

# A Topography of Climate Change Research - Methods

Max Callaghan<sup>1,2</sup>

<sup>1</sup>Mercator Research Institute on Global Commons and Climate Change, Torgauer Straße, 10829 Berlin, Germany

<sup>2</sup>School of Earth and Environment, University of Leeds, Leeds LS2 9JT, United Kingdom

Draft current January 3, 2019

## 1 Data

2 This study reproduces the query developed by [1], which is carried out on the Web of Science core  
3 collection. Though not exhaustive, it gives a good coverage of the literature in major peer-reviewed  
4 journals. Each document is assigned to an assessment period according to the timeline shown in  
5 table 1.

6 We use the references scraped from IPCC assessment reports from [2], and attempt to match  
7 these with the results from the web of science. Table [x] shows the percentage of IPCC citations  
8 matched in each working group for each assessment report.

name	years
AR1	1988-1989
AR2	1990-1994
AR3	1996-2000
AR4	2001-2006
AR5	2007-2013
AR6	2014-

Table 1: Assessment period time windows

## 9 Pre-processing

10 Data quality in earlier Web of Science results is poorer, and some documents have missing abstracts.  
11 In the quantification of the size of the literature and its vocabulary in table [], titles are substituted  
12 for abstracts where they are not available. The words of the documents are lemmatized/stemmed,  
13 replacing different forms of the same word (i.e. word/words) with a single instance. Commonly  
14 occurring words, or “stopwords” are removed, as are all words shorter than 3 characters, and all  
15 words containing only punctuation or numbers.

16 For each period, the documents are transformed into a document-term matrix, each row repre-  
17 sents a document, and each column represents a unique word. Each cell contains the number of that  
18 column’s terms in that document. Only terms which occur more than once are considered.

19 For the calculation of the topic model, documents with missing abstracts are ignored, and the  
20 document term matrix is transformed into a document frequency-inverse document frequency (tf-idf)

21 matrix, where scores are scaled according to the frequency of their occurrence in the corpus. This  
 22 gives more weight to terms which appear in few documents, and less weight to those which appear  
 23 in many.

$$tf(t, d) = f_{t,d}, \quad idf(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|} \quad (1)$$

## 24 Dynamic Non-negative Matrix Factorisation

25 Non-negative Matrix Factorisation (NMF) is an approach to topic modelling which factorises the  
 26 term-frequency-inverse document frequency matrix  $V$  into the matrices  $W$ , the topic-term matrix,  
 27 and  $H$  the document-topic matrix, whose product approximates  $V$ :

$$V_{i\mu} \approx (WH)_{i\mu} = \sum_{a=1}^r W_{ia} H_{a\mu} \quad (2)$$

28 As demonstrated in Figure 1, each topic is represented as a set of word scores, and each document  
 29 a set of topic scores. The combination of the two give the word scores in the document. For clarity  
 30 in the figure, these are shown as simple counts, but in the model these are scaled according to each  
 31 term’s frequency within the corpus.

32 In Dynamic Non-negative Matrix Factorisation, proposed by Greene [3], a separate topic model  
 33 is run for each period. These are then joined through another topic model, which takes the topic-  
 34 term matrices of the all periods as  $V$ , and produces dynamic topics, which describe the window  
 35 topics according to the words which occur in them. Similar topics across and within time periods  
 36 are thereby grouped together.

37 While Greene uses an automatic approach to deciding on topic numbers within time periods, we  
 38 found the number of topics derived from topic coherence scores (as used by Greene) to be noisy, and  
 39 instead opt for identifying (subjectively) an optimal number of topics for each window. We do this  
 40 by comparing topic lists with increasing numbers of topics, where similar topics are automatically  
 41 placed next to each other. We similarly compare different numbers of Dynamic topics.

42 We settle on  $[x, x, x, x, x]$  topics for ARs 1 to 6, and  $[x]$  dynamic topics.

43 Topics are calculated using the scikitlearn library [4]

## 44 Topic Representation and Newness

45 To calculate topic representation in IPCC reports we divide each topic’s share in the subsample of  
 46 documents cited by IPCC reports by its share in the whole corpus.

47 We calculate a dynamic topic’s total score as the sum of document-dynamic topic scores (which in  
 48 turn are made of the product of all document-window topic and window topic-dynamic topic scores).  
 49 A dynamic topic’s window score is the sum of document-dynamic scores considering only documents  
 50 in the given time window. To represent a dynamic topic’s newness, we multiply each assessment  
 51 period number by the share of it’s total score occurring in that window, and take the mean of these  
 52 scores. A topic in which 100% of documents which make it up occurred in assessment period 1 (6)  
 53 would thereby receive a score of 1 (6), while a topic evenly distributed across all assessment periods  
 54 would receive a score of 3.5.

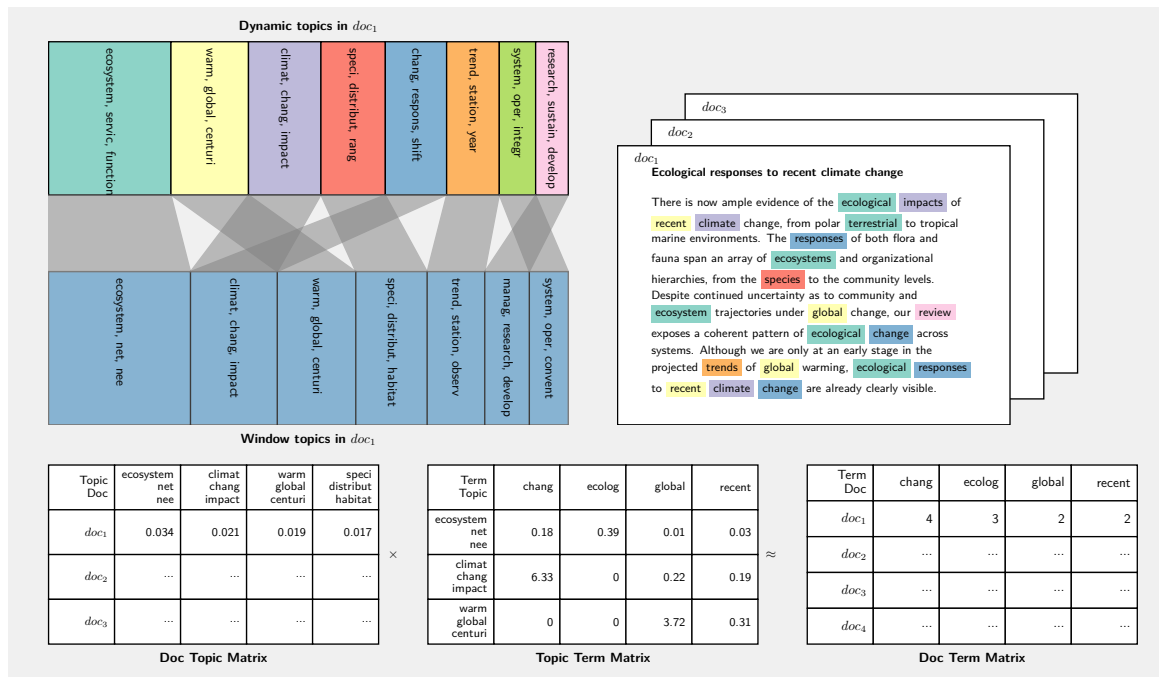


Figure 1: SI Topic make up of a single document

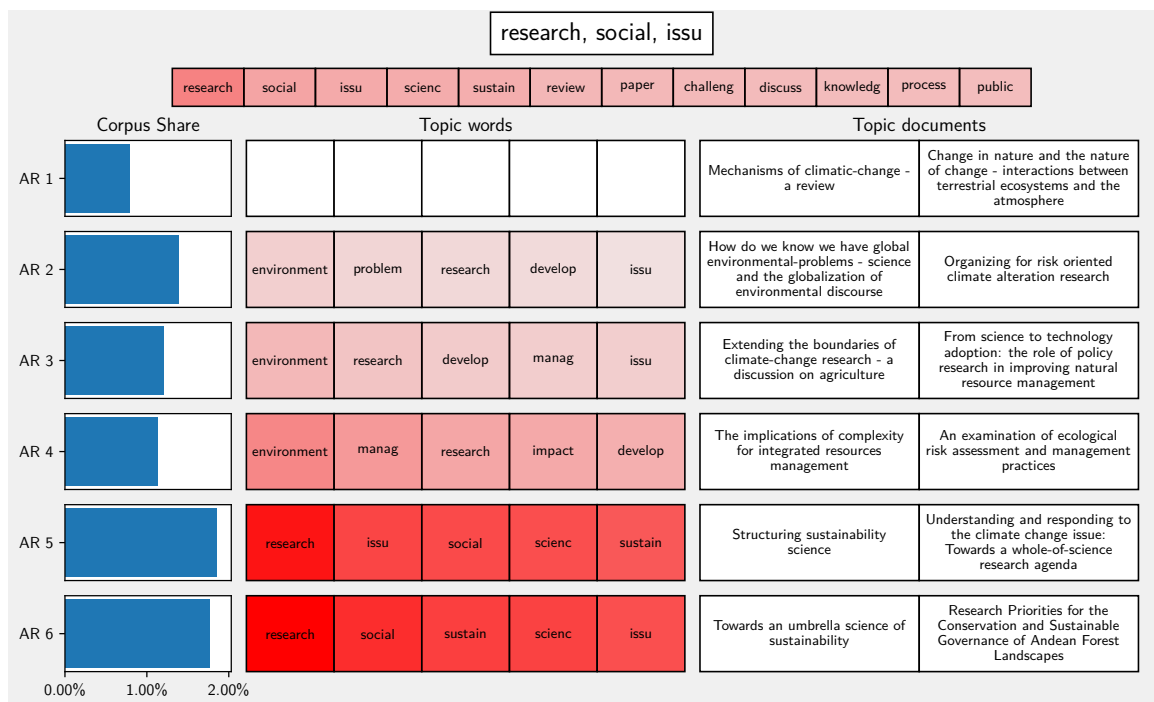


Figure 2: Word and document development of the “Research” dynamic topic

## References

- [1] Michael Grieneisen and Minghua Zhang. The Current Status of Climate Change Research. *Nature Climate Change*, 1:72–73, 2011.
- [2] Jan C. Minx, Max Callaghan, William F. Lamb, Jennifer Garard, and Ottmar Edenhofer. Learning about climate change solutions in the IPCC and beyond. *Environmental Science & Policy*, 2017.
- [3] Derek Greene and James P Cross. Exploring the Political Agenda of the European Parliament Using a Dynamic Topic Modeling Approach. pages 1–47, 2016.
- [4] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Mattheiu Perrot, and Édouard Duchesnay. Scikit-learn: Machine Learning in Python Fabian. *Journal of Machine Learning Research*, 12:2825–2830, 2011.