

# Social Topic Distributions

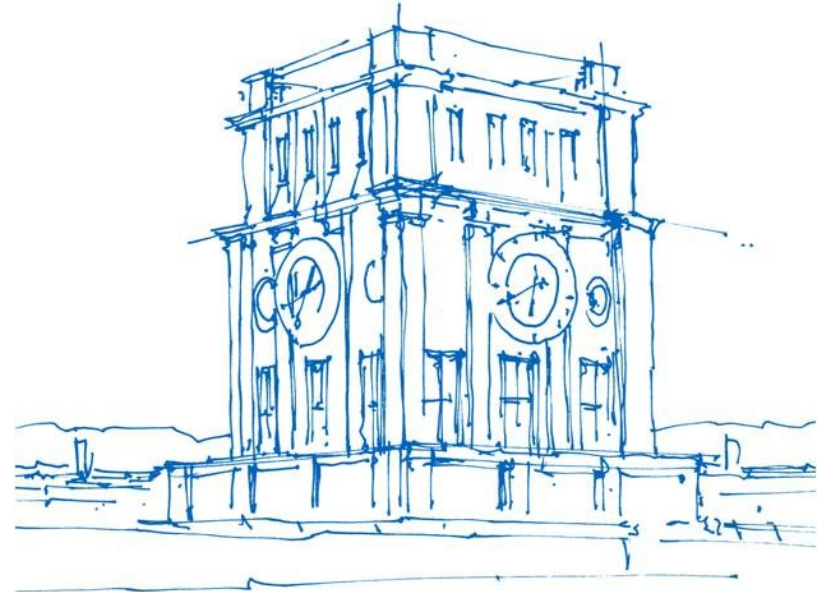
Mürüvvet Hasanbaşoğlu

Hakan Akyürek

Technical University of Munich

Faculty of Informatics

Munich, 20. July 2020



*Uhrenturm der TUM*

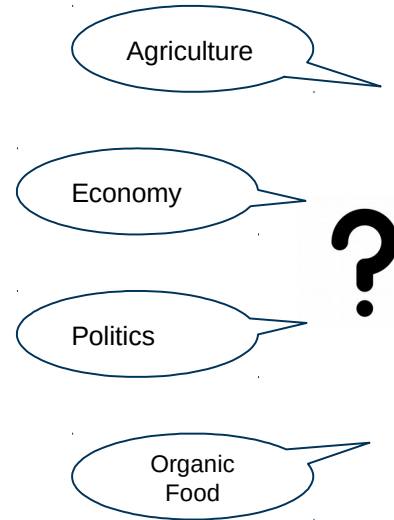
# Outline

- Motivation
- Data
- Experiments
- Results

# Outline

- **Motivation**
- Data
- Experiments
- Results

# Motivation



## Amazon to buy Whole Foods Market in deal valued at \$13.7 billion

Amazon.com jolted the grocery industry Friday when it announced plans to buy Whole Foods Market, introducing new uncertainty to a sector already struggling to keep up with growing competition. The \$13.7 billion deal heightens a years-long battle between Amazon, the Internet darling, and powerhouse merchants such as Walmart, which recently beefed up its online operations with a \$3.3 billion purchase of an Amazon competitor. Now Seattle-based Amazon - which for years has been testing grocery innovations in quiet corners - could lay claim to a fleet of more than 460 stores throughout the United States, Canada and Britain.....

---

**Comments:**

Guy above shows that this is not the scenario I thought it was. Still annoying though - Yeah, I definitely missed that development. I wonder if that will eventually be deemed illegal...

I enjoy shopping with Amazon. I am a Prime member and their shipping and customer service is However, I am considering not renewing next year as I have split some of my shopping over to Jet.com (Which Walmart acquired last year)....



There are I mean literally very few. But the dishes are priced, is one of the leading Organic caterers in Bangalore. With over a 100 organic dishes,....



You can find many online grocery stores like Dalbasket which is well-established and renowned for providing the different types of pulses...

# Outline

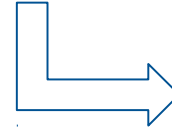
- Motivation
- **Data**
- Experiments
- Results

# Data

Curated Organic Dataset		# of articles	# of users	
Biased		Facebook	5,013	4,705
		Food Babe	15	15
		Food Revolution	78	60
		Organic Authority	66	0
		Organic Consumers	64	0
Unbiased	Forum	Cafe Mom	86	85
		Disqus	36	36
		Quora	567	523
		Reddit	81	78
		US Message Board	0	0
	News sites	Chicago Tribune	2,283	78
		Huffington Post	880	0
		LA Times	1,522	77
		NY Post	106	0
		NY Times	438	137
		USA Today	95	22
		Washington Post	1,563	943

12,893

164,387



1. Tokenization: `r"[a-zA-Z0-9]+|\.|\\?|\\!"`
2. Lowercase
3. Stop word removal
4. Rare word removal

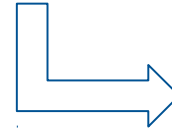
Tokens	Vocabulary
~19M	88,119

# Data

Curated Organic Dataset		# of articles	# of users	
Biased		Facebook	5,013	4,705
		Food Babe	15	15
		Food Revolution	78	60
		Organic Authority	66	0
		Organic Consumers	64	0
Unbiased	Forum	Cafe Mom	86	85
		Disqus	36	36
		Quora	567	523
		Reddit	81	78
		US Message Board	0	0
	News sites	Chicago Tribune	2,283	78
		Huffington Post	880	0
		LA Times	1,522	77
		NY Post	106	0
		NY Times	438	137
		USA Today	95	22
		Washington Post	1,563	943

12,893

164,387



1. Tokenization: `r"[a-zA-Z0-9]+|\.|\\.|\\?|\\!|\\|"`
2. Lowercase
3. Stop word removal
4. Rare word removal

Tokens	Vocabulary
~19M	88,119

Word embeddings → GloVe



Sentence embeddings → Universal Sentence

Encoder TensorFlow Hub



# Data

Curated Organic Dataset		# of articles	# of users	
Biased		Facebook	5,013	4,705
		Food Babe	15	15
		Food Revolution	78	60
		Organic Authority	66	0
		Organic Consumers	64	0
Unbiased	Forum	Cafe Mom	86	85
		Disqus	36	36
		Quora	567	523
		Reddit	81	78
		US Message Board	0	0
	News sites	Chicago Tribune	2,283	78
		Huffington Post	880	0
		LA Times	1,522	77
		NY Post	106	0
		NY Times	438	137
		USA Today	95	22
		Washington Post	1,563	943



# Data

Curated Organic Dataset		# of articles	# of users		
Biased		Facebook	5,013	4,705	→ 150 Articles
		Food Babe	15	15	
		Food Revolution	78	60	
		Organic Authority	66	0	
		Organic Consumers	64	0	
Unbiased	Forum	Cafe Mom	86	85	→ 150 Articles
		Disqus	36	36	
		Quora	567	523	
		Reddit	81	78	
		US Message Board	0	0	
	News sites	Chicago Tribune	2,283	78	
		Huffington Post	880	0	
		LA Times	1,522	77	
		NY Post	106	0	
		NY Times	438	137	
		USA Today	95	22	
		Washington Post	1,563	943	

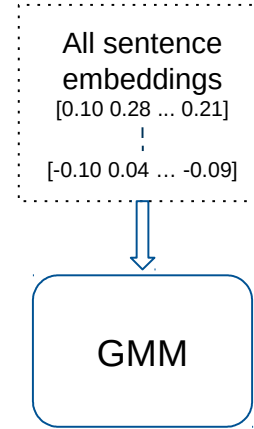
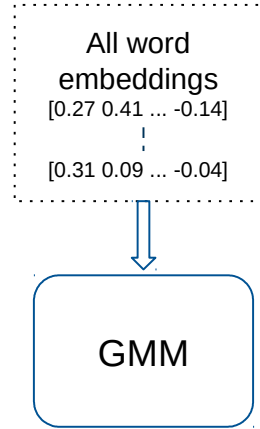
User datasets:

- Answering Users
- Random Users:
  1. General
  2. Biased
  3. Unbiased
  4. Forum
  5. News sites

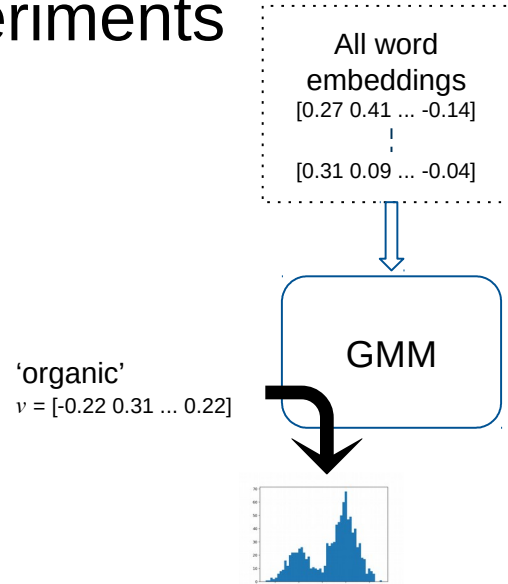
# Outline

- Motivation
- Data
- **Experiments**
- Results

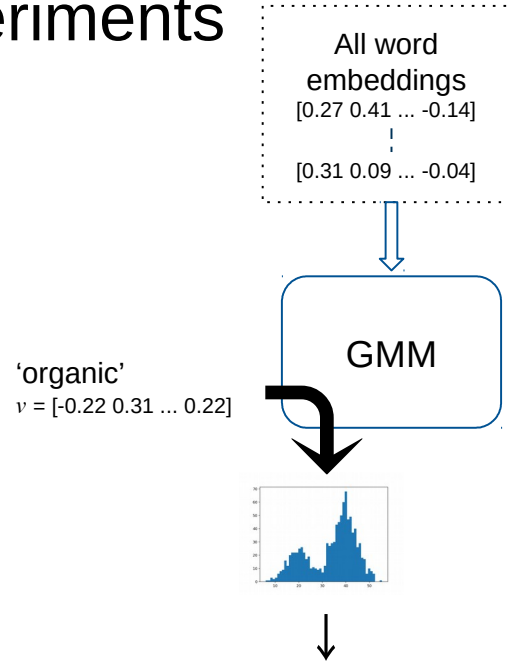
# Experiments



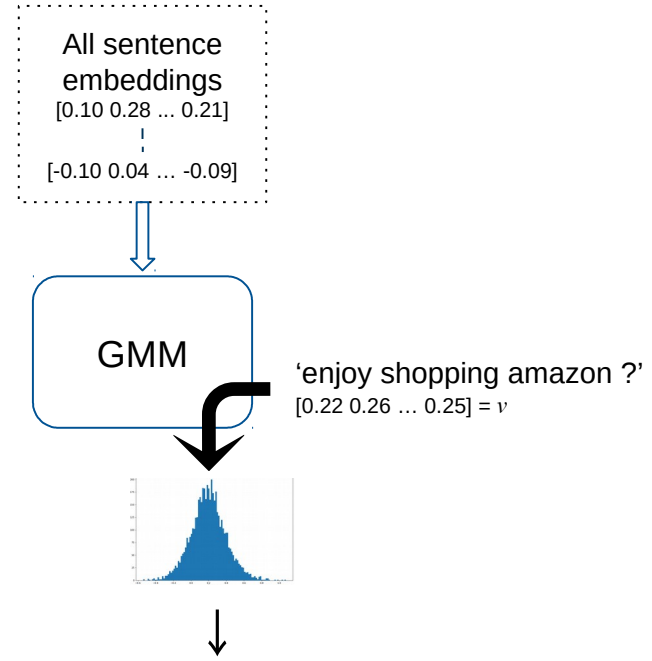
# Experiments



# Experiments



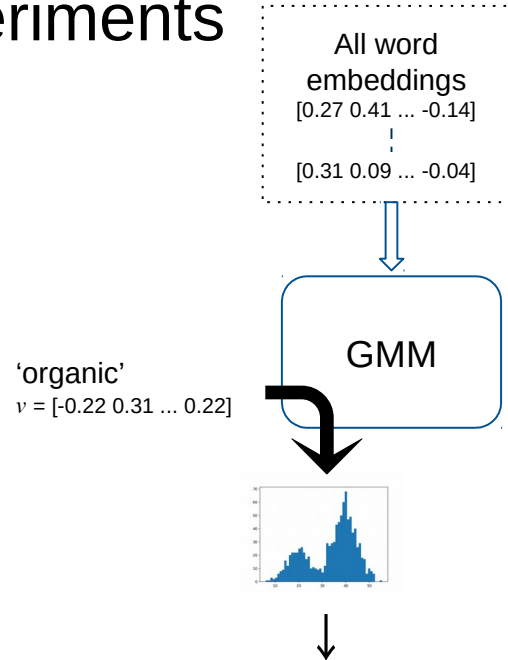
$$P(d_l) = \frac{CF * IDF(c_i, d_l, D)}{\sum_{i=1}^n \sum_{j=1}^k v_{ij}} = \frac{\sum_{j=1}^k v_{ij} * \log \frac{|D|}{|d \in D; c_i \in d|}}{\sum_{i=1}^n \sum_{j=1}^k v_{ij}}$$



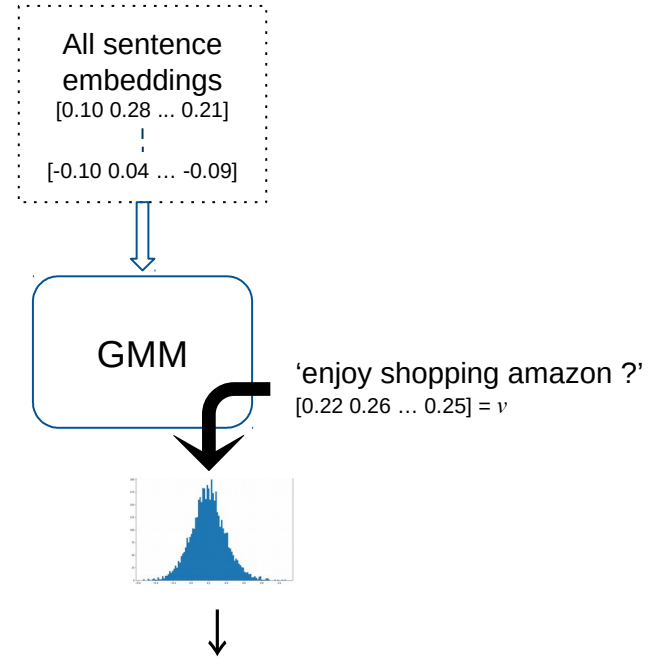
$$P(d_l) = \frac{\sum_{j=1}^k v_{ij}}{\sum_{i=1}^n \sum_{j=1}^k v_{ij}}$$

$l$  : # of documents  
 $k$  : # of words / sentences in a document  
 $n$  : # of topic clusters

# Experiments



$$P(d_l) = \frac{CF * IDF(c_i, d_l, D)}{\sum_{i=1}^n \sum_{j=1}^k v_{ij}} = \frac{\sum_{j=1}^k v_{ij} * \log \frac{|D|}{|d \in D; c_i \in d|}}{\sum_{i=1}^n \sum_{j=1}^k v_{ij}}$$



$$P(d_l) = \frac{\sum_{j=1}^k v_{ij}}{\sum_{i=1}^n \sum_{j=1}^k v_{ij}}$$

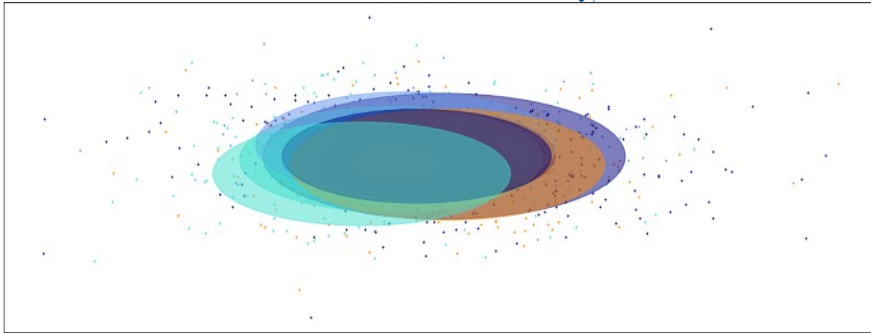
$l$  : # of documents  
 $k$  : # of words / sentences in a document  
 $n$  : # of topic clusters

?

# Experiments

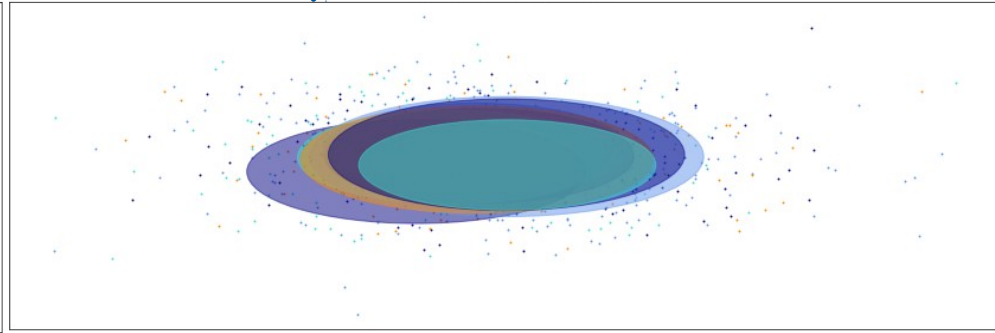
All word  
embeddings  
[0.27 0.41 ... -0.14]  
⋮  
[0.31 0.09 ... -0.04]

10 cluster sample diag GMM



All sentence  
embeddings  
[0.10 0.28 ... 0.21]  
⋮  
[-0.10 0.04 ... -0.09]

10 cluster sample diag GMM

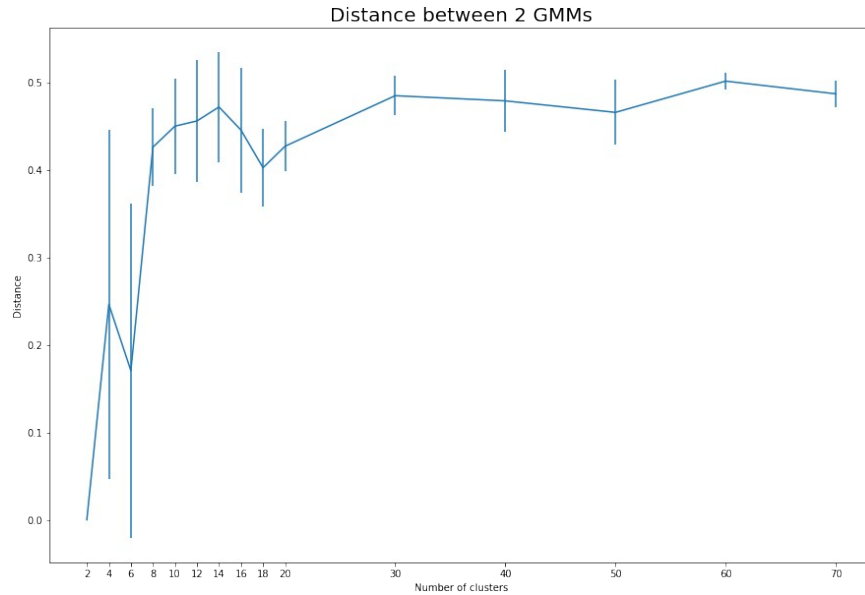


$l$  : # of documents  
 $k$  : # of words / sentences in a document  
 $n$  : # of topic clusters

?

# Experiments

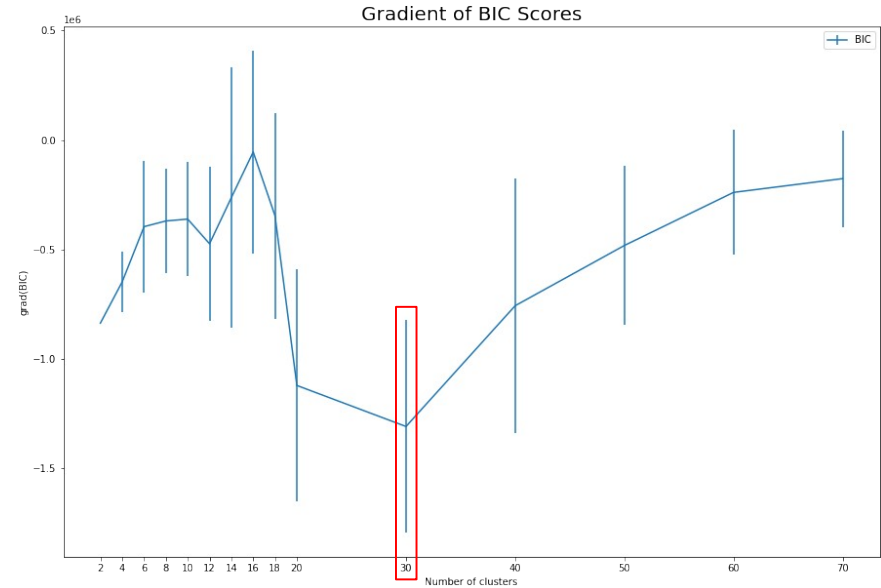
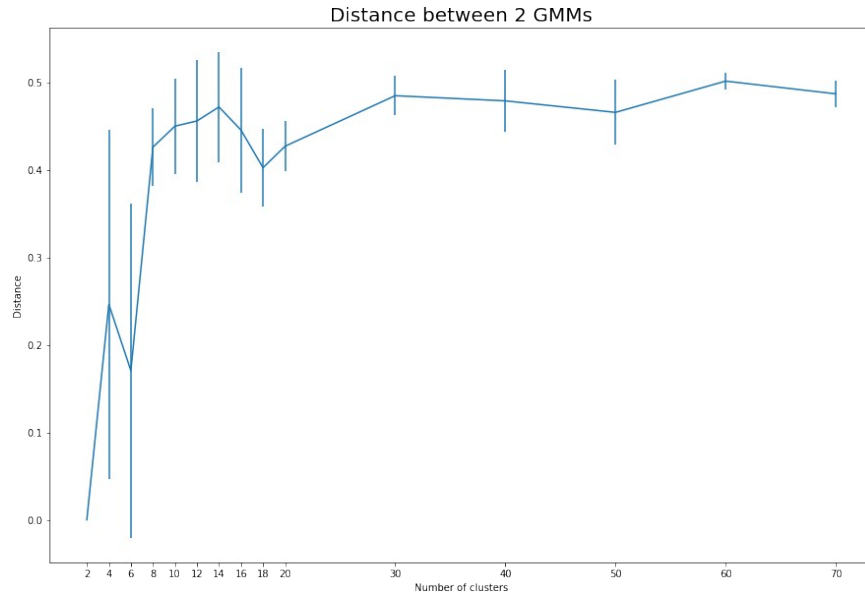
Optimal number of clusters for word embeddings:





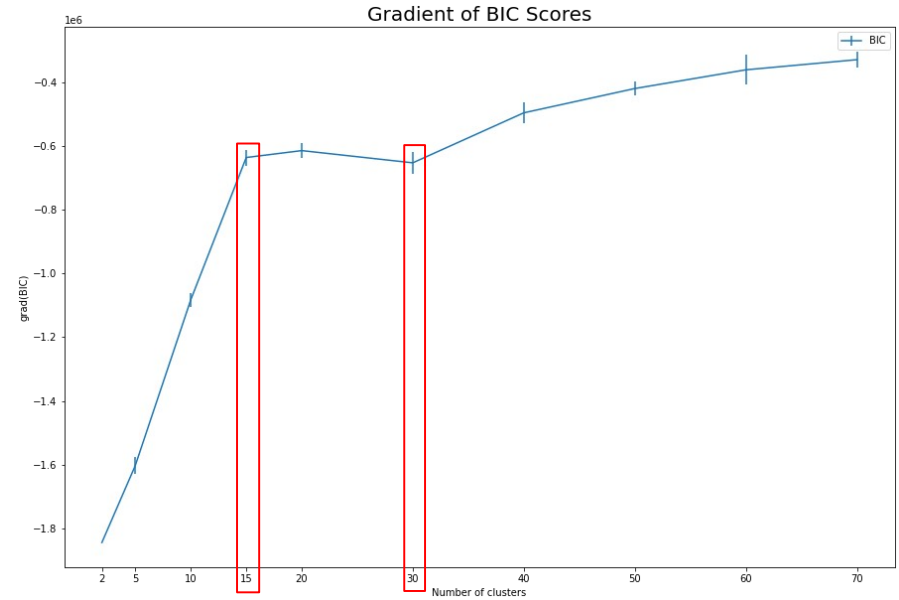
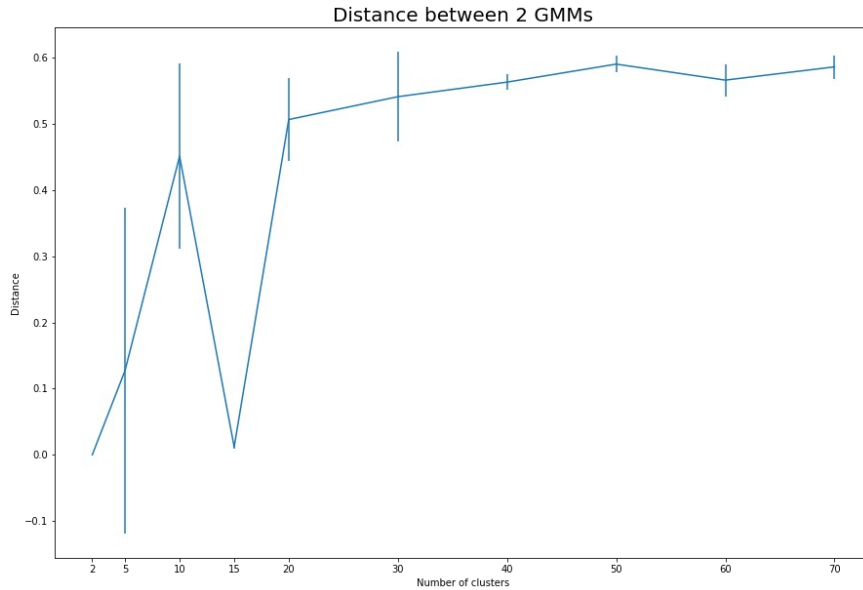
# Experiments

Optimal number of clusters for word embeddings:



# Experiments

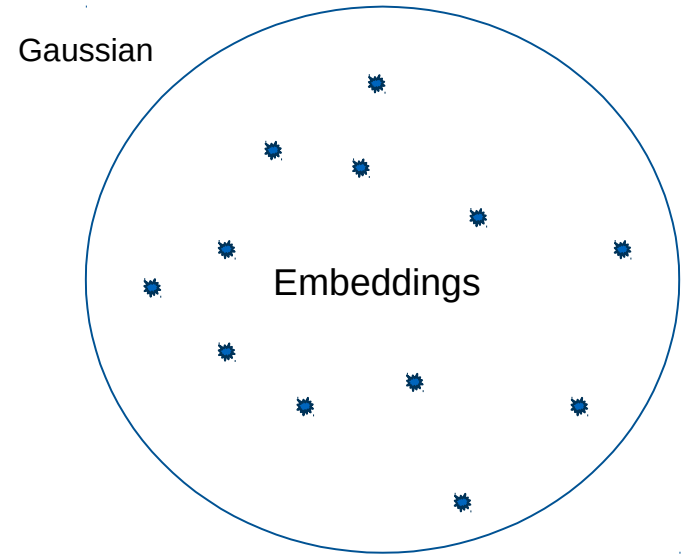
Optimal number of clusters for sentence embeddings:



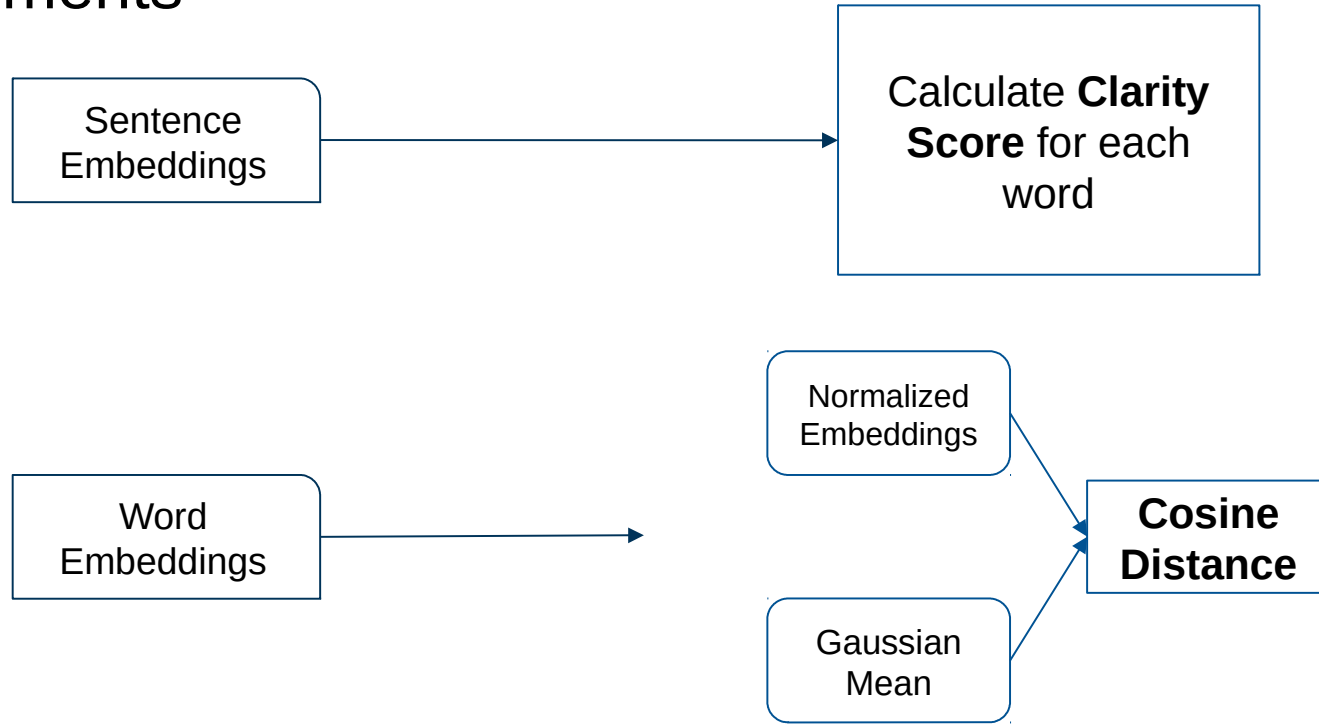
# Experiments

How do we label Gaussians in a GMM model?

- Find the most representative words
- Manually select a topic



# Experiments



# Outline

- Motivation
- Data
- Experiments
- **Results**

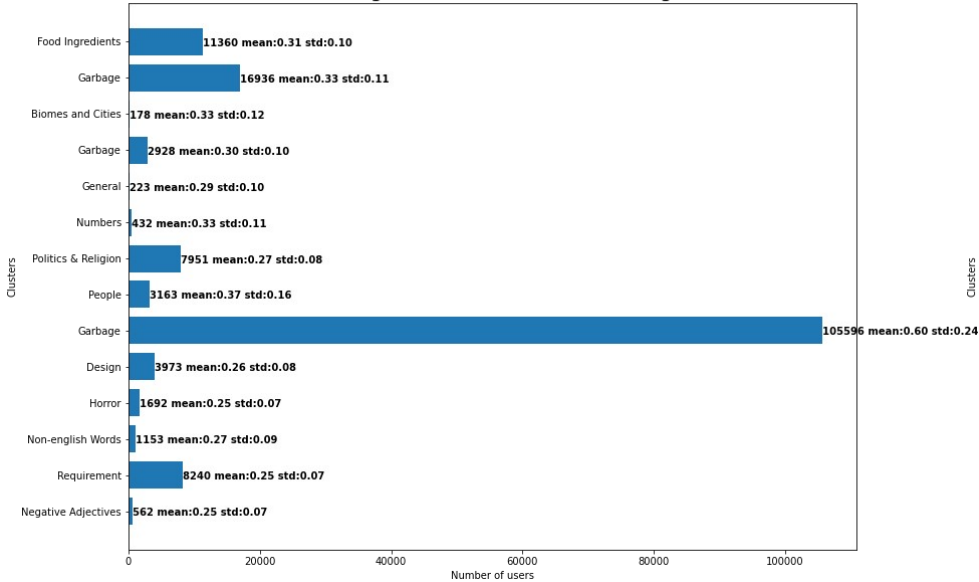
# Results

How the users are distributed?

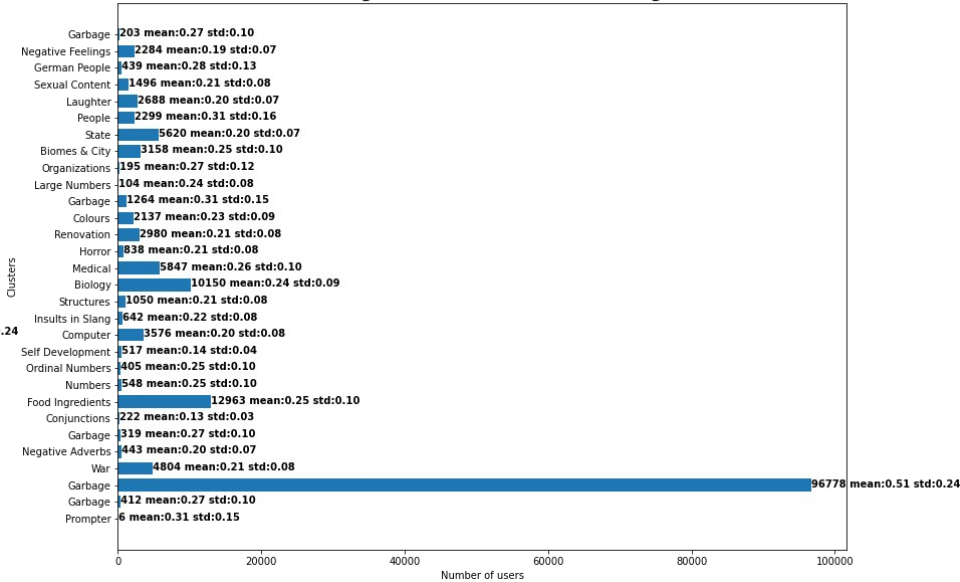
- in Word embedding space
- in Sentence embedding space

# Results

User clustering with Glove word embeddings, N=15

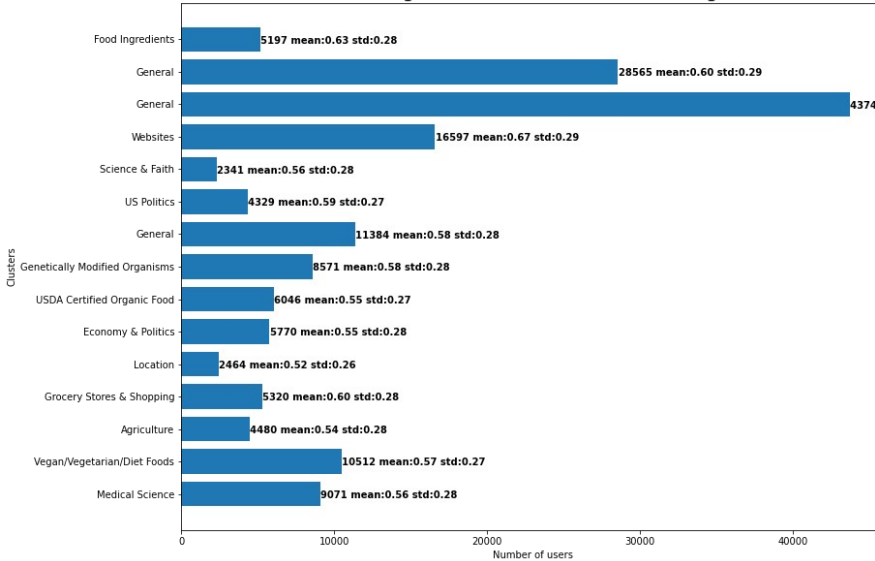


User clustering with Glove word embeddings, N=30

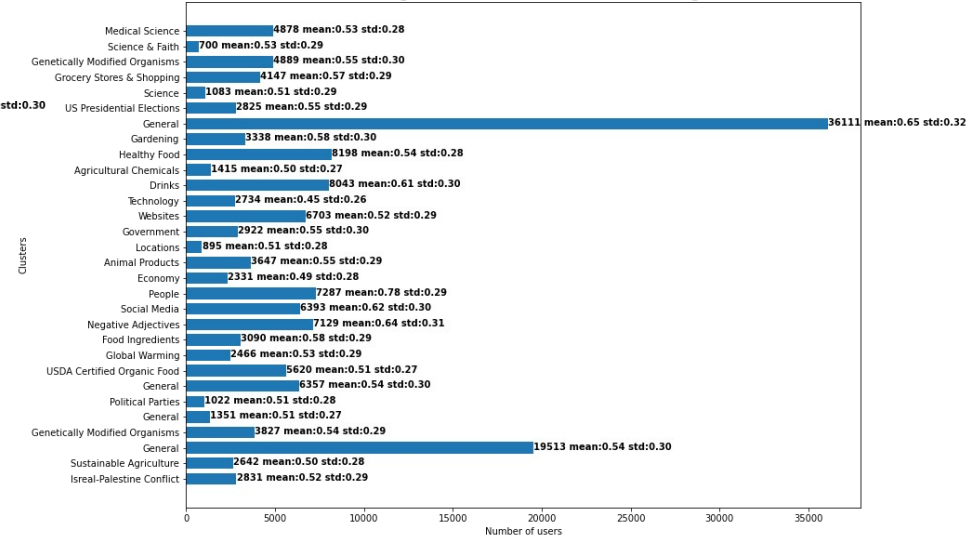


# Results

User clustering with USE sentence embeddings, N=15



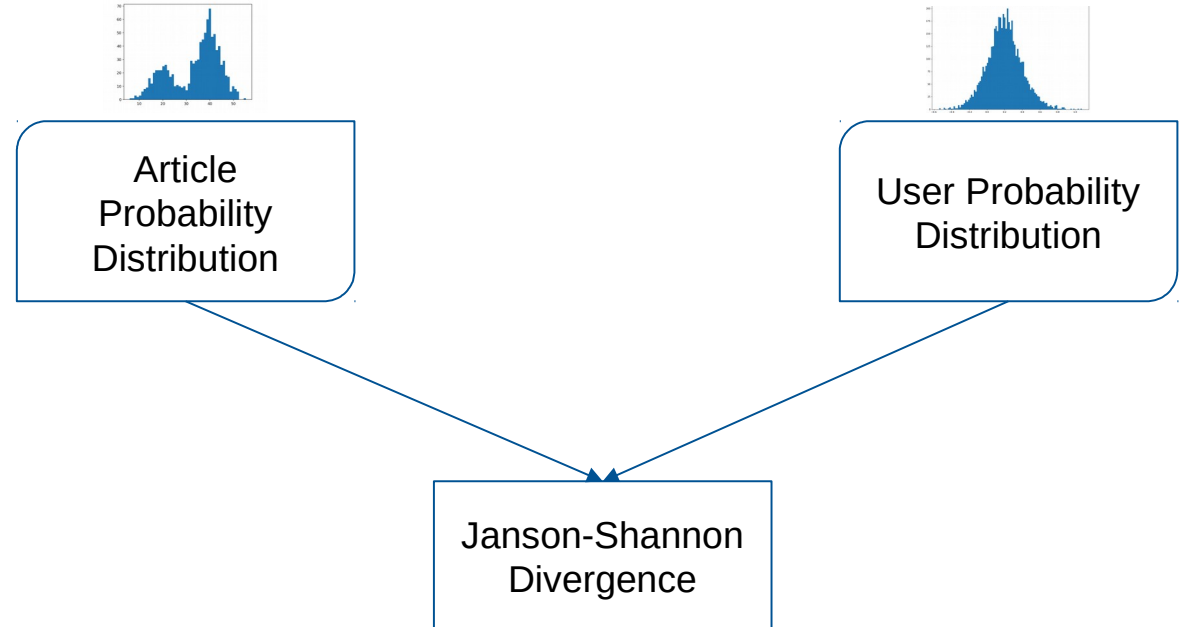
User clustering with USE sentence embeddings, N=30



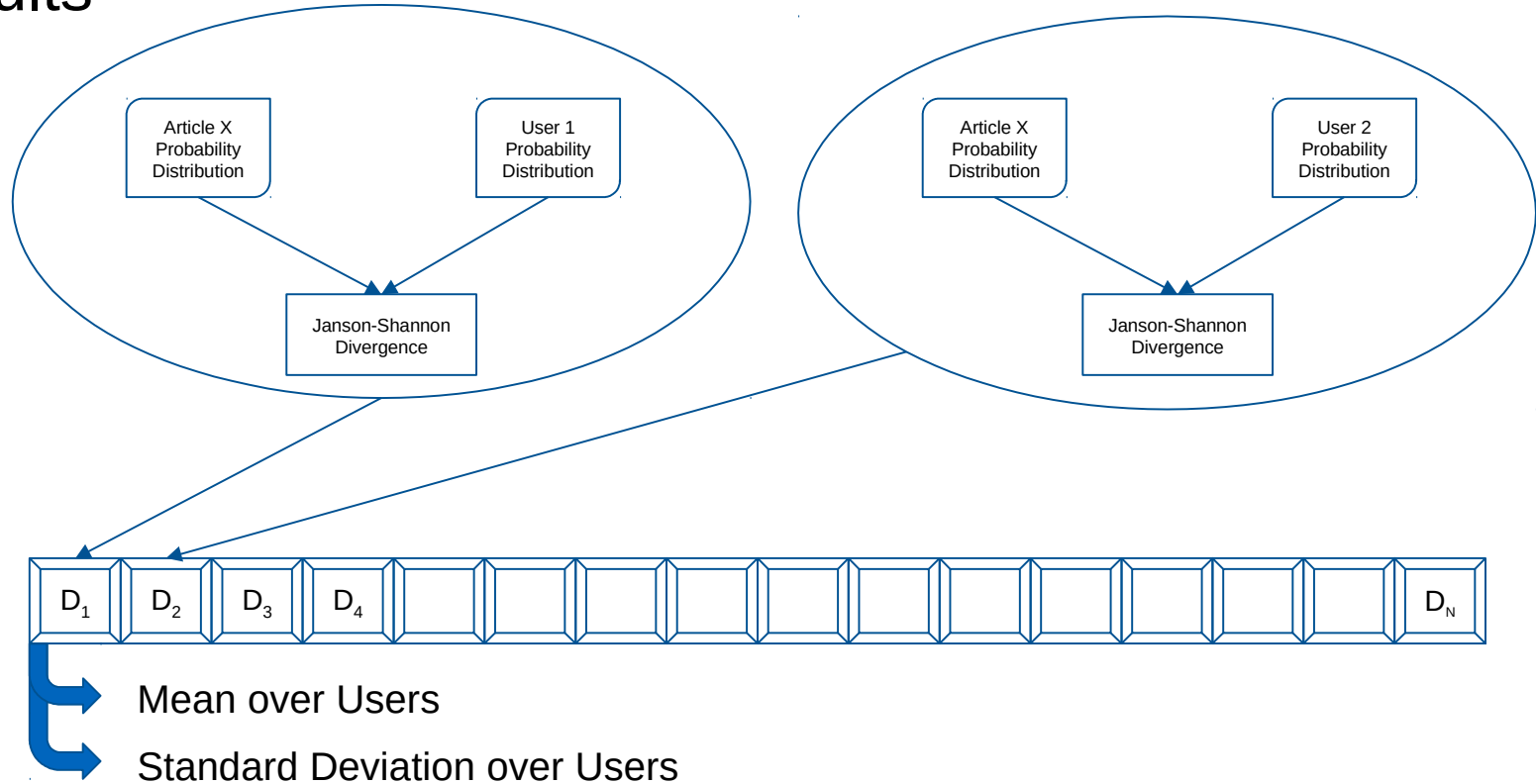


# Results

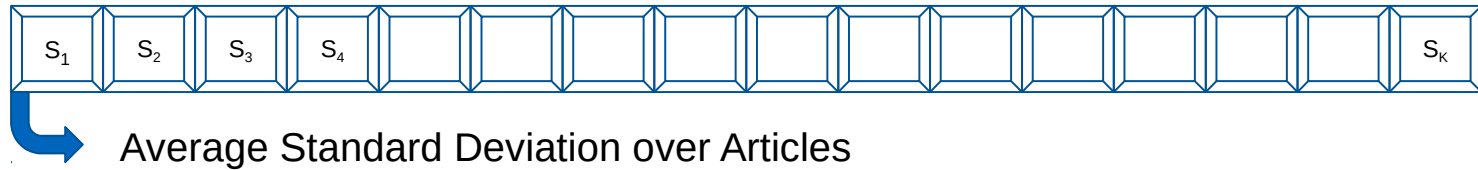
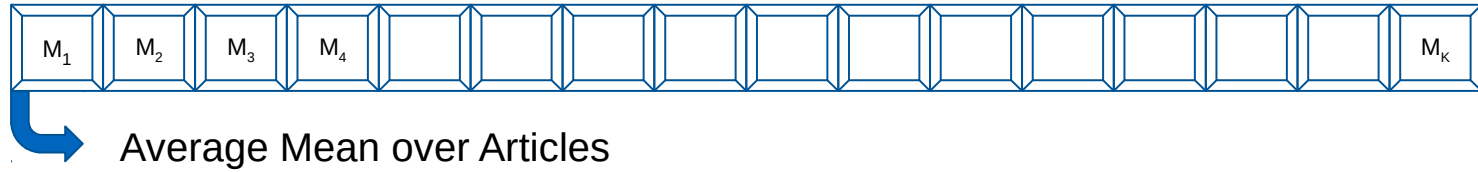
**Hypothesis:** Answering users should be closer to the article than random users.



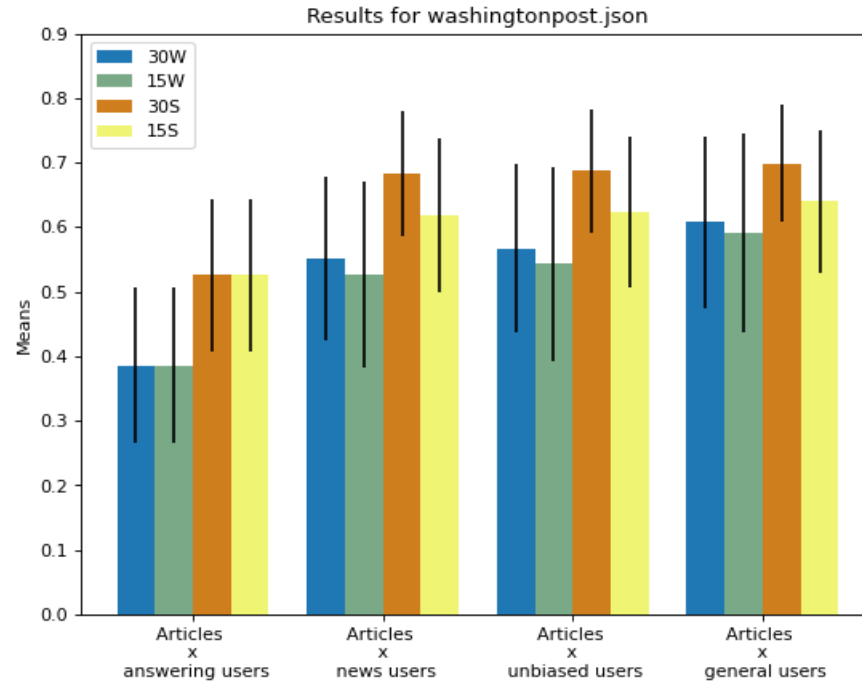
# Results



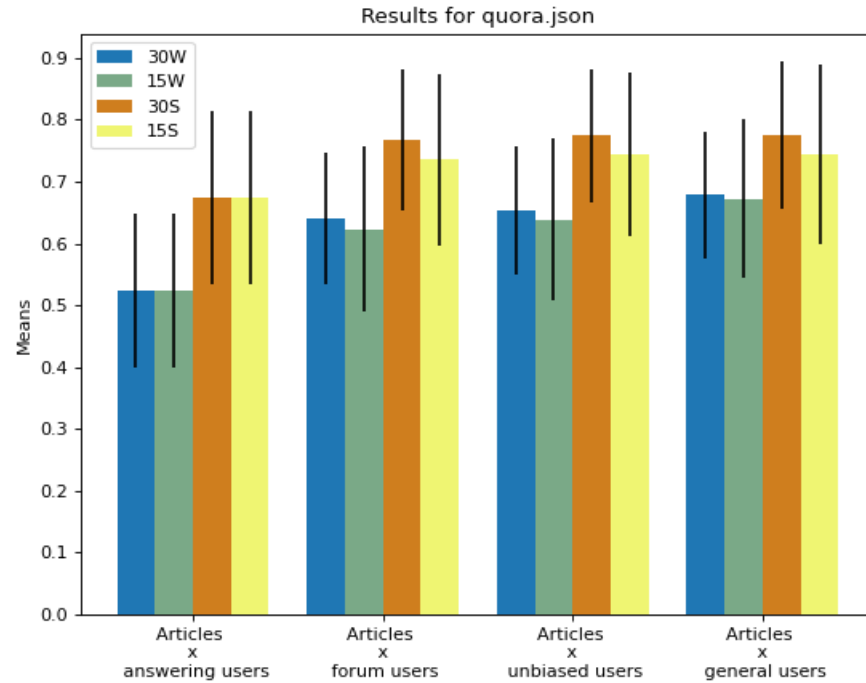
# Results



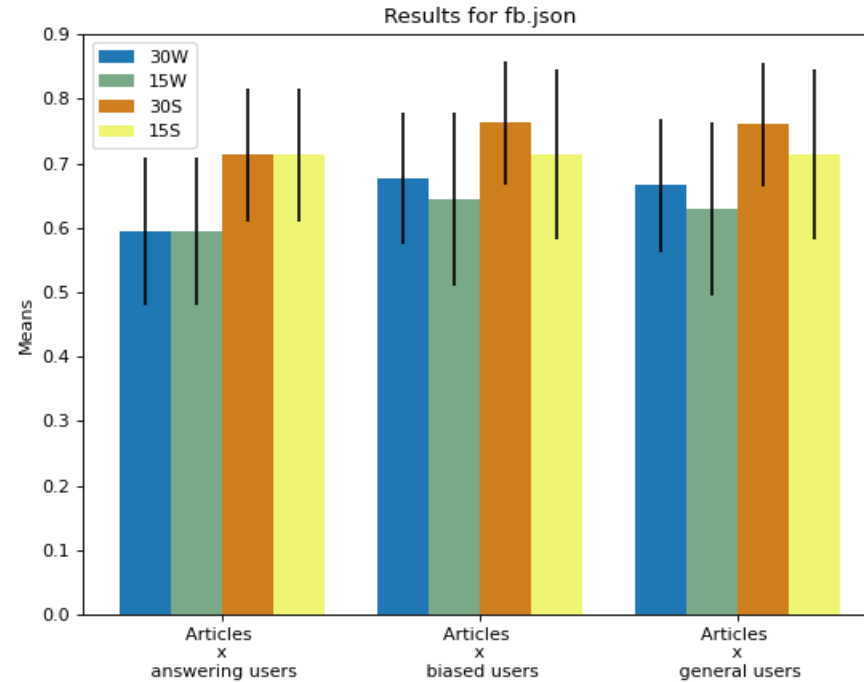
# Results



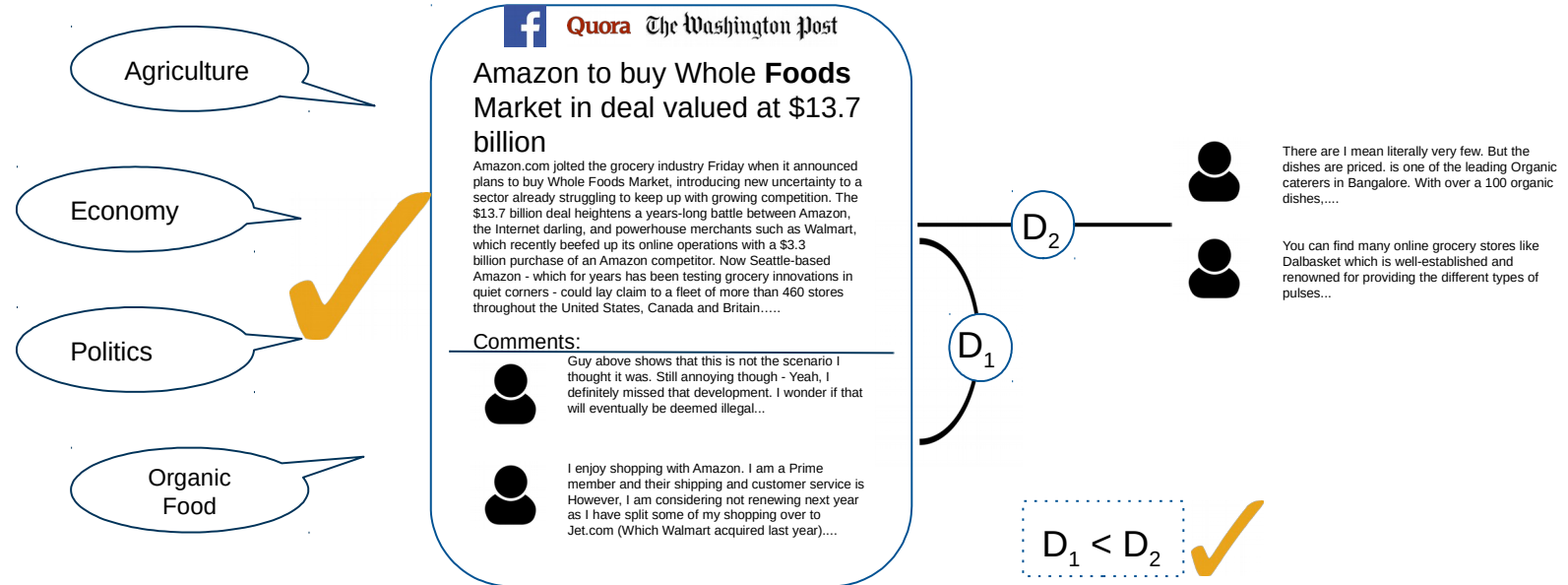
# Results



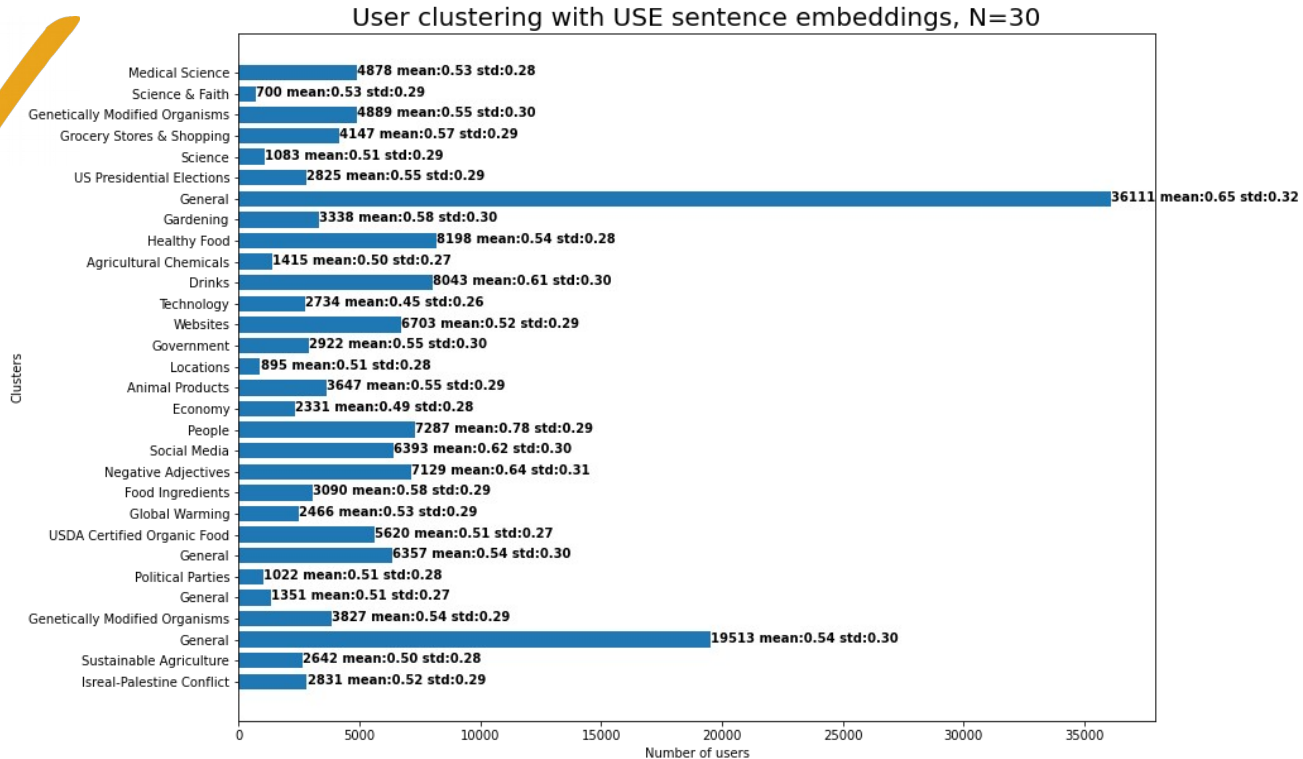
# Results



# Results



# Results





Thanks for listening!

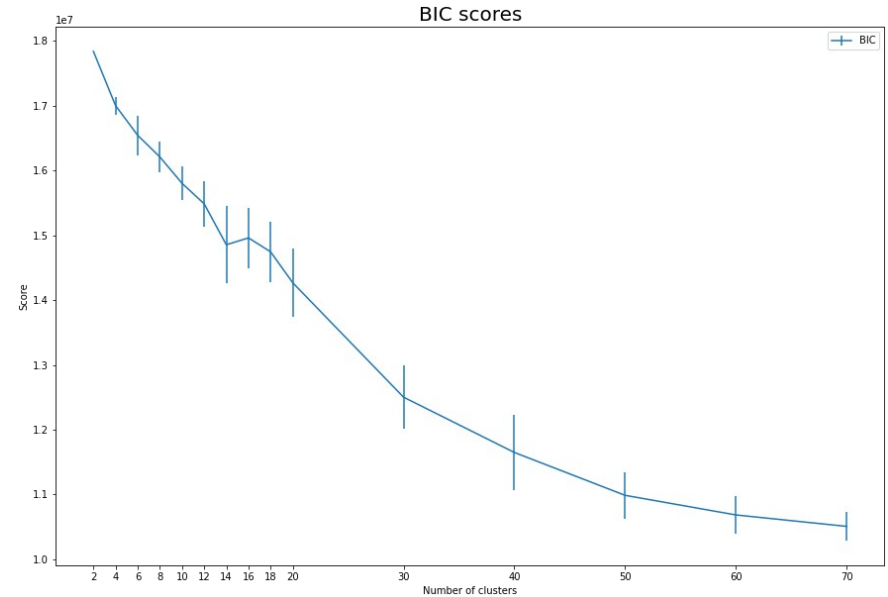
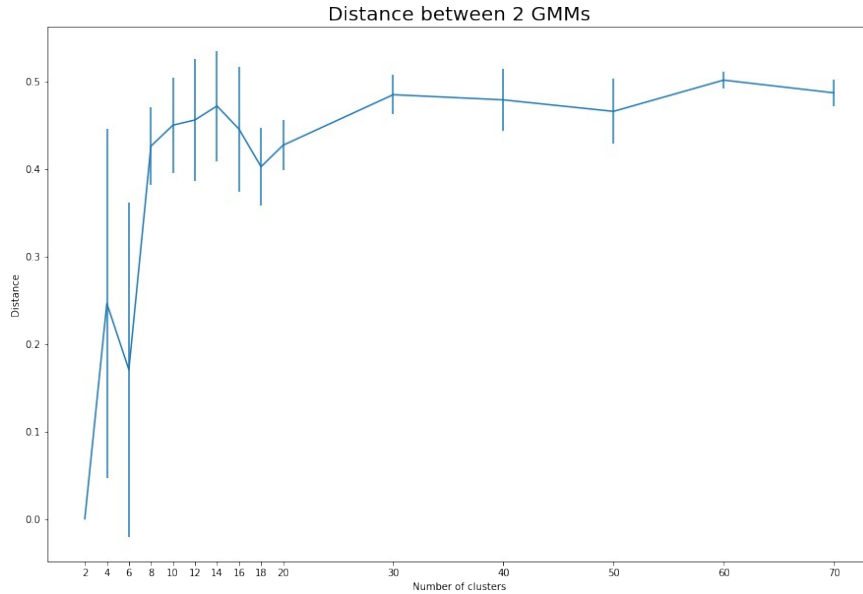
Questions ?

# References

- Pennington, Jeffrey, Richard Socher, and Christopher D. Manning. "Glove: Global vectors for word representation." *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014.
- Sridhar, Vivek Kumar Rangarajan. "Unsupervised topic modeling for short texts using distributed representations of words." *Proceedings of the 1st workshop on vector space modeling for natural language processing*. 2015.
- [Fine tune GloVe embeddings using Mittens](#)
- [A fast implementation of GloVe, with optional retrofitting](#)
- [Basics of using pre-trained glove vectors in python](#)
- [Sklearn Gaussian Mixture Model](#)
- [Gaussian Mixture Model Clusterization how to select the number of components clusters](#)
- [KL divergence of two gmms](#)
- [How to find the similarity between two probability distributions using python](#)
- Cer, D., Yang, Y., Kong, S. Y., Hua, N., Limtiaco, N., John, R. S., ... & Sung, Y. H. (2018). Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.
- Angelidis, Stefanos, and Mirella Lapata. "Summarizing opinions: Aspect extraction meets sentiment prediction and they are both weakly supervised." *arXiv preprint arXiv:1808.08858* (2018).
- Kim, Han Kyul, Hyunjoong Kim, and Sungzoon Cho. "Bag-of-concepts: Comprehending document representation through clustering words in distributed representation." *Neurocomputing* 266 (2017): 336-352.

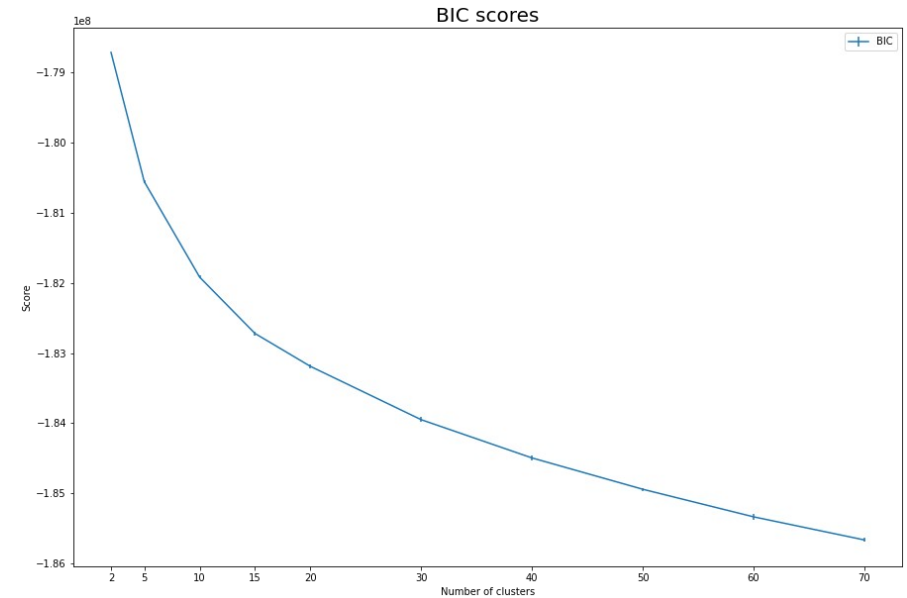
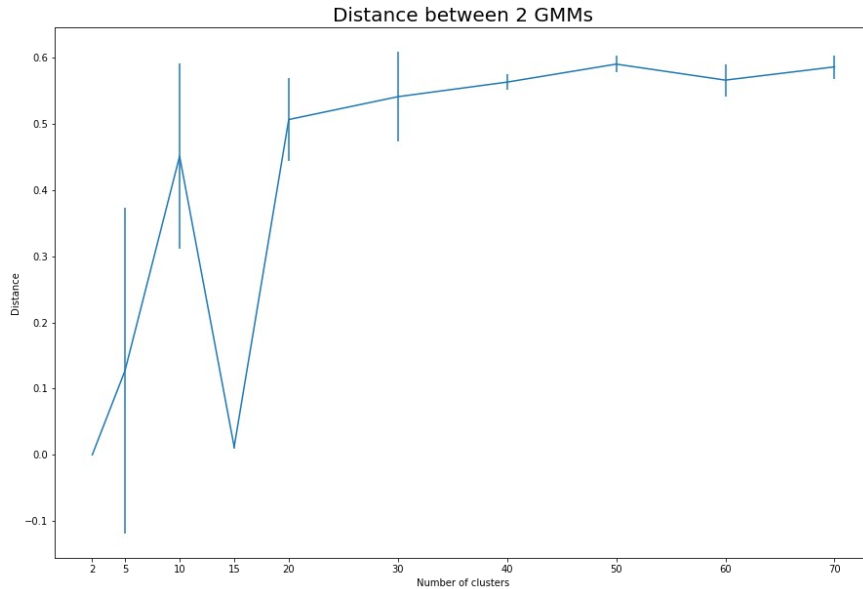
# Experiments

Optimal number of clusters for word embeddings:



# Experiments

Optimal number of clusters for sentence embeddings:



# Clarity Score

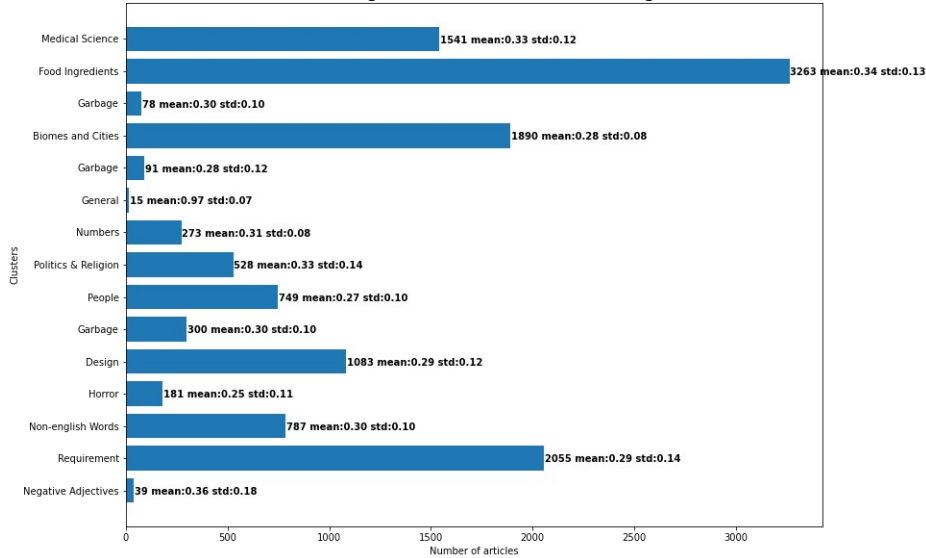
**Definition:** Clarity score measures how much more likely it is to observe word  $w$  in the subset of segments(sentences in our case) that discuss aspect(cluster, gaussian)  $a$ , compared to the corpus as a whole.

$$\text{score}_a(w) = t_a(w) \log_2 \frac{t_a(w)}{t(w)}$$

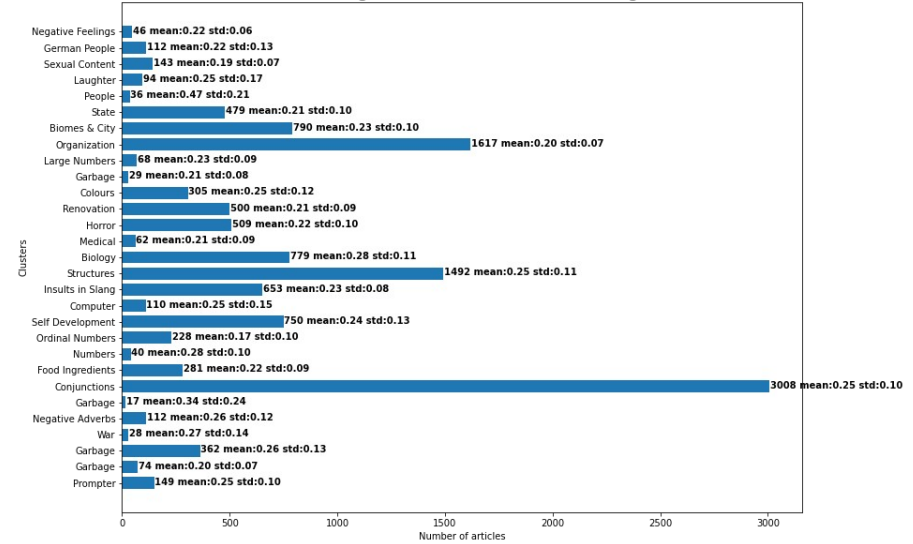
$t_a(w)$ : l1-normalized tf-idf score of  $w$  in the segments annotated with aspect  $a$ .  
 $t(w)$ : l1-normalized tf-idf score of  $w$  in all of the segments.

# Article Clustering

Article clustering with Glove word embeddings, N=15

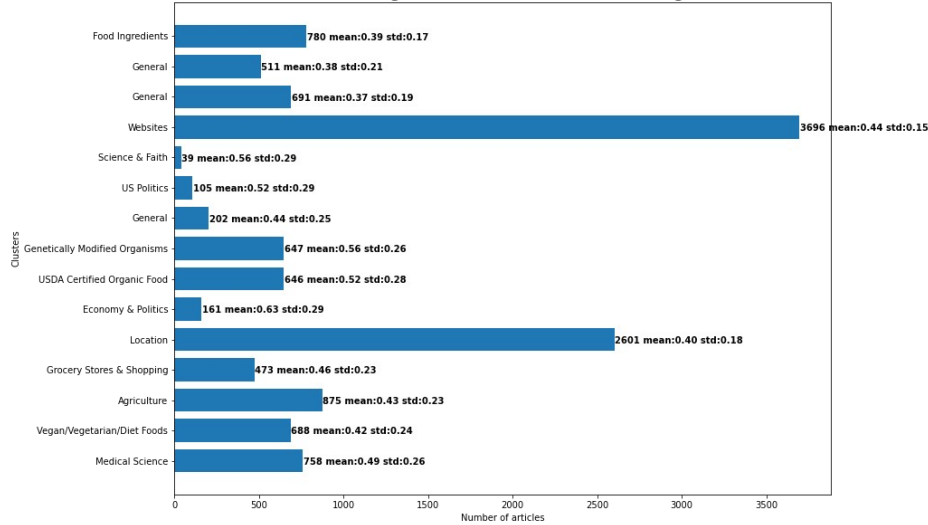


Article clustering with Glove word embeddings, N=30



# Article Clustering

Article clustering with USE sentence embeddings, N=15



Article clustering with USE sentence embeddings, N=30

