

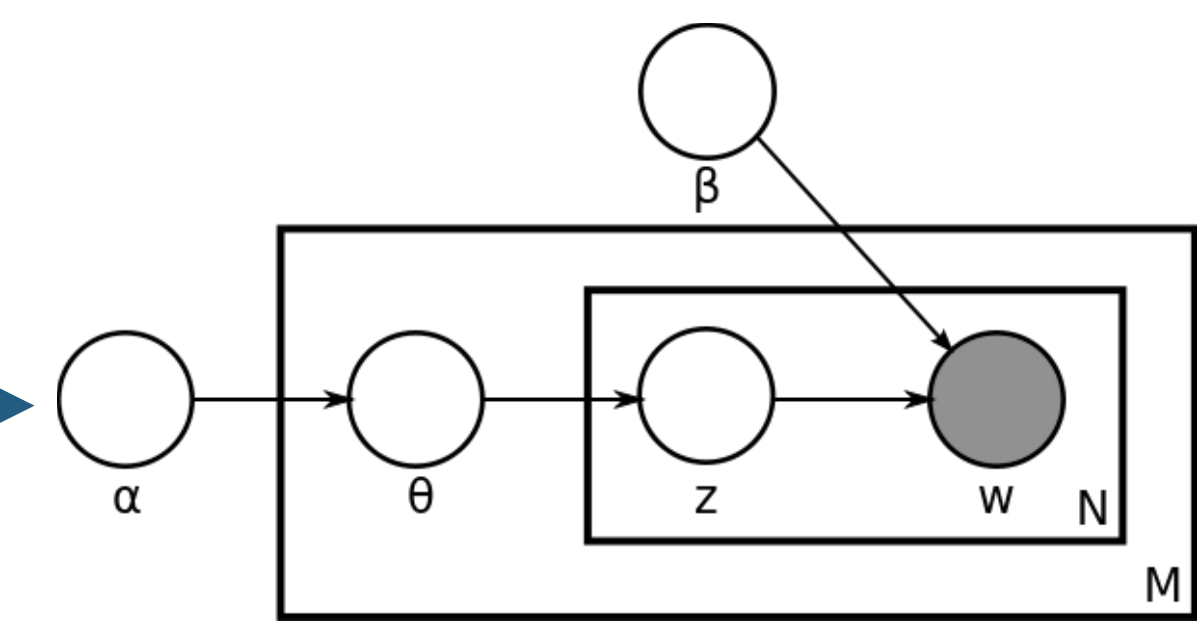
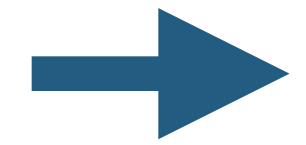
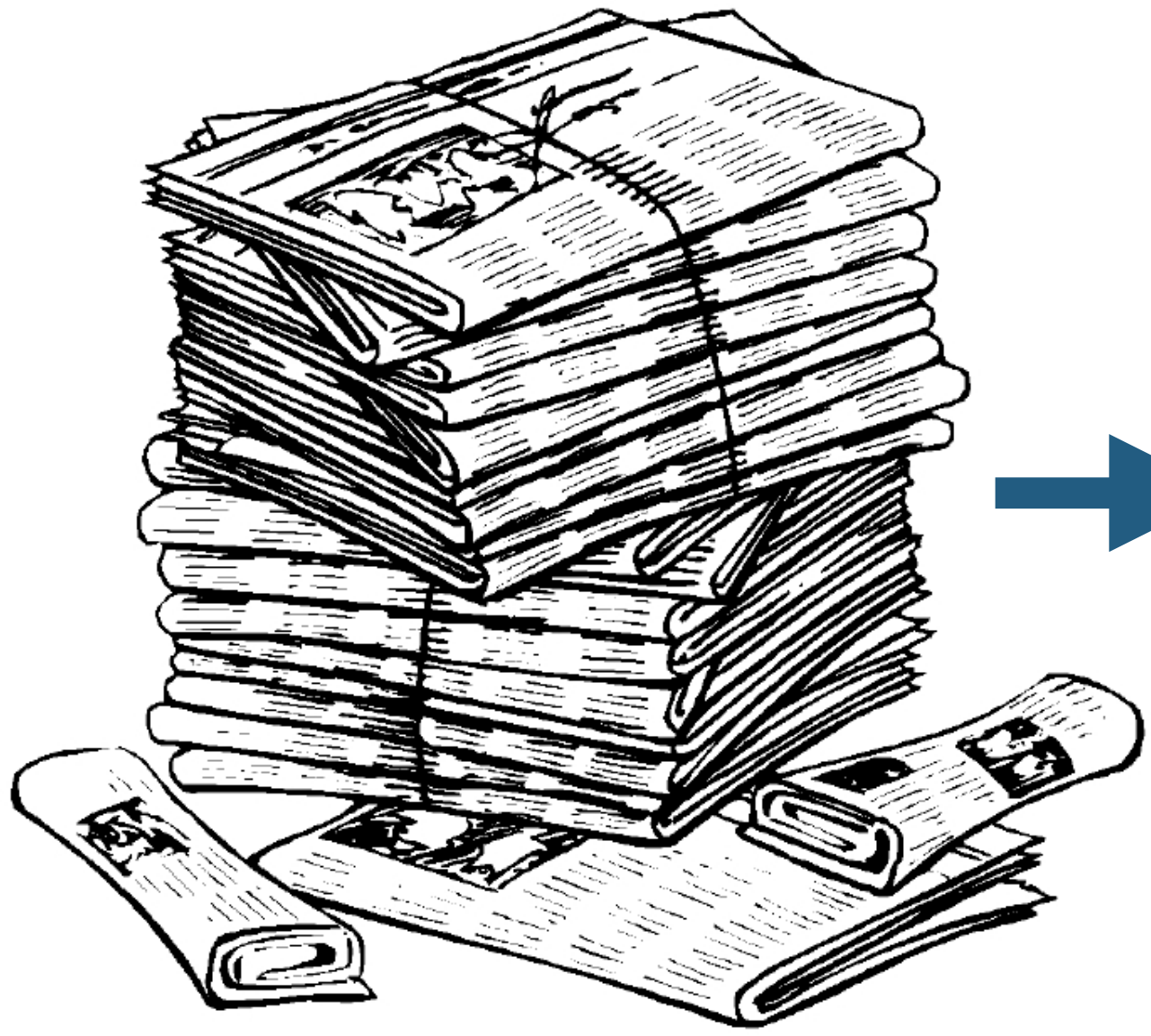
Visualizing Topic Models

Ben Mabey
@bmabey

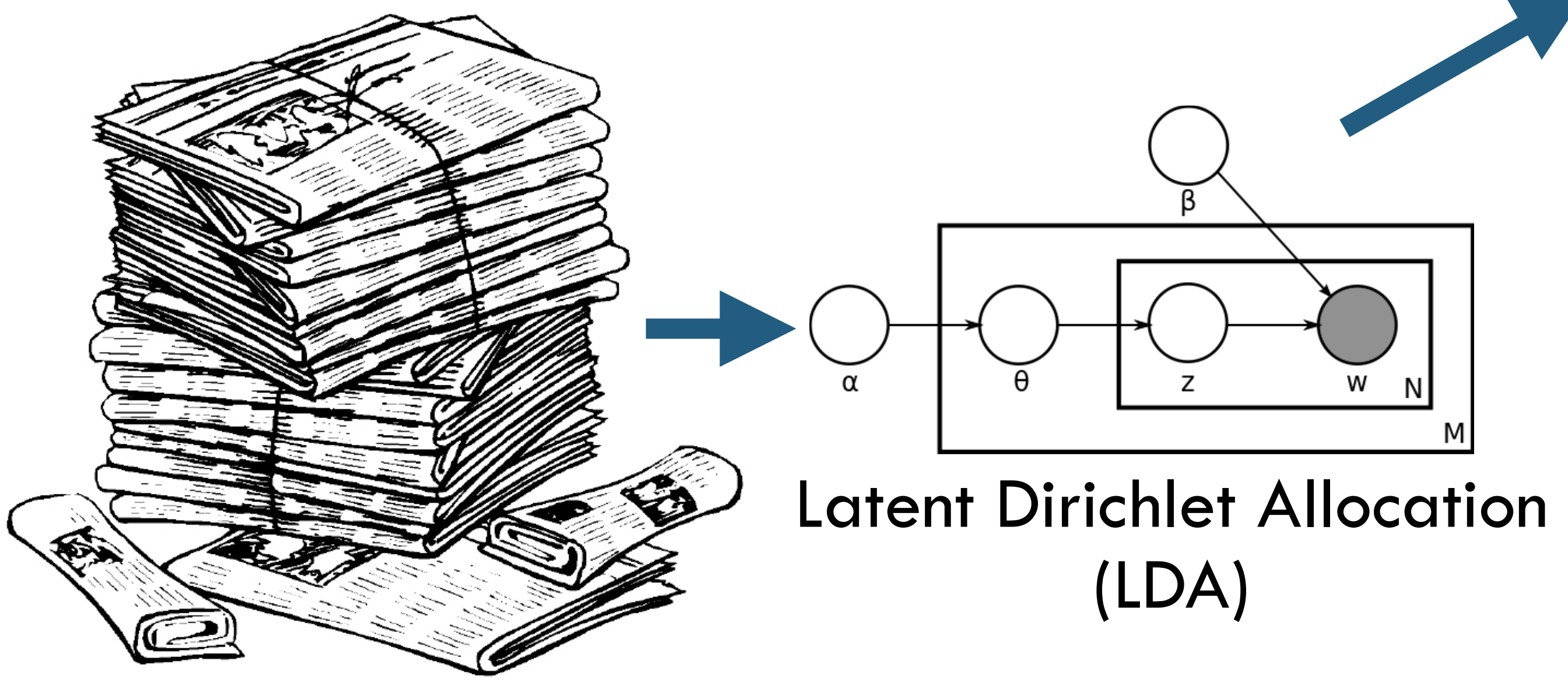


DATA SCIENCE SUMMIT &
DATO CONFERENCE 2015



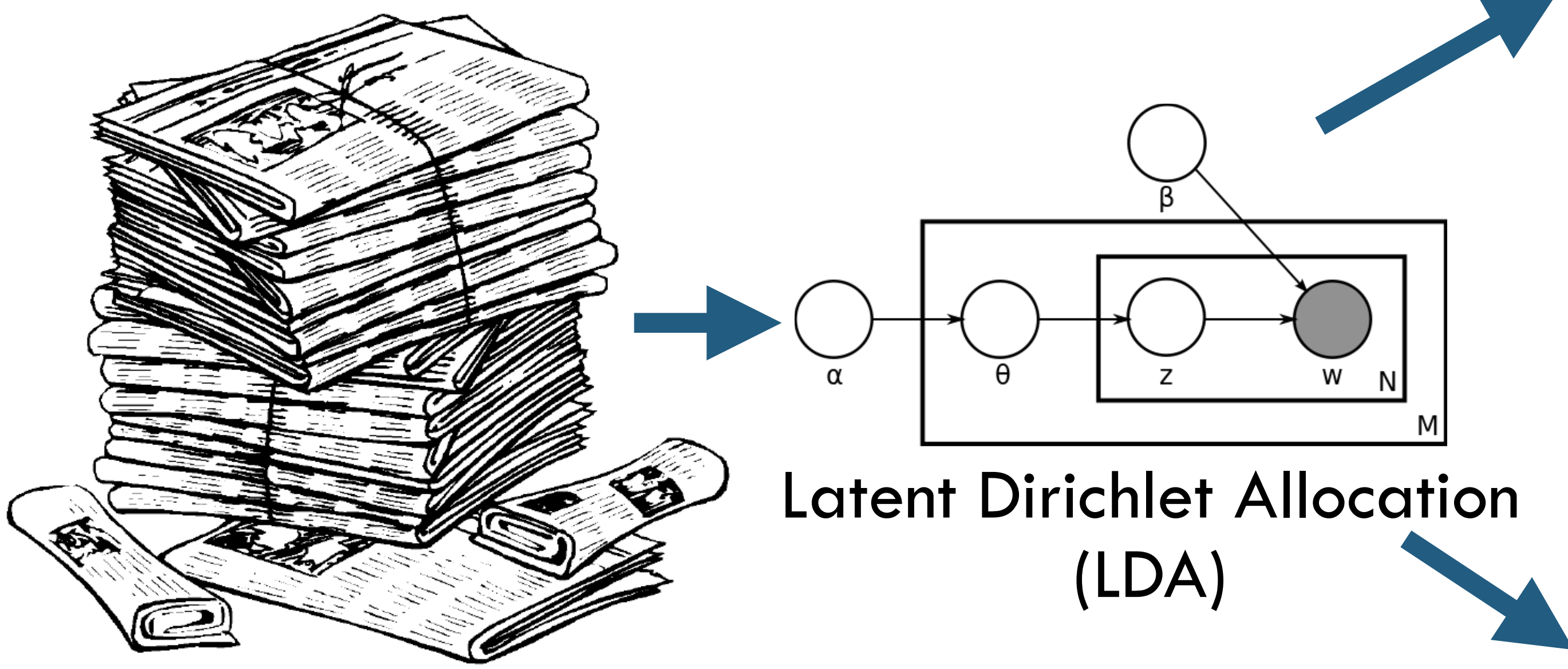


Latent Dirichlet Allocation (LDA)



Document-Topic Distributions

	<i>0</i>	<i>1</i>	...	<i>k</i>
<i>doc a</i>	0.25	0.14	...	0.02
<i>doc b</i>	0.01	0.30	...	0.09
...	0.31
<i>doc D</i>	0.13	0.07	...	0.01



Document-Topic Distributions

	0	1	...	k
<i>doc a</i>	0.25	0.14	...	0.02
<i>doc b</i>	0.01	0.30	...	0.09
...	0.31
<i>doc D</i>	0.13	0.07	...	0.01

Term-Topic Distributions

	0	1	...	k
<i>bird</i>	0.002	0.01	...	0.004
<i>coffee</i>	0.001	0.003	...	0.009
...	0.031
<i>work</i>	0.002	0.006	...	0.021





250k+ stories
July 2007 - May 2014



250k+ stories
July 2007 - May 2014



POS tagging w/spaCy



250k+ stories
July 2007 - May 2014



POS tagging w/spaCy



Phrase detection w/Gensim



Hacker News

250k+ stories
July 2007 - May 2014



POS tagging w/spaCy



Phrase detection w/Gensim



Stopword removal &
only kept nouns or phrases with nouns



250k+ stories
July 2007 - May 2014



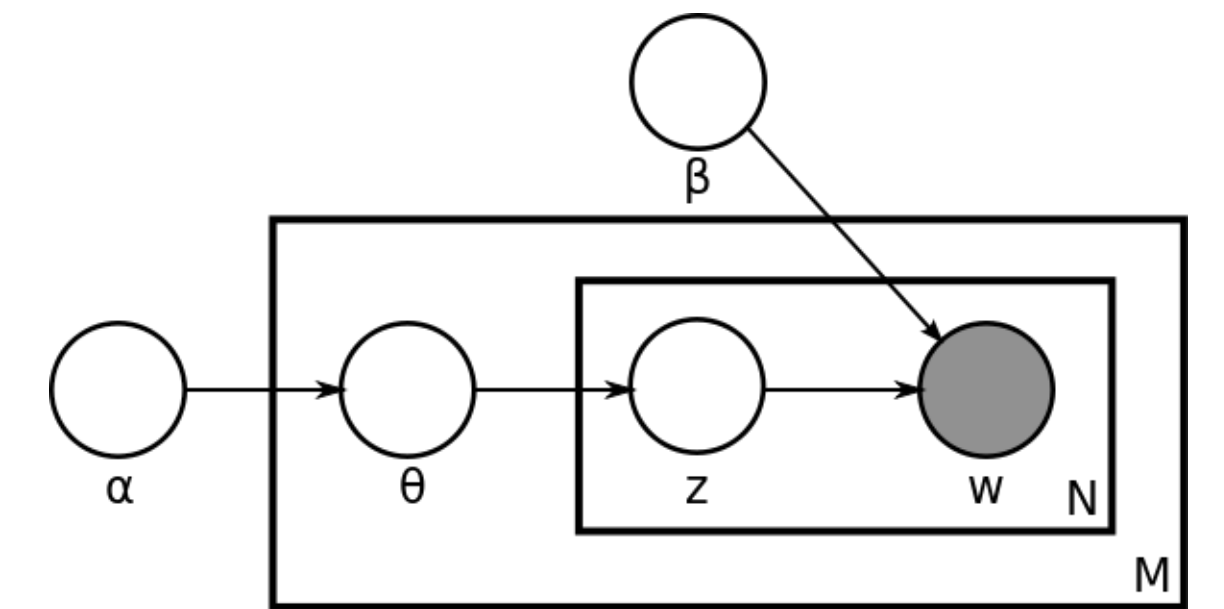
POS tagging w/spaCy



Phrase detection w/Gensim



Stopword removal &
only kept nouns or phrases with nouns



Fit LDA models varying
the number of topics



Hacker News



Game written by 14 year old passes Angry Birds as the top free iphone app



Game written by 14 year old passes Angry Birds as the top free iphone app

Document-Topic Distribution

Topic	$P(T D)$
58	0.19
38	0.14
16	0.06
...	...



Game written by 14 year old passes Angry Birds as the top free iphone app

Document-Topic Distribution

<i>Topic</i>	<i>P(T D)</i>
58	0.19
38	0.14
16	0.06
...	...

Sorted Topic-Term Distributions

58	38	16
<i>app</i>	<i>game</i>	<i>language</i>
<i>developer</i>	<i>player</i>	<i>code</i>
<i>mobile</i>	<i>video game</i>	<i>programming</i>
<i>user</i>	<i>gaming</i>	<i>java</i>
<i>app store</i>	<i>developer</i>	<i>programmer</i>



Game written by 14 year old passes Angry Birds as the top free iphone app

Document-Topic Distribution

<i>Topic</i>	<i>$P(T D)$</i>
<i>mobile apps</i>	<i>0.19</i>
<i>38</i>	<i>0.14</i>
<i>16</i>	<i>0.06</i>
<i>...</i>	<i>...</i>

Sorted Topic-Term Distributions

<i>mobile apps</i>	<i>38</i>	<i>16</i>
<i>app</i>	<i>game</i>	<i>language</i>
<i>developer</i>	<i>player</i>	<i>code</i>
<i>mobile</i>	<i>video game</i>	<i>programming</i>
<i>user</i>	<i>gaming</i>	<i>java</i>
<i>app store</i>	<i>developer</i>	<i>programmer</i>



Game written by 14 year old passes Angry Birds as the top free iphone app

Document-Topic Distribution

<i>Topic</i>	<i>$P(T D)$</i>
<i>mobile apps</i>	<i>0.19</i>
<i>video games</i>	<i>0.14</i>
<i>16</i>	<i>0.06</i>
<i>...</i>	<i>...</i>

Sorted Topic-Term Distributions

<i>mobile apps</i>	<i>video games</i>	<i>16</i>
<i>app</i>	<i>game</i>	<i>language</i>
<i>developer</i>	<i>player</i>	<i>code</i>
<i>mobile</i>	<i>video game</i>	<i>programming</i>
<i>user</i>	<i>gaming</i>	<i>java</i>
<i>app store</i>	<i>developer</i>	<i>programmer</i>




Game written by 14 year old passes Angry Birds as the top free iphone app

Document-Topic Distribution

<i>Topic</i>	<i>$P(T D)$</i>
<i>mobile apps</i>	<i>0.19</i>
<i>video games</i>	<i>0.14</i>
<i>programming</i>	<i>0.06</i>
...	...


Sorted Topic-Term Distributions

<i>mobile apps</i>	<i>video games</i>	<i>programming</i>
<i>app</i>	<i>game</i>	<i>language</i>
<i>developer</i>	<i>player</i>	<i>code</i>
<i>mobile</i>	<i>video game</i>	<i>programming</i>
<i>user</i>	<i>gaming</i>	<i>java</i>
<i>app store</i>	<i>developer</i>	<i>programmer</i>



Interpreting Topic Models


What is the meaning of each topic?



Interpreting Topic Models

What is the meaning of each topic?

How prevalent is each topic?




Interpreting Topic Models

What is the meaning of each topic?

How prevalent is each topic?

How do the topics relate to each other?



Interpreting Topic Models

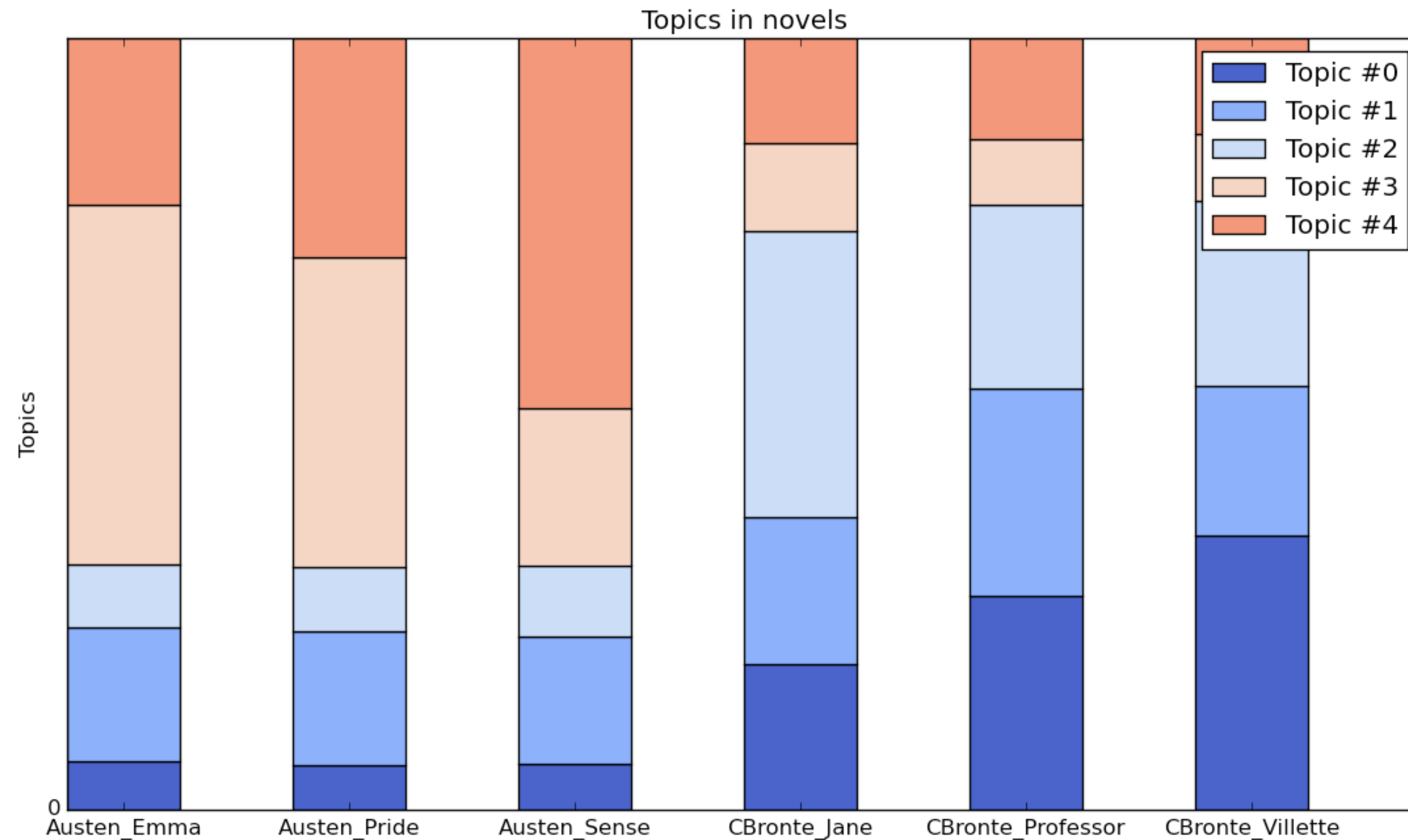
What is the meaning of each topic?

How prevalent is each topic?

How do the topics relate to each other?

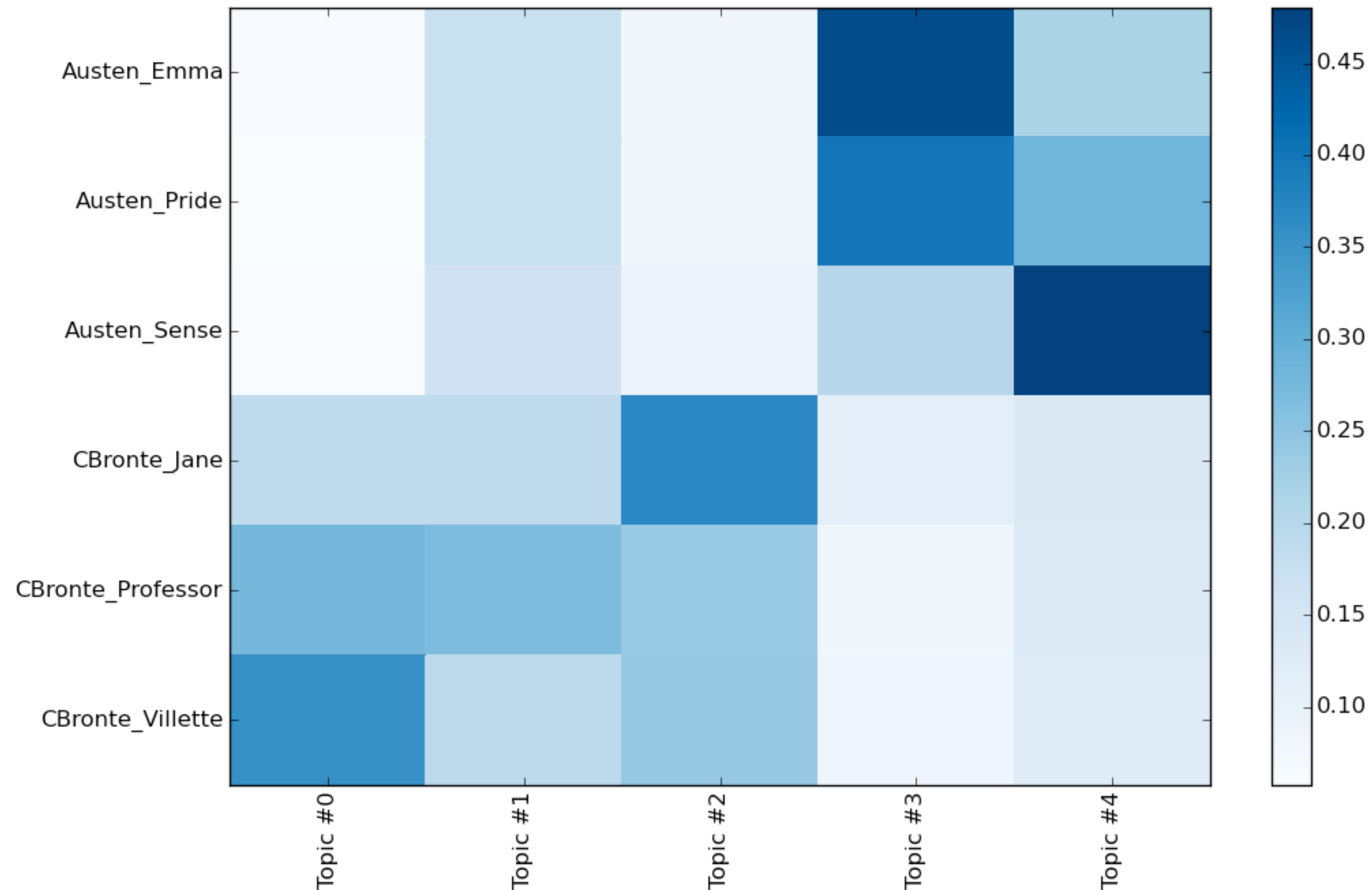
How do the documents relate to each other?

Visualizing Topic Models



https://de.dariah.eu/tatom/topic_model_visualization.html

Visualizing Topic Models

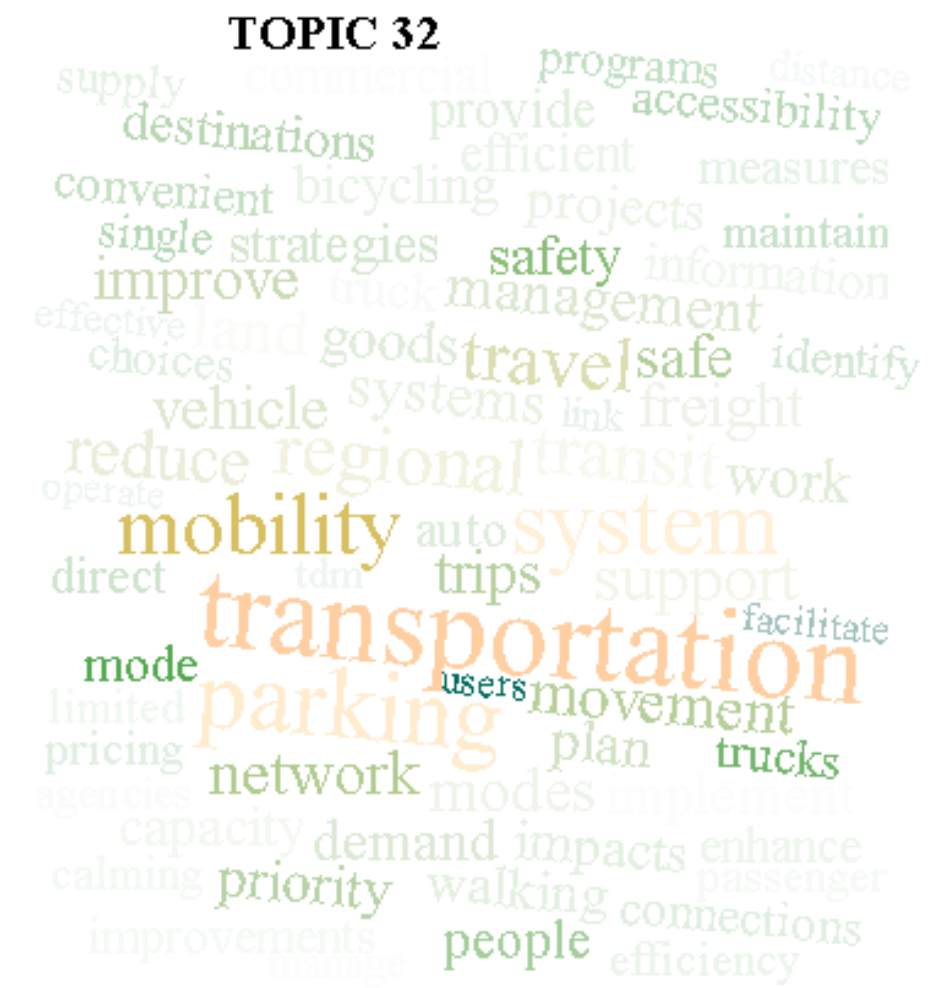


https://de.dariah.eu/tatom/topic_model_visualization.html

Visualizing Topic Models

Topic #0	Topic #1	Topic #2	Topic #3	Topic #4
hand	looked	night	mr	elinor
good	found	room	miss	mother
madame	side	door	mrs	sister
life	speak	long	emma	marianne
heart	girl	house	jane	time
thought	gave	rochester	good	mrs
de	word	round	elizabeth	felt
day	made	hour	thing	letter
monsieur	sense	heard	dear	make
eye	eyes	back	great	john

https://de.dariah.eu/tatom/topic_model_visualization.html

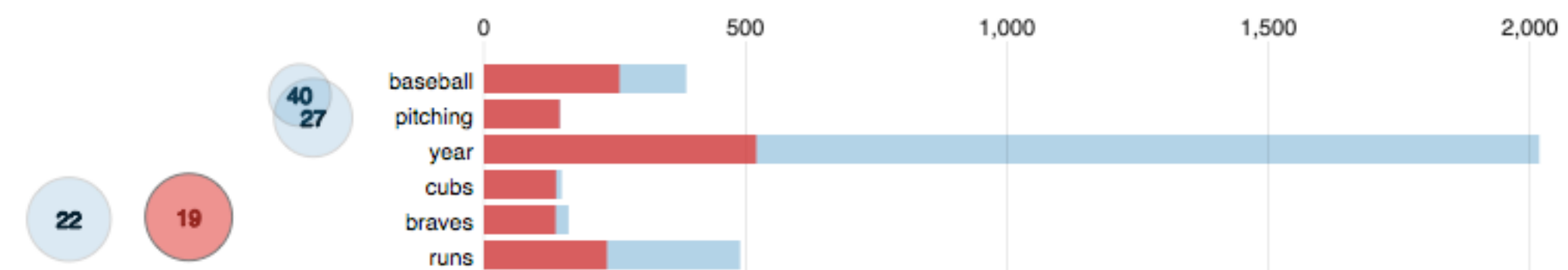


Please don't...

<https://dhs.stanford.edu/algorithmic-literacy/using-word-clouds-for-topic-modeling-results/>



LDavis

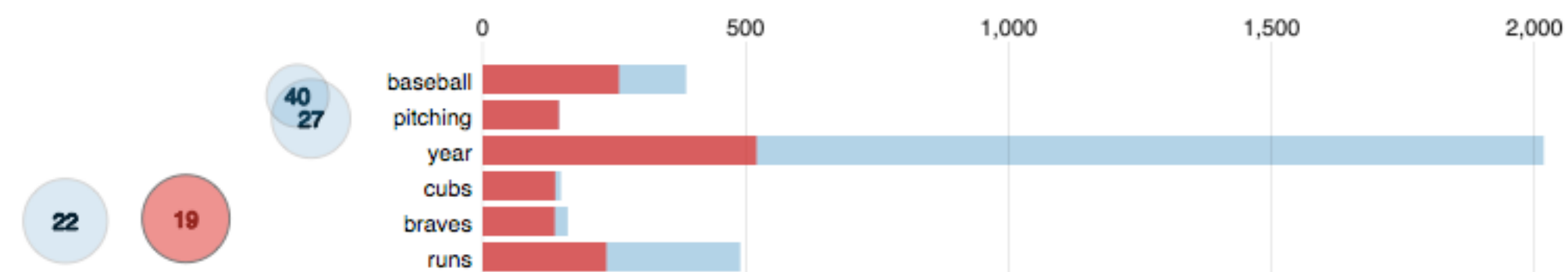


<https://github.com/cpsievert/LDAvis>



~~Put a **py** bird on it~~

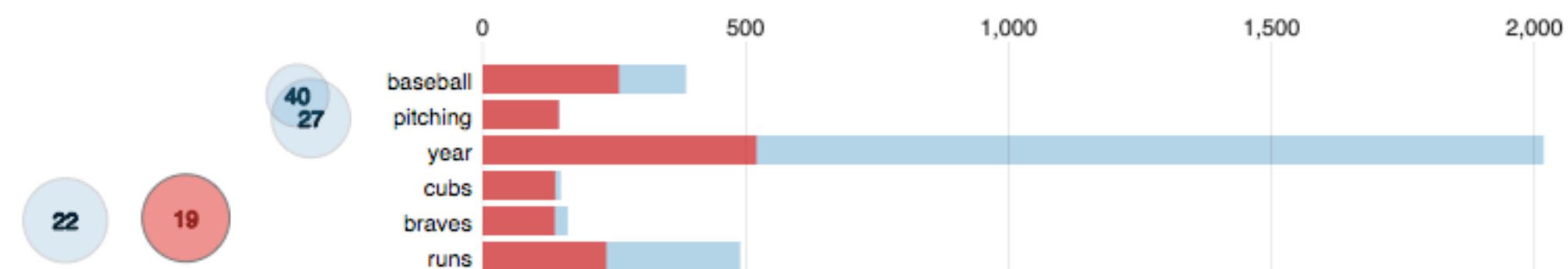
pyLDAvis



<https://github.com/bmabey/pyLDAvis>



~~Put a **py** bird on it~~
pyLDAvis



<https://github.com/bmabey/pyLDAvis>

IP[y]



gensim



Demo Time!



Distinctiveness & Saliency

Termite: Visualization Techniques for Assessing Textual Topic Models
Jason Chuang, Christopher D. Manning and Jeffrey Heer. 2012

measure how much information a term conveys about topics



Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>
<i>game</i>	10	10	50
<i>apple</i>	20	40	20
<i>angry birds</i>	1	1	30
<i>python</i>	50	5	10

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>
<i>game</i>	10	10	50
<i>apple</i>	20	40	20
<i>angry birds</i>	1	1	30
<i>python</i>	50	5	10

$$\text{distinctiveness}(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>
<i>game</i>	10	10	50
<i>apple</i>	20	40	20
<i>angry birds</i>	1	1	30
<i>python</i>	50	5	10

<i>P(T game)</i>	0.14	0.14	0.71
<i>P(T apple)</i>	0.25	0.50	0.25
<i>P(T angry birds)</i>	0.03	0.03	0.94
<i>P(T python)</i>	0.77	0.08	0.15
<i>P(T)</i>	0.33	0.23	0.45

$$distinctiveness(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>
<i>game</i>	10	10	50
<i>apple</i>	20	40	20
<i>angry birds</i>	1	1	30
<i>python</i>	50	5	10

<i>P(T game)</i>	0.14	0.14	0.71
<i>P(T apple)</i>	0.25	0.50	0.25
<i>P(T angry birds)</i>	0.03	0.03	0.94
<i>P(T python)</i>	0.77	0.08	0.15
<i>P(T)</i>	0.33	0.23	0.45

$$\text{distinctiveness}(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>	<i>distinctiveness</i>
<i>game</i>	10	10	50	
<i>apple</i>	20	40	20	
<i>angry birds</i>	1	1	30	0.56
<i>python</i>	50	5	10	

<i>P(T game)</i>	0.14	0.14	0.71
<i>P(T apple)</i>	0.25	0.50	0.25
<i>P(T angry birds)</i>	0.03	0.03	0.94
<i>P(T python)</i>	0.77	0.08	0.15
<i>P(T)</i>	0.33	0.23	0.45

$$distinctiveness(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>	<i>distinctiveness</i>
<i>game</i>	10	10	50	0.15
<i>apple</i>	20	40	20	0.18
<i>angry birds</i>	1	1	30	0.56
<i>python</i>	50	5	10	0.41

<i>P(T game)</i>	0.14	0.14	0.71
<i>P(T apple)</i>	0.25	0.50	0.25
<i>P(T angry birds)</i>	0.03	0.03	0.94
<i>P(T python)</i>	0.77	0.08	0.15
<i>P(T)</i>	0.33	0.23	0.45

$$\text{distinctiveness}(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>	<i>distinctiveness</i>
<i>game</i>	10	10	50	0.15
<i>apple</i>	20	40	20	0.18
<i>angry birds</i>	1	1	30	0.56
<i>python</i>	50	5	10	0.41

$$\text{saliency}(w) = P(w) \times \text{distinctiveness}(w)$$

distinctiveness weighted by the term's overall frequency

$$\text{distinctiveness}(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>	<i>distinctiveness</i>	<i>P(w)</i>
<i>game</i>	10	10	50	0.15	0.28
<i>apple</i>	20	40	20	0.18	0.32
<i>angry birds</i>	1	1	30	0.56	0.13
<i>python</i>	50	5	10	0.41	0.26

$$\text{saliency}(w) = P(w) \times \text{distinctiveness}(w)$$

distinctiveness weighted by the term's overall frequency

$$\text{distinctiveness}(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics

Distinctiveness & Saliency

	<i>coding</i>	<i>tech news</i>	<i>video games</i>	<i>distinctiveness</i>	<i>P(w)</i>	<i>saliency</i>
<i>game</i>	10	10	50	0.15	0.28	0.04
<i>apple</i>	20	40	20	0.18	0.32	0.06
<i>angry birds</i>	1	1	30	0.56	0.13	0.07
<i>python</i>	50	5	10	0.41	0.26	0.11

$$saliency(w) = P(w) \times distinctiveness(w)$$

distinctiveness weighted by the term's overall frequency

$$distinctiveness(w) = \sum_T P(T|w) \log \frac{P(T|w)}{P(T)}$$

computes the KL divergence between the distribution of topics given a term and the marginal distribution of topics



Distinctiveness & Saliency

measure how much information a term conveys about topics...

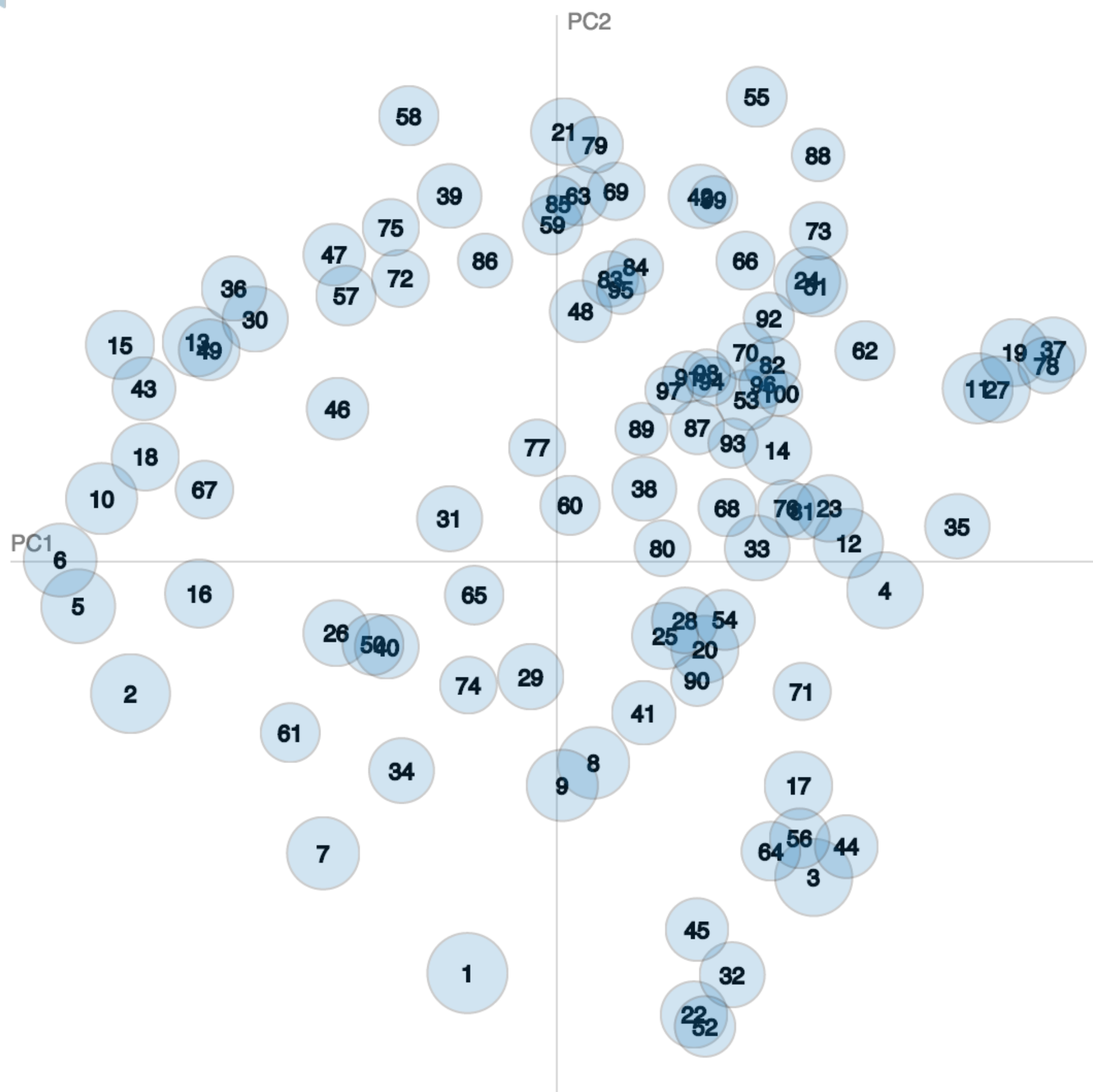


Distinctiveness & Saliency

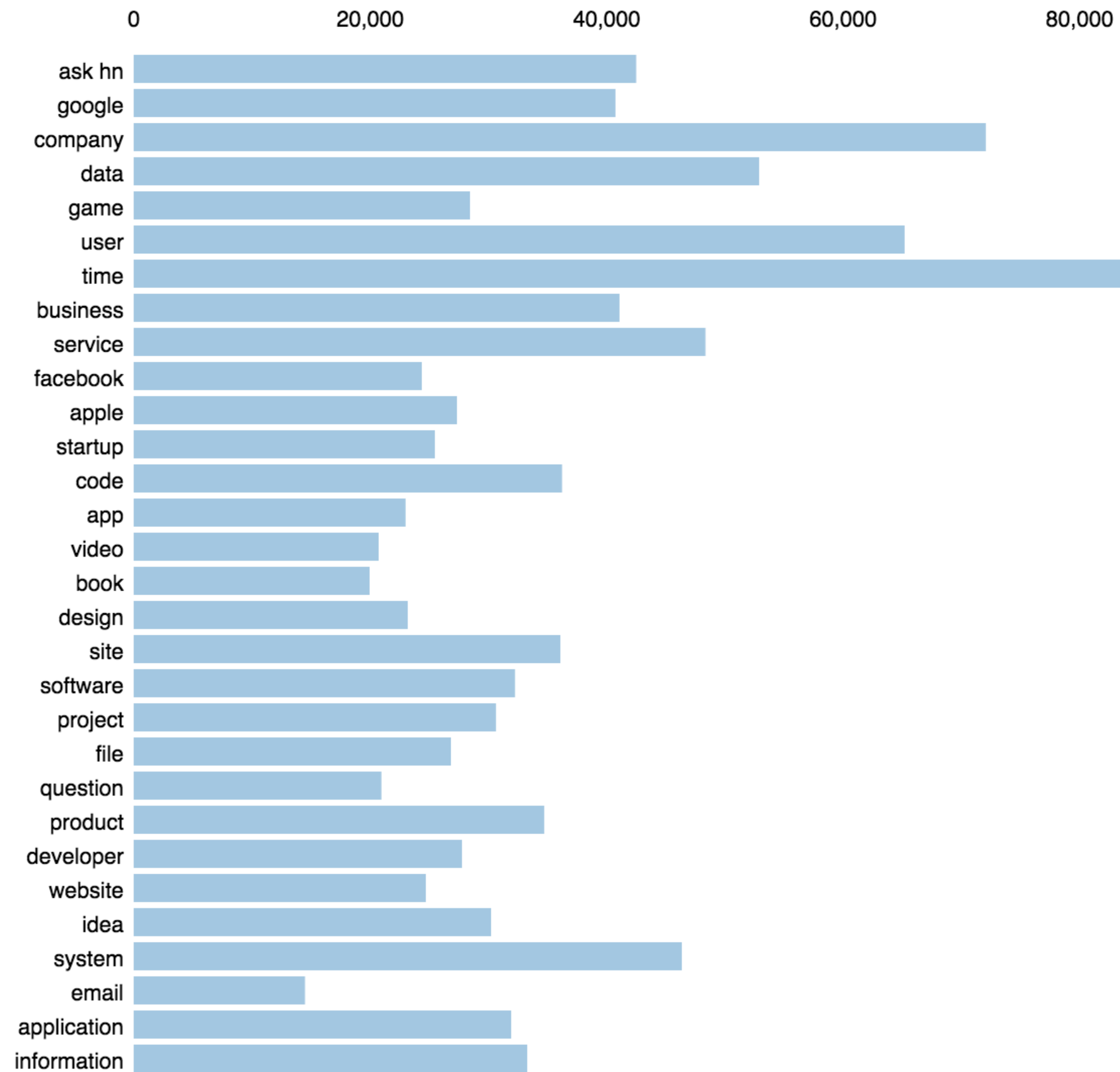
measure how much information a term conveys about topics...

globally

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Salient Terms¹



Thank you!

Learn more at <http://github.com/bmabey/pyLDAvis>
http://nbviewer.ipython.org/github/bmabey/hacker_news_topic_modelling/

Ben Mabey
@bmabey



DATA SCIENCE SUMMIT &
DATO CONFERENCE 2015