



Yet More Topic Modelling

Constructing a topic modelling/clustering,
validation & summarization pipeline

*Martin Wood – Data Scientist, Home Office
(+ ONS, + ESSnet + Brookes...)*

- 1) Outline
- 2) Introduction
- 3) Optimal Clustering
- 4) Validating Clusters
- 5) Summarizing Clusters
- 6) Conclusions



2) Introduction



Latent Dirichelet Allocation
(BOW)

HDBSCAN (Doc2Vec)

OPTIMAL clusters

Coherence

Perplexity

VALIDATE clusters

Centroid separation?

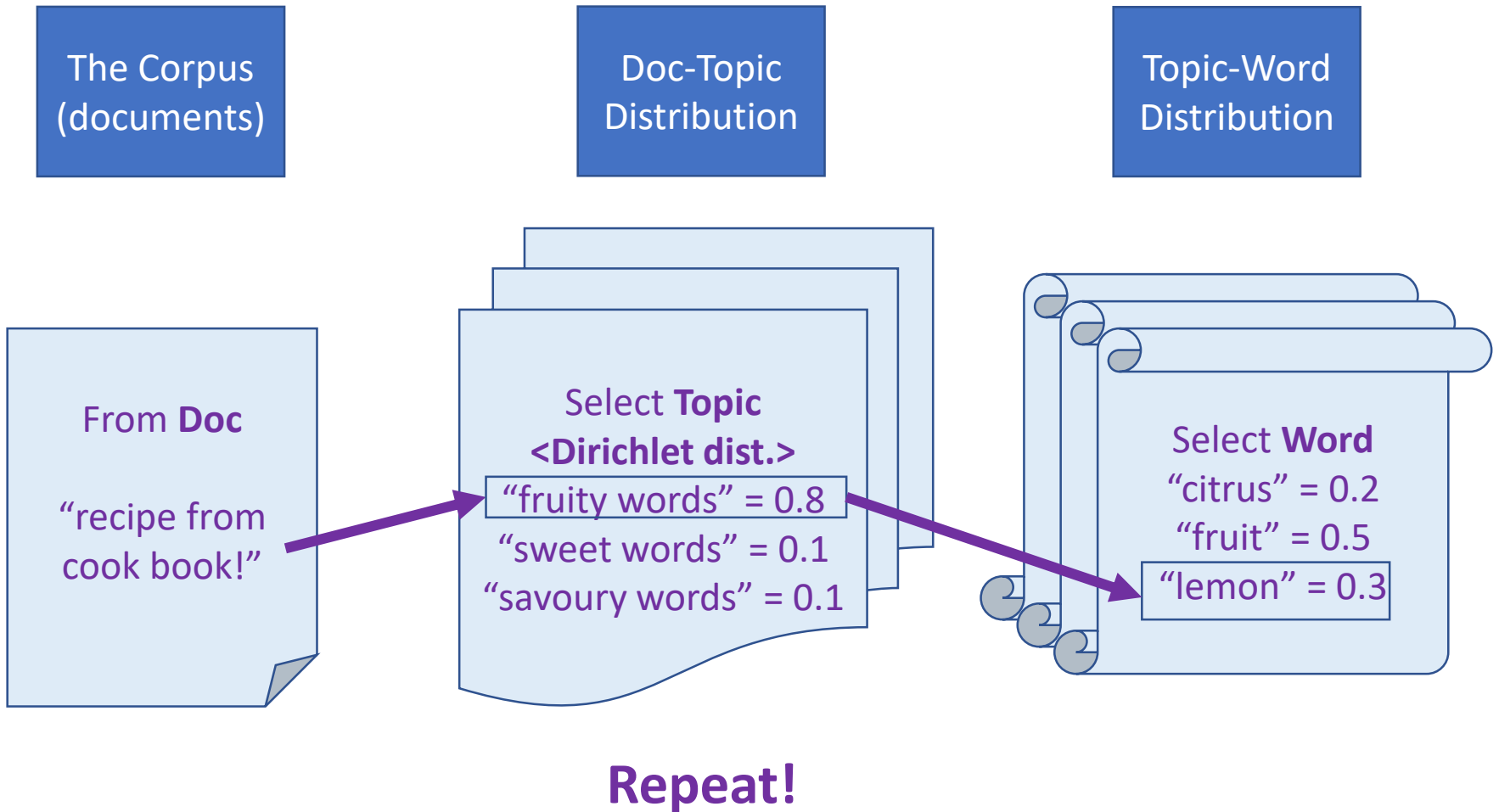
Coherence measure?

Extractive
(representative sentences)

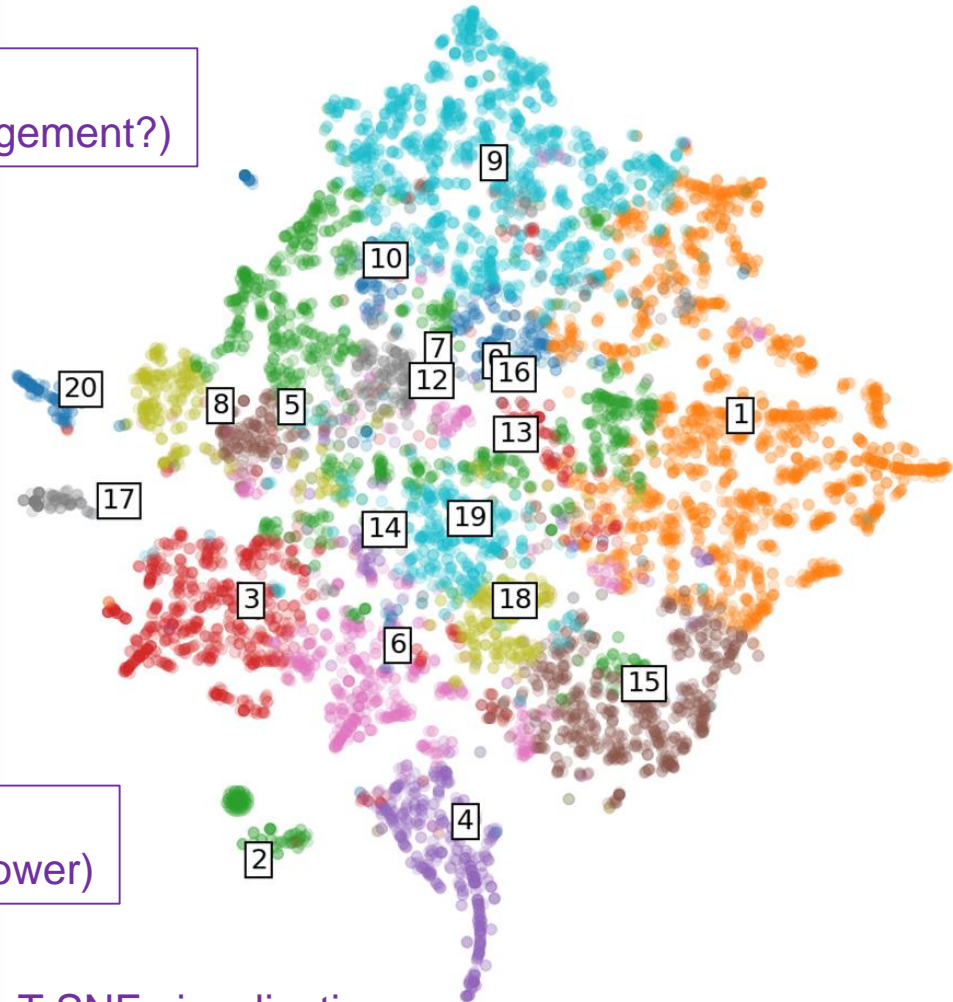
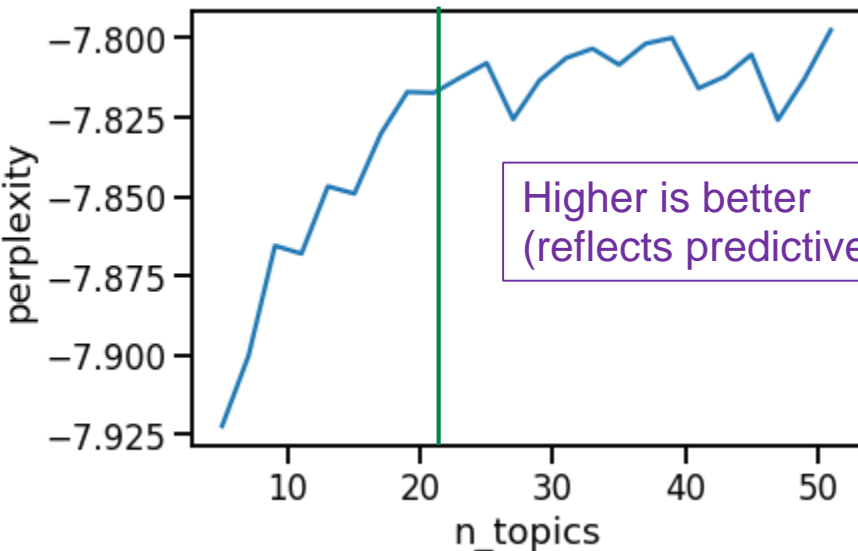
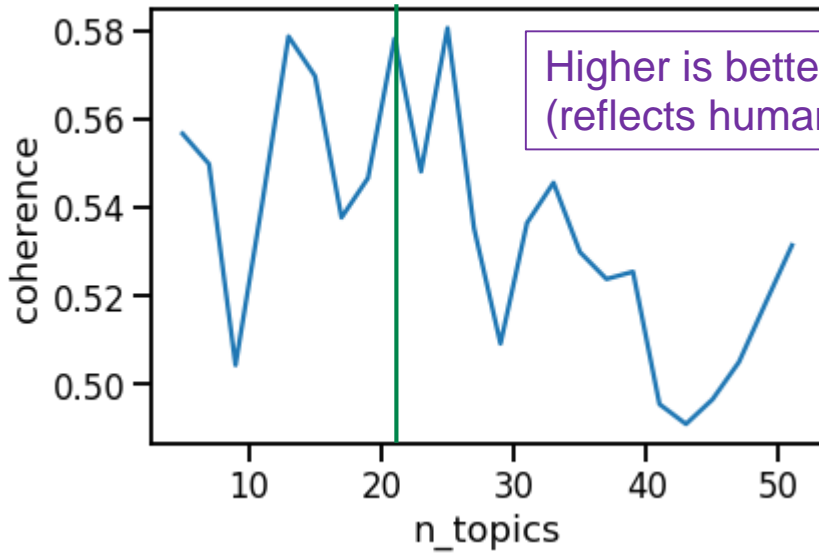
Augmented
(named entity recognition?)

SUMMARIZE clusters

3) Optimising Latent Dirichlet Allocation



3) Optimising Latent Dirichlet Allocation



T-SNE visualisation

No “true” topics, decide how many clusters are useful to you, then locally optimise?

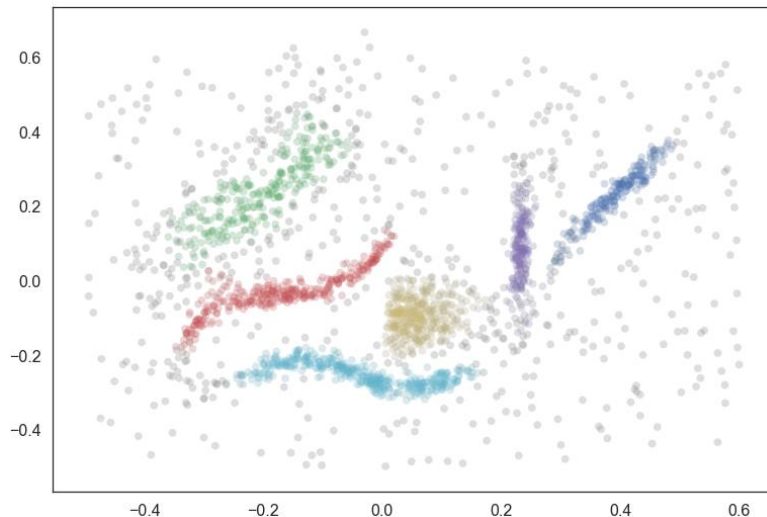
3) Optimising HDBSCAN



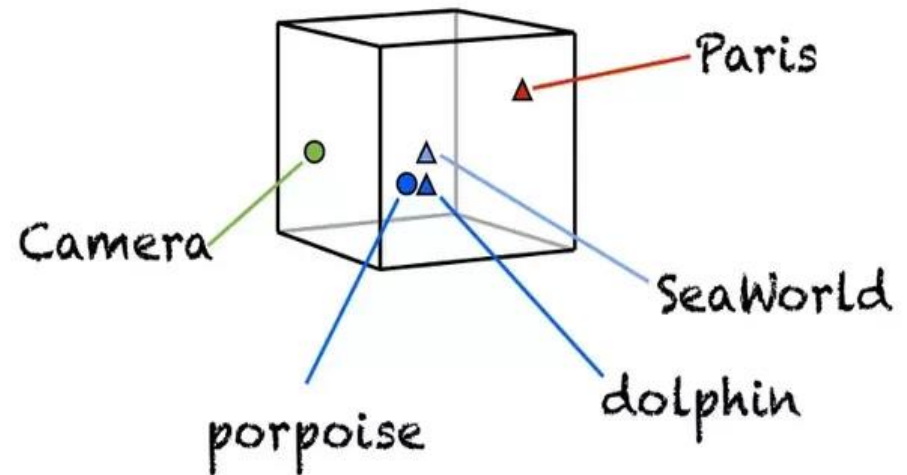
<Hasn't produced pretty graphs yet>

Doc2Vec and then find high-dimensional clusters

Good for very short records, not enough “innate” data to link – use pre-trained models



<https://hdbscan.readthedocs.io/en/latest/index.html>



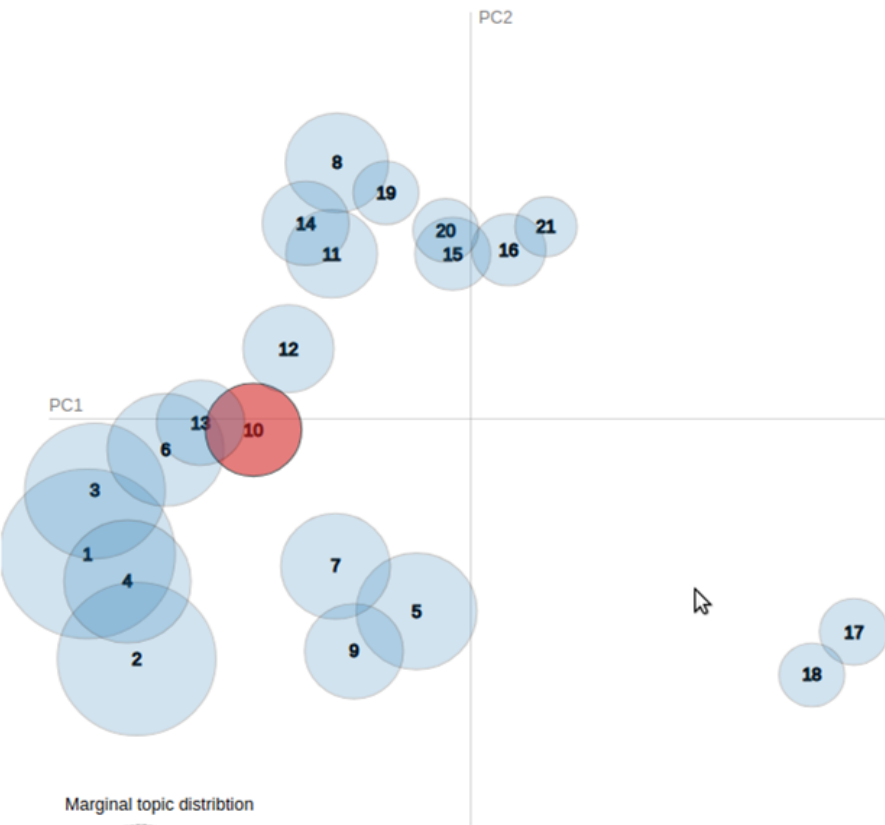
<https://www.kaggle.com/sbongo/do-pretrained-embeddings-give-you-the-extra-edge>

4) Validation – Valid LDA topics

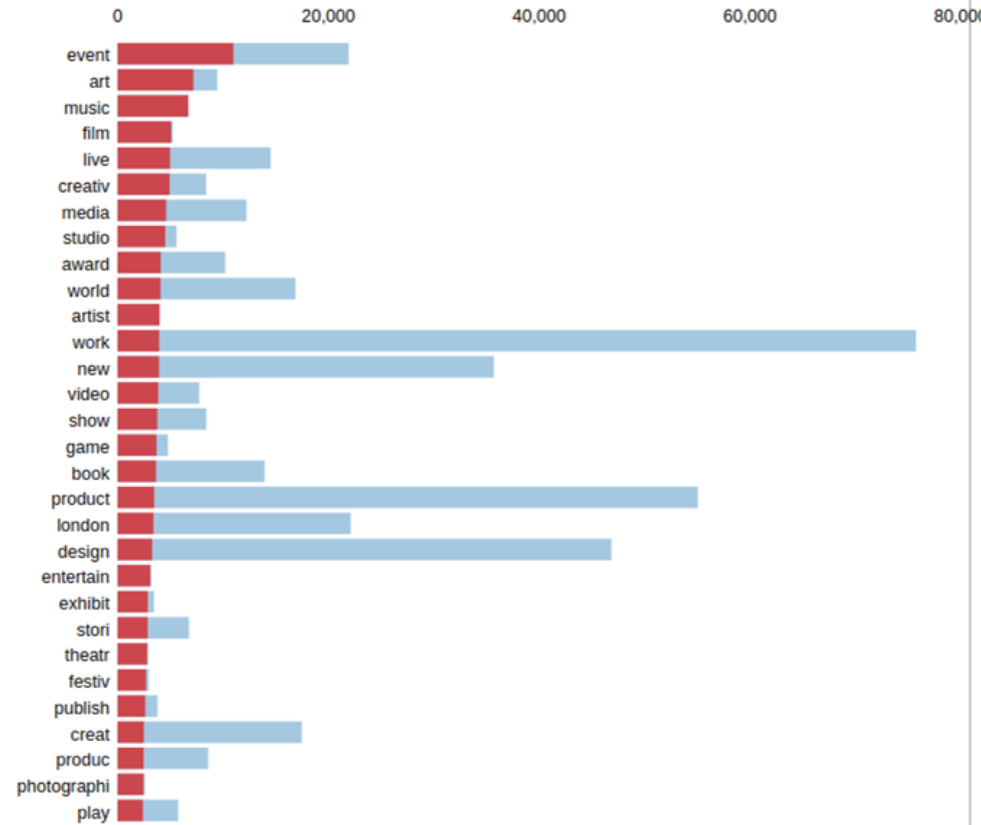


From a corpus of business's descriptions of themselves on the internet:

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 10 (4.1% of tokens)

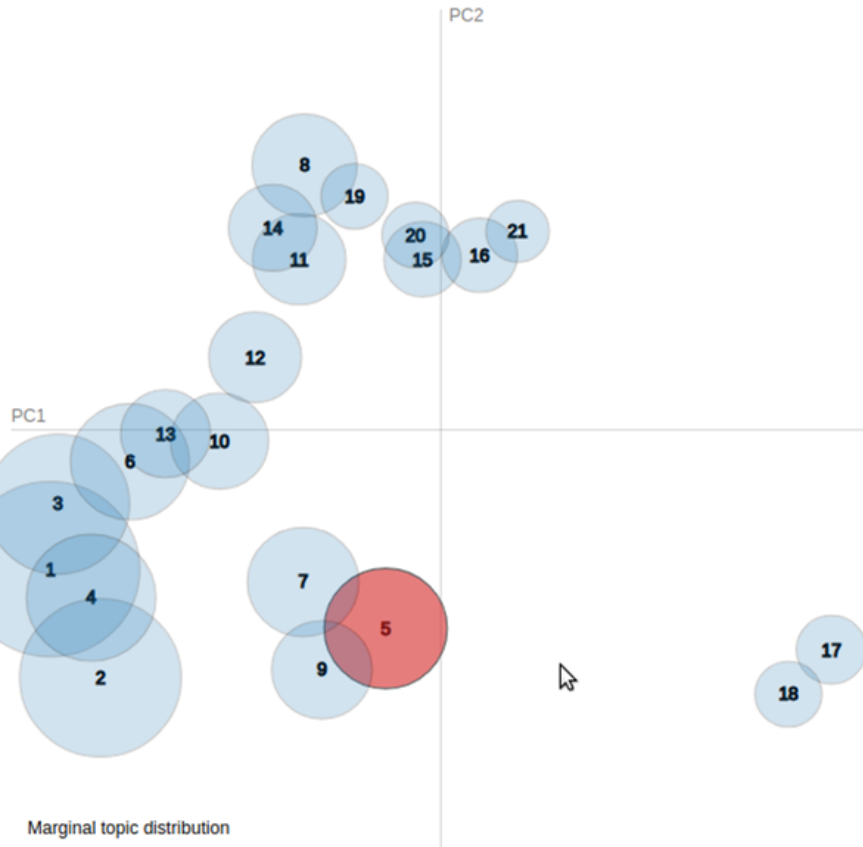


1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

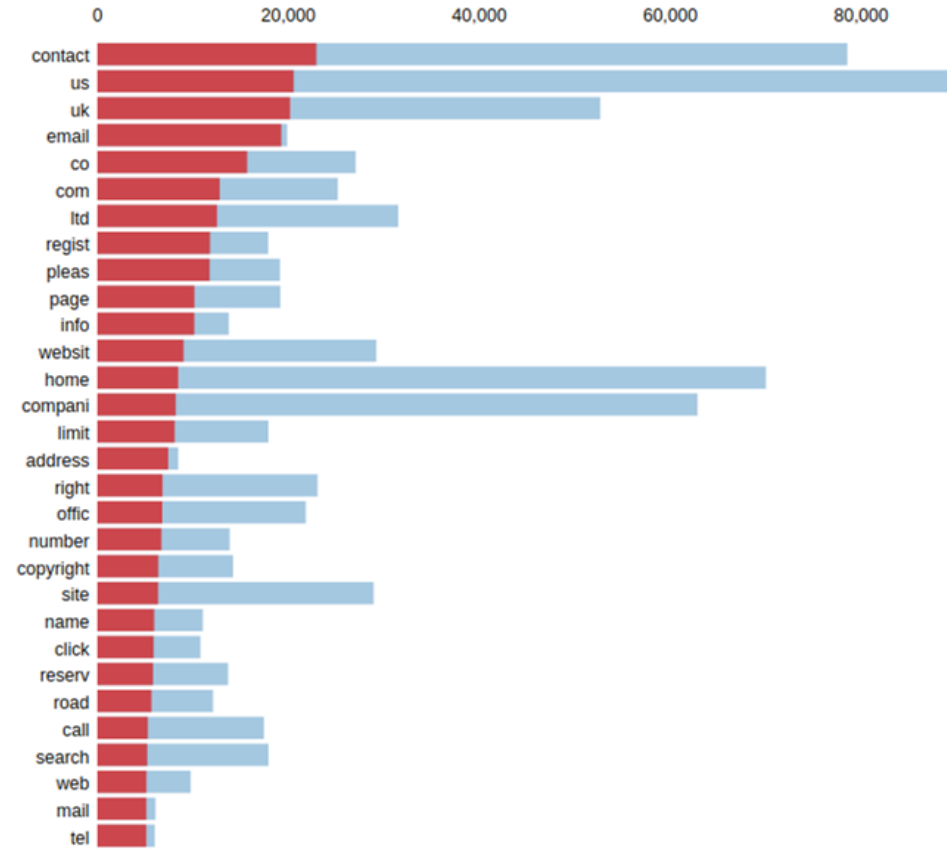
4) Validation – Invalid LDA topics



Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 5 (6.4% of tokens)



Overall term frequency
 Estimated term frequency within the selected topic

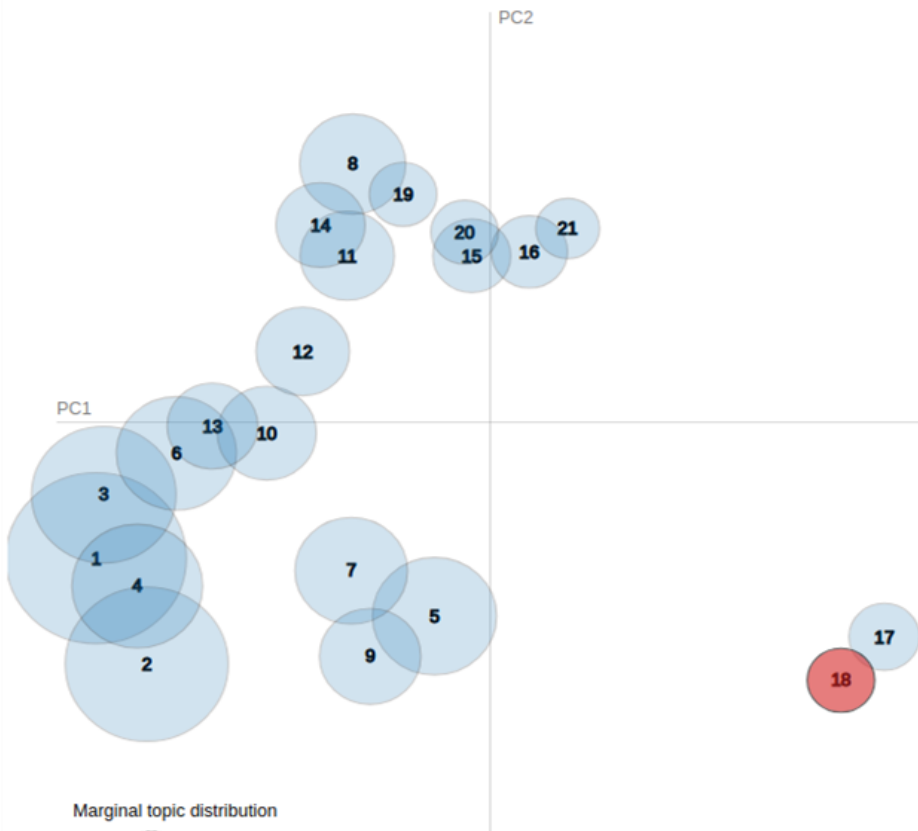
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
 2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

4) Validation – Invalid DATA



Invalid topics WELL SEPARATED in topic/vocabulary space, discriminate automatically?

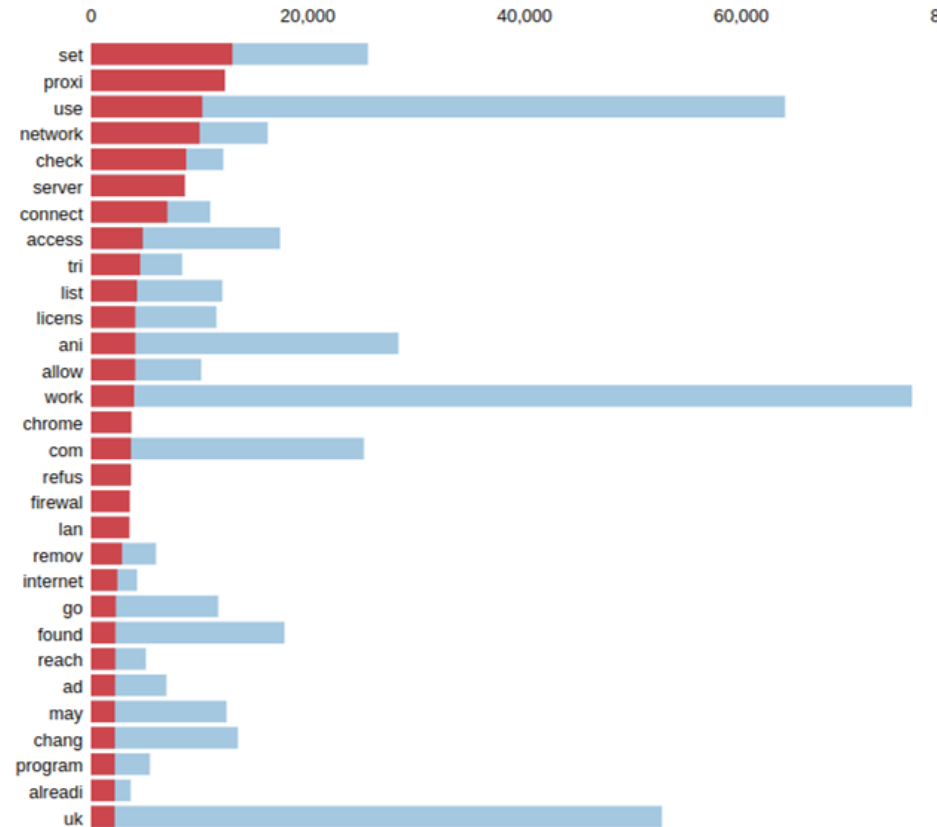
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



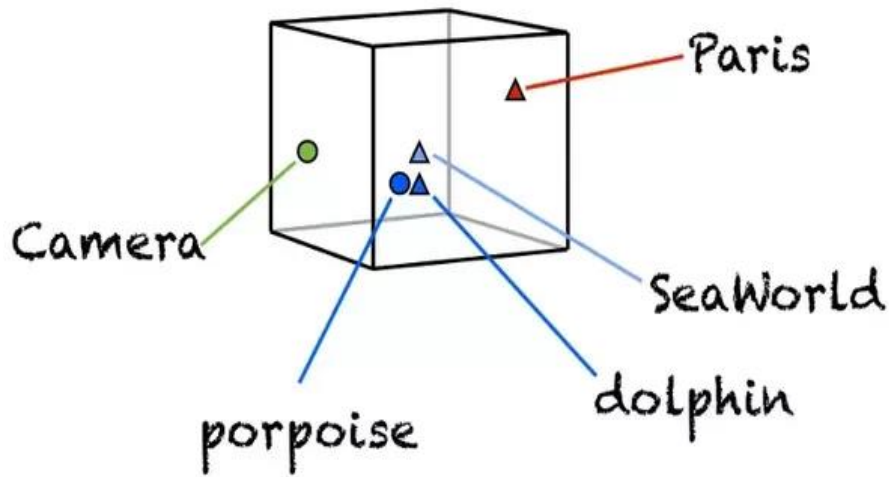
Top-30 Most Relevant Terms for Topic 18 (1.9% of tokens)



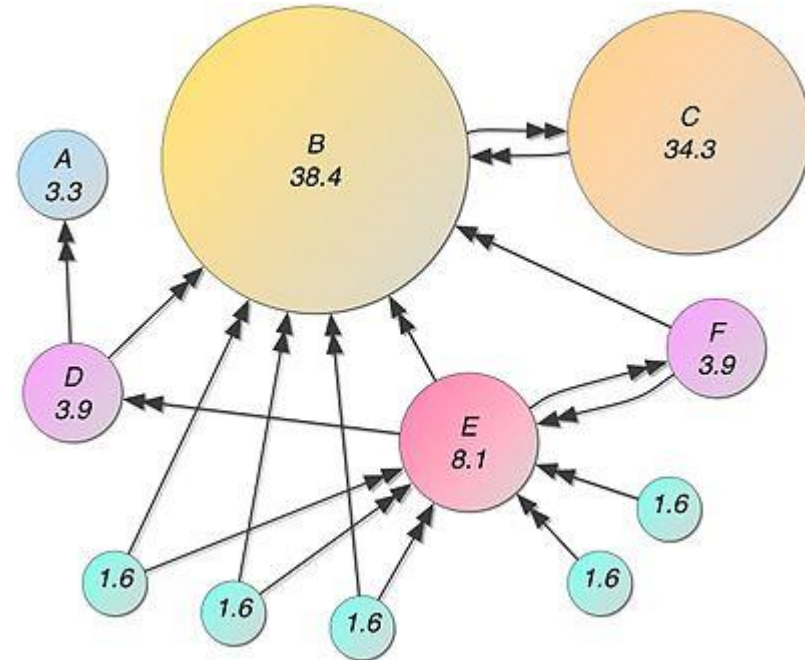
Overall term frequency (blue bar)
Estimated term frequency within the selected topic (red bar)

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)

5) Summarizing Clusters



<https://www.kaggle.com/sbongo/do-pretrained-embeddings-give-you-the-extra-edge>



<https://en.wikipedia.org/wiki/PageRank>

TextRank – variant of Google PageRank Algorithm

Create representation of sentences (summed Word2Vec, Doc2Vec, TF-IDF, edit distance, combo?)

Create giant similarity matrix (cosine distance)

Extract X sentences MOST similar to all other sentences

5) Summarizing Clusters



Chosen text representation = implicit assumptions

Summed Word2Vec
(GloVe,
from Stanford NLP)

'Many researchers predict that such narrow AI work in different individual domains will eventually be incorporated into a machine with artificial general intelligence (AGI), combining most of the narrow skills mentioned in this article and at some point even exceeding human ability in most or all these areas.'

'Some of the learners described below, including Bayesian networks, decision trees, and nearest-neighbor, could theoretically, if given infinite data, time, and memory, learn to approximate any function, including whatever combination of mathematical functions would best describe the entire world.'

Doc2Vec

(Lau & Baldwin,
2016)

'They solve most of their problems using fast, intuitive judgements.'

'Emergent behavior such as this is used by evolutionary algorithms and swarm intelligence.'

'Many learning algorithms use search algorithms based on optimization.'

5) Summarizing Clusters



Some Humans!

'In computer science, artificial intelligence (AI), sometimes called machine intelligence, is intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans and other animals. Computer science defines AI research as the study of "intelligent agents": any device that perceives its environment and takes actions that maximize its chance of successfully achieving its goals.[1] More in detail, Kaplan and Haenlein define AI as “a system’s ability to correctly interpret external data, to learn from such data, and to use those learnings to achieve specific goals and tasks through flexible adaptation”.[2] Colloquially, the term "artificial intelligence" is applied when a machine mimics "cognitive" functions that humans associate with other human minds, such as "learning" and "problem solving".[3] The scope of AI is disputed...'

Alternative - Abstractive Summarization?

Nallapati et al – 2016 – “*Abstractive Text Summarization Using Sequence-to-Sequence RNN’s and Beyond*”

Really hard + “state of the art” results don’t look great (so far)

5) Summarizing Clusters



text	POS	count
US	GPE	51
Donald Trump	PERSON	29
U.S.	GPE	25
Brexit	ORG	17
British	NORP	15
French	NORP	15
Russian	NORP	15
Syria	GPE	14
China	GPE	12

Named entity recognition RSS feeds and SpaCy...

- Prepackaged NN's label
Parts of Speech (POS)
- Pick out common entities
- Distinguish between people,
countries, orgs...
- Pick out adjectives?
Sentiment?

Conclusions

- A lot of work to be done
- Clustering can only be locally optimised
- Concept of “true” topics waiting for us is misleading!
- True abstractive summarization out of reach for now
- Aim to improve extractive to point where humans don't have to open the source documents

Future Work

- HDBSCAN, more robust for sparse datasets?
- Use Doc2Vec features, similarities to “landmark” statements, statistical properties of vectors of descriptive sentences...

Questions?

My email: Martin.Wood5@homeoffice.gov.uk

My git thingy: <https://github.com/ozwaldcavendish>

My map of wine tasting notes:

