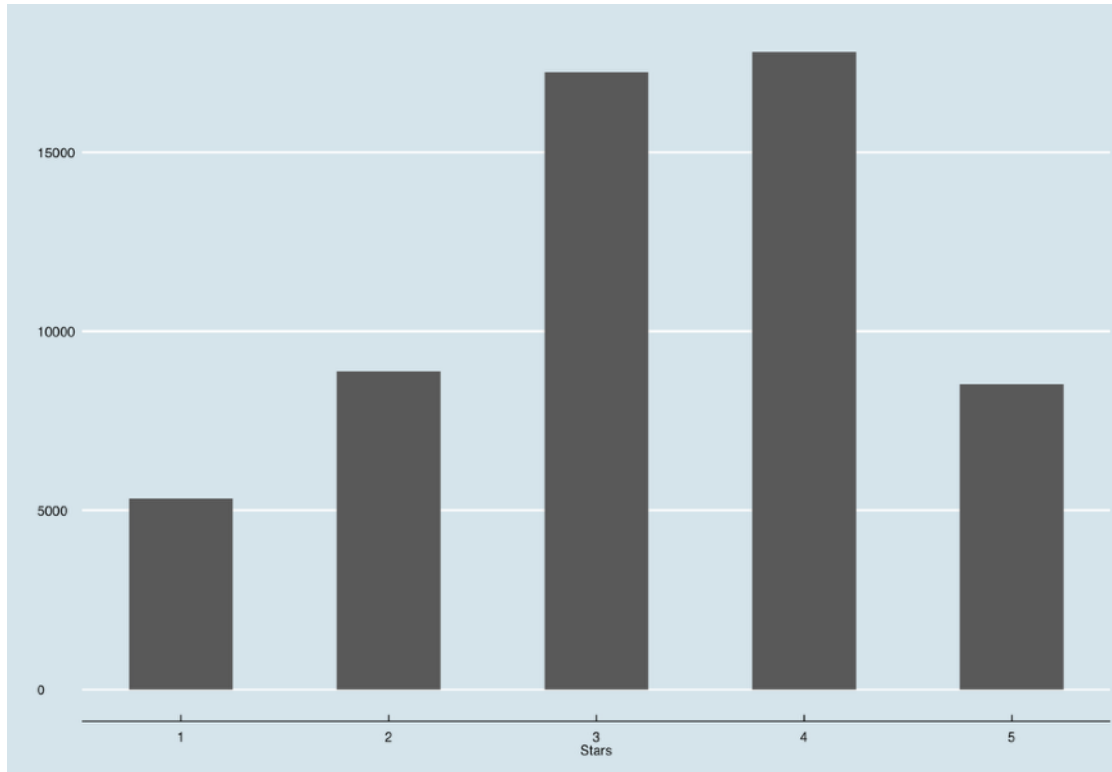
- Sub-sample
 - Finance industry only
- Number of documents
 - 57,765
- Number of words (tokens)
 - 1,608,259
- Number of unique words
 - 1,740



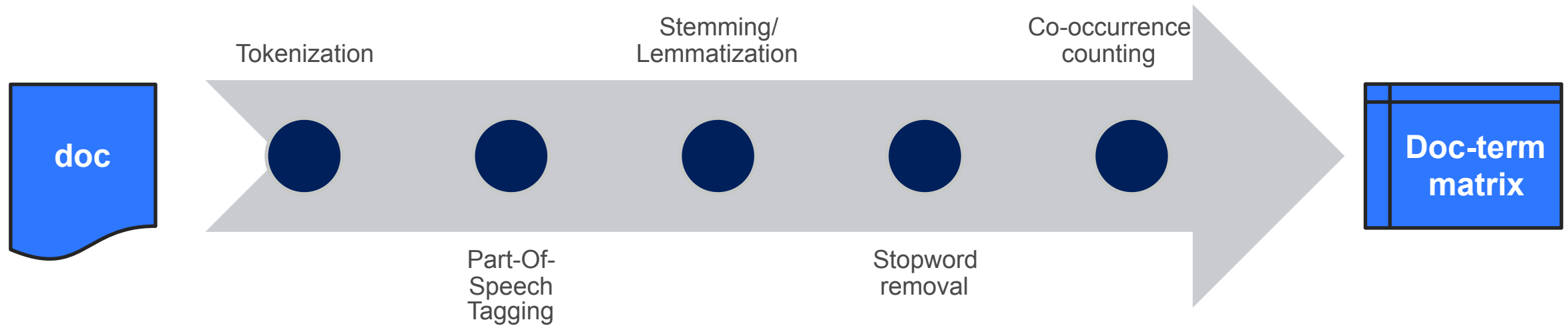
Data Understanding



Data Understanding

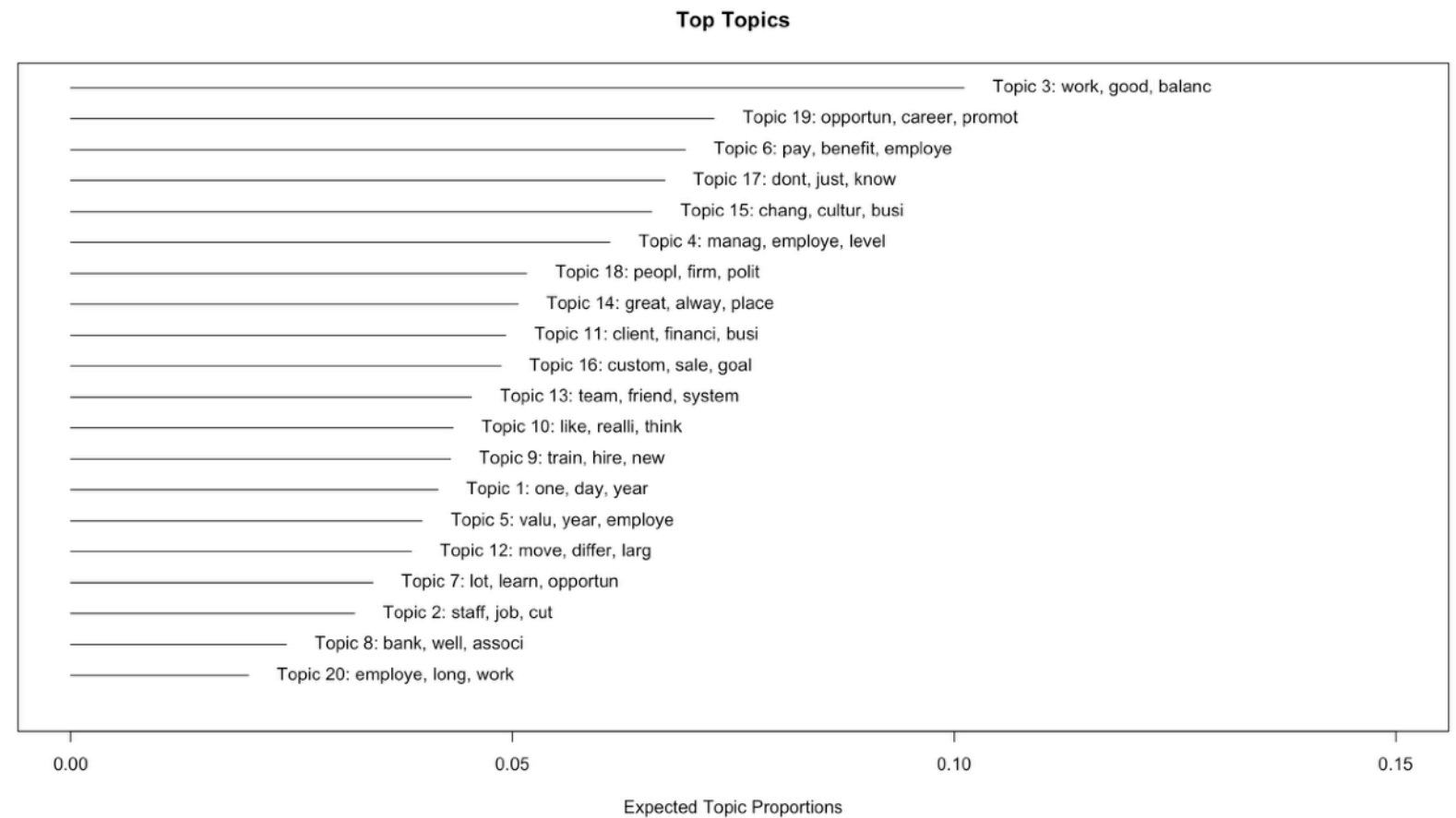


Data Preparation



Modeling and Evaluation – Iteration #1

- Estimate 30 most prevalent topics
 - Takes approx. 10 minutes on a MacBook Pro
 - Wait for the “aha” effect



Modeling and Evaluation – Iteration #1

- Estimate 30 most prevalent topics
 - Takes approx. 10 minutes on a MacBook Pro
 - Wait for the “aha” effect

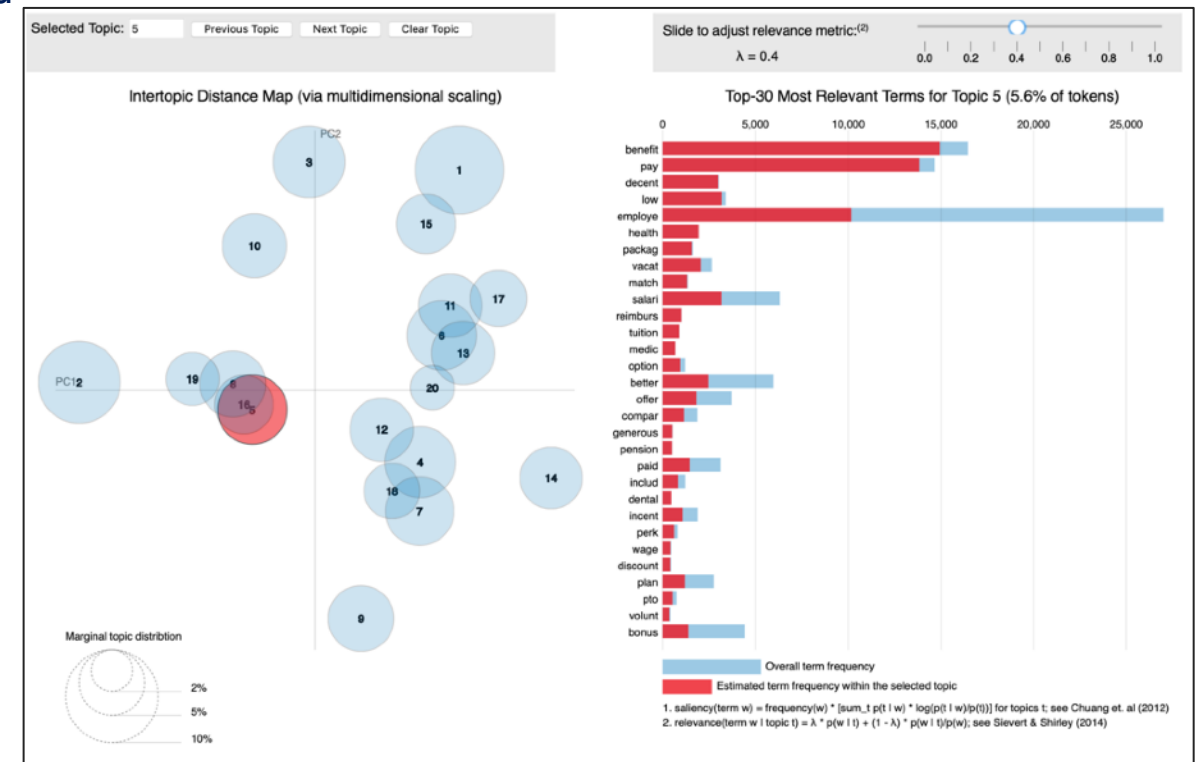
Top-3 Documents for Topic 6

Health insurance, Vacation, sick pay, paid maternity leave 12 weeks. Every year perks decrease and are eliminated. Uneducated people with their nose in the air.
Decent benefits, decent bonuses, decent vacation/off time. Inadequate pay and inadequate coverage.
Huge Annual Bonus amount, Paid Overtime amount can be earned, 1 Time free meal & free transport facility. Less On paper CTC offered. Should include the approximate Annual Bonus & Gratuity in the offered CTC on-paper.

Modeling and Evaluation – Iteration #2

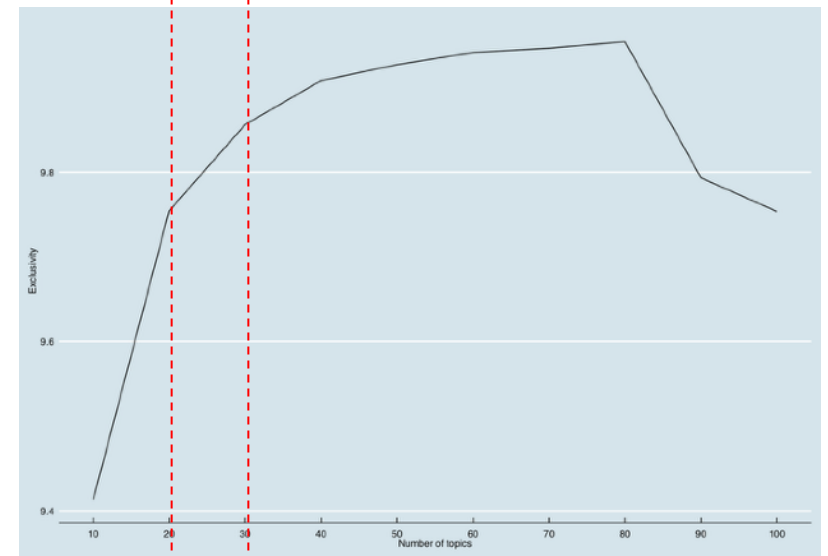
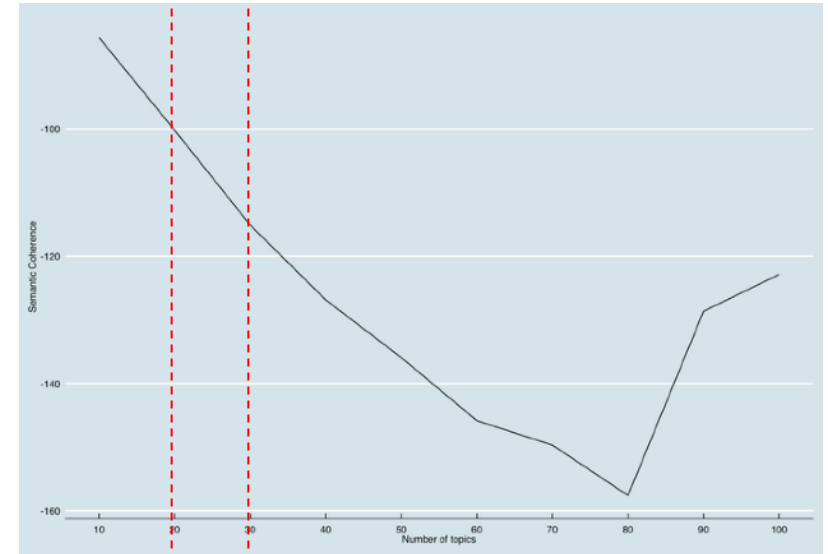
- Find the right number of topics
 - Manual** investigation
 - Are topics coherent? No duplicate topics? No fused topics?
 - Increase or reduce number of topics

LDAvis: A R package for interactive topic model visualization.



Modeling and Evaluation – Iteration #2

- Find the right number of topics
 - **Automated** search
 - e.g.: From 10 to 100 topics, in steps of 10
 - Takes several hours on a MacBook Pro
 - Evaluate models with regards to **Semantic Coherence** and **Exclusivity**



Modeling and Evaluation – Iteration #3

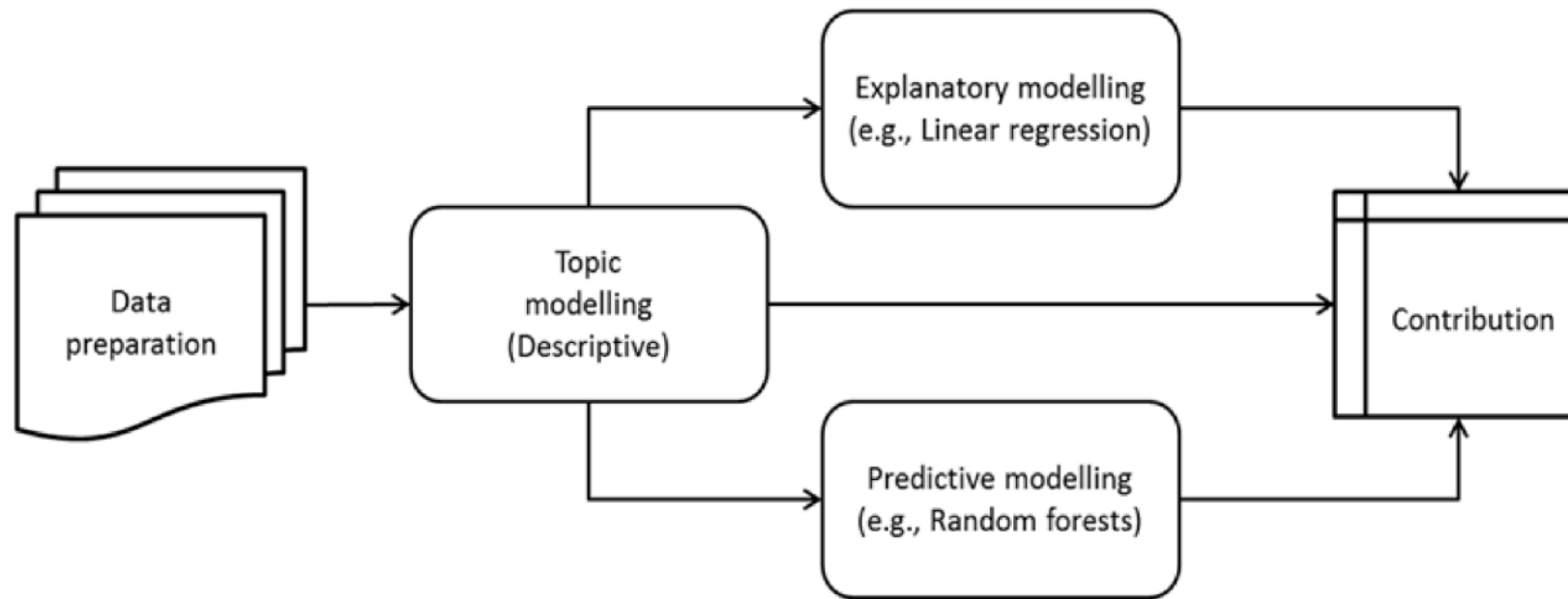
- Final experimental evaluation through human coders
 - Word intrusion task
 - Topic intrusion task

<http://etc.ch/bRsu>

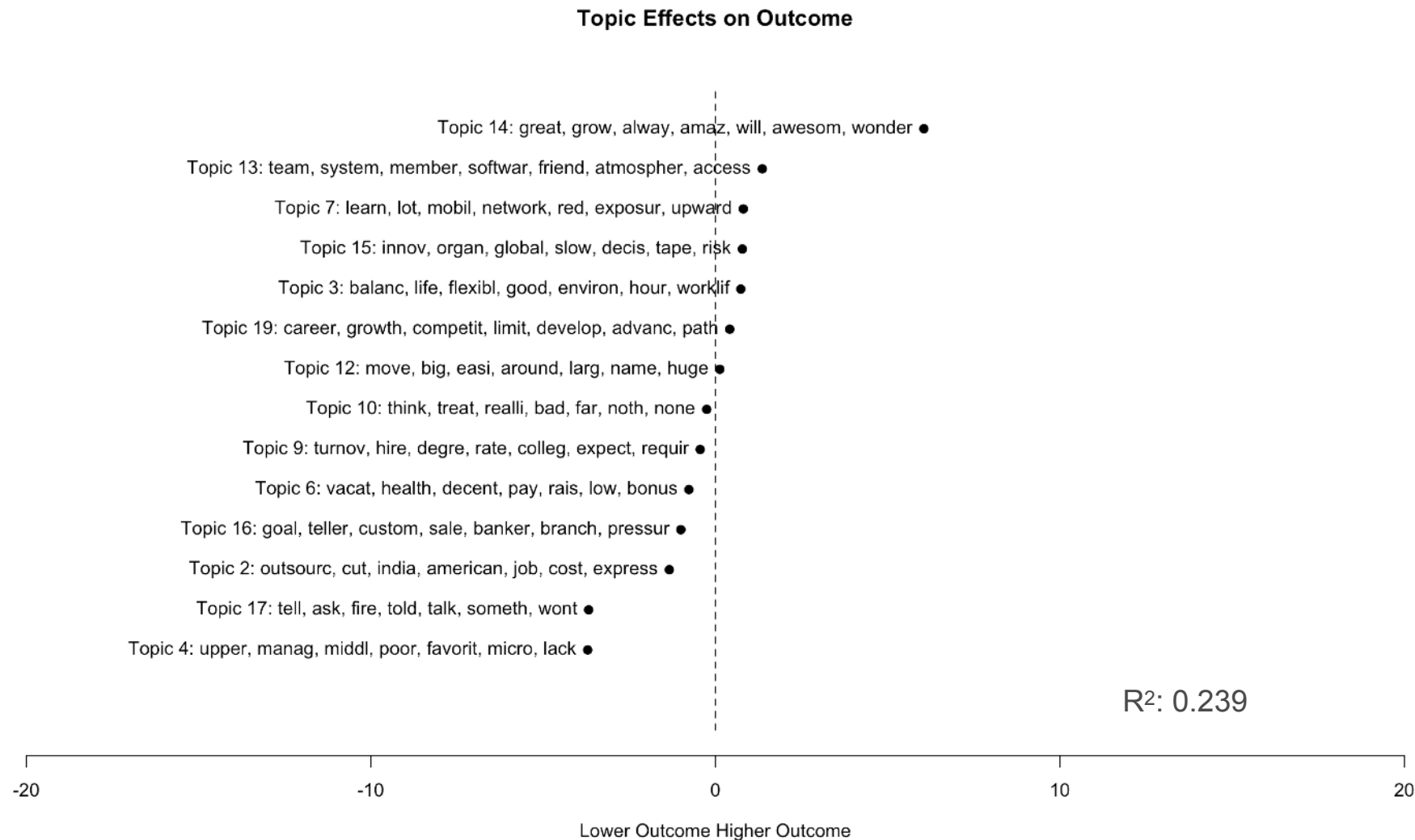


Modeling and Evaluation – Iteration #4

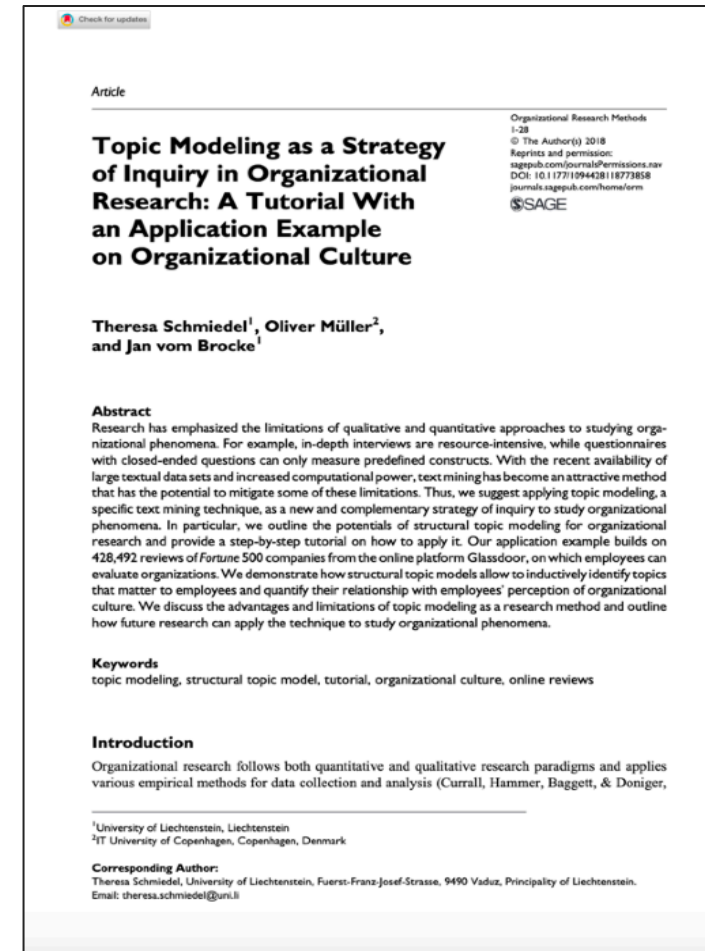
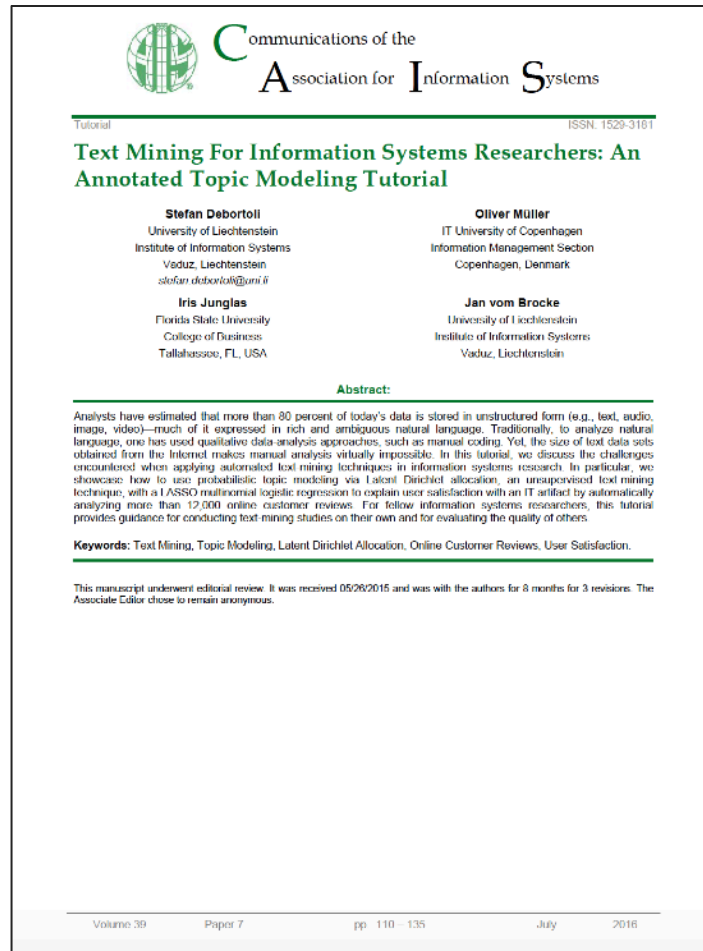
- Modeling the relationship between topics and stars



Modeling and Evaluation – Iteration #4

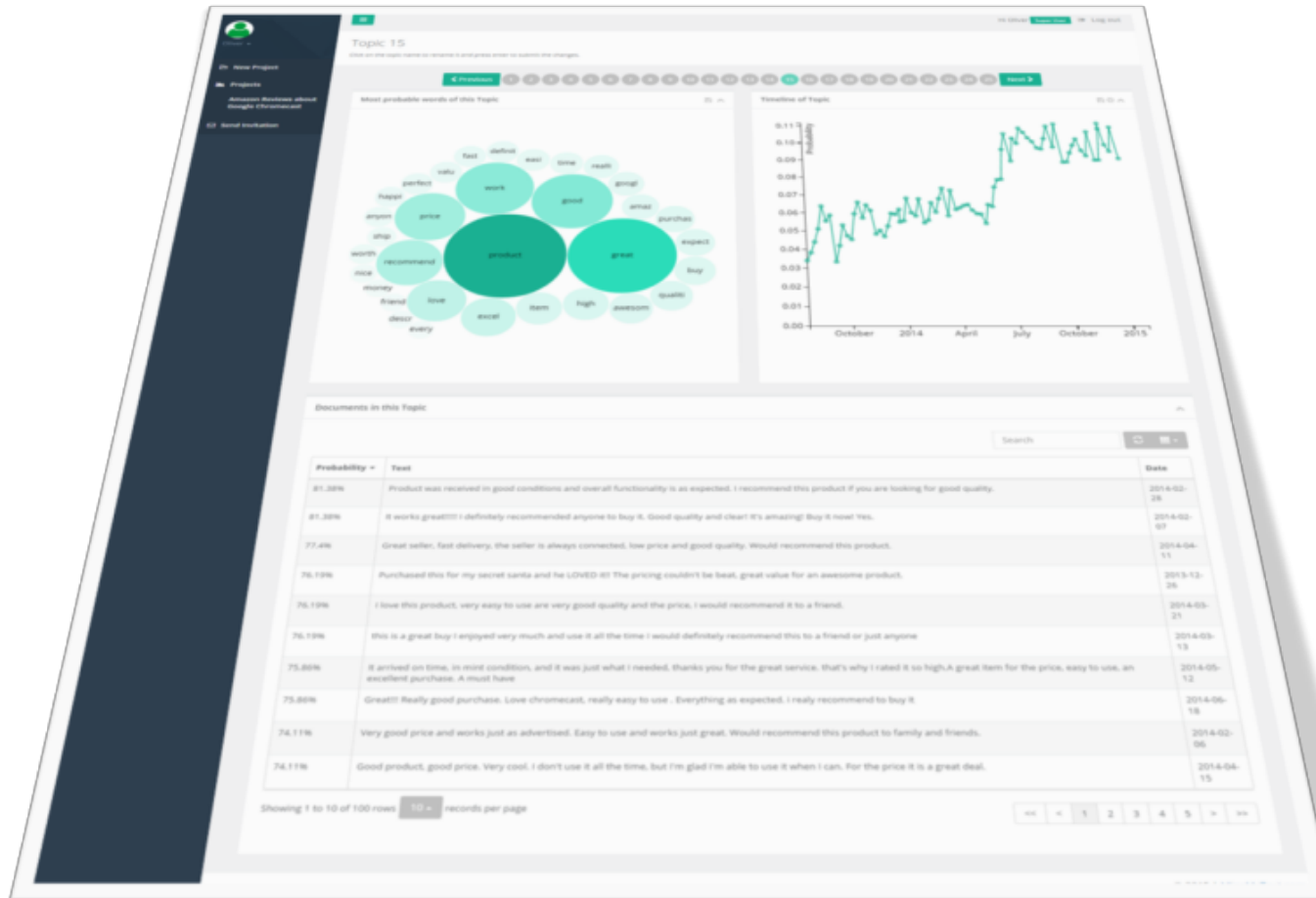


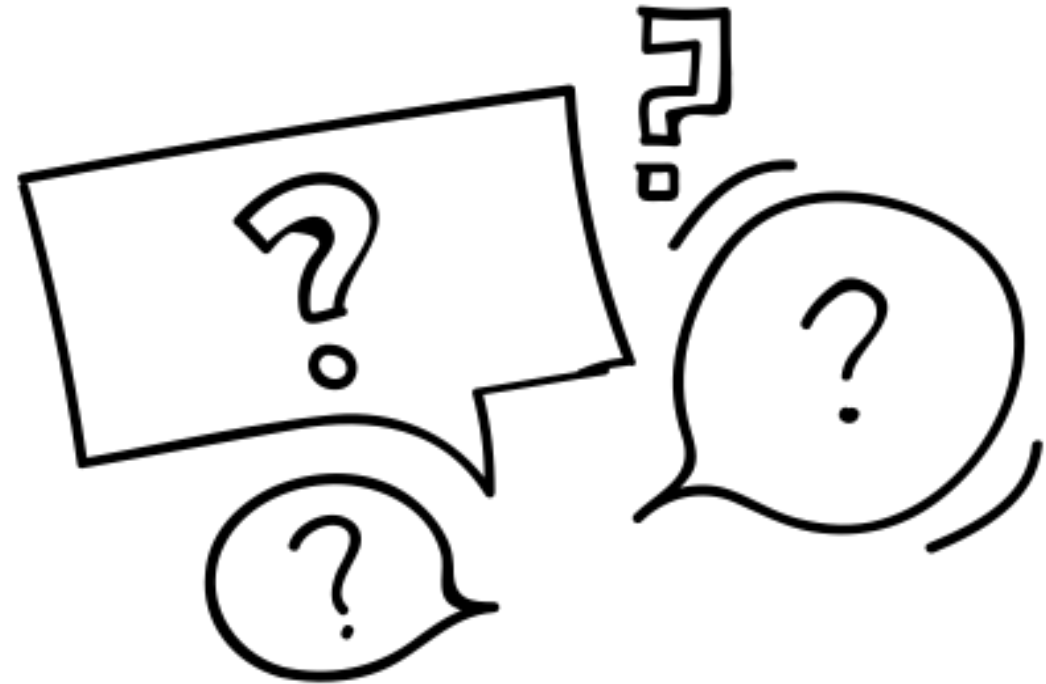
Communication



TOPIC MODELING WALKTHROUGH

www.MineMyText.com





Prof. Dr. Oliver Müller

Lehrstuhl für Wirtschaftsinformatik, insb. Data Analytics
Universität Paderborn

Warburger Str. 100, 33098 Paderborn

R: Q2.457

E: oliver.mueller@uni-paderborn.de

T: +49-5251-605245

W: <https://wiwi.uni-paderborn.de/dep3/mueller/>

- Cox M, Ellsworth D (1997) Application-controlled demand paging for out-of-core visualization. Proceedings of the 8th Conference on Visualization, 235–244.
- Laney, D. (2001). 3D data management: Controlling data volume, velocity and variety. META Group Research Note, 6(70).
- IBM. (2012). Analytics: The real-world use of big data. Retrieved from <http://www-935.ibm.com/services/us/gbs/thoughtleadership/ibv-big-data-at-work.html>
- Dhar, V. (2013). Data Science and Prediction. Communications of the ACM, 56 (12), 64–73.
- Zur Mühlen, M., & Shapiro, R. (2010). Business process analytics. In J. Vom Brocke & M. Rosemann (Eds.), Handbook on Business Process Management (Vol. 2). Springer.
- Mitchell, T. M. (1997). Machine learning. McGraw Hill.
- Shearer C. (2000). The CRISP-DM model: the new blueprint for data mining, Journal of Data Warehousing, 5(4), 13-22.
- Debortoli, S., Müller, O., Junglas, I. A., & vom Brocke, J. (2016). Text mining for information systems researchers: an annotated topic modeling tutorial. Communication of the AIS, 39, 7.
- Schmiedel, T., Müller, O., & vom Brocke, J. (2018). Topic modeling as a strategy of inquiry in organizational research: A tutorial with an application example on organizational culture. Organizational Research Methods, 1094428118773858.
- Icons: questions by Depb Dew, Jemis Mali from the Noun Project