

Imperial College London

Department of Life Sciences, Faculty of Natural Sciences

The Analysis of Topic Bias in Four Academic Funding Agencies using Machine-Learning Approaches

Author:

Jingkai SUN

Date:

26 August 2021

CID:

01991822

Word Count:

5882

A thesis submitted for the partial fulfillment of the requirements for the degree of Master of Science at Imperial College London

Submitted for the MSc in Computational Methods in Ecology and Evolution

Declaration

I can declare that I am the only person who is responsible for data collection, data cleaning, the development of methodology for the topic analysis, text mining and machine learning models used in this project. I herewith certify that all works in my project are done by myself. David Orme provided solutions to some statistical problems, James Rosindell shared the idea of using topic modelling techniques with me and provided some scientific papers related to my project. He also gave feedback on the draft of my project. Samraat Pawar provided supports for the technical questions in terms of Python, R and Latex and also offered useful suggestions for my draft of the project.

Acknowledgement

Firstly, I would like to thank my supervisors for solving all my questions during the weekly meeting. David Orme offered me an opportunity to study this interesting topic and many supports on statistical analysis in the project. Samraat Pawar provided lots of computing modules and solutions to my technical questions so that I have the ability to process the data and create machine-learning models using Python, Shell and R. James Rosindell provided me numerous ideas from the start of the project and suggestions for assigning topic labels in my topic model. Additionally, he arranged multiple meetings for providing me many useful feedbacks based on the draft of the project. In addition, I also would like to thank Yewshen Lin, who is one of my colleagues, for sharing his suggestions and experiences of using text-mining techniques. Moreover, a big thank you to Marina Papadopoulou, a student who graduated with an MSc in Computational Methods in Ecology and Evolution at Imperial College London, for the methods of the selection and assessment of topic models she introduced in the dissertation. Last but not least, many thanks to everyone in James' lab meeting for listening to my presentation and giving suggestions to me.

Contents

List of Figures	4
1 Abstract	5
2 Introduction	6
3 Methods	8
3.1 Experimental Data	8
3.2 Latent Dirichlet Allocation	9
3.3 Coherence Measurement	10
3.4 Model Selection	11
3.5 Prediction of Funding Amount	12
3.6 Computing Tools	13
3.7 Data Availability	14
4 Results	14
4.1 Topic Selection	14
4.2 LDA Model Fitting and Topic Visualisation	16
4.3 Topic Analysis	18
4.4 Prediction of Funding Amount	22
5 Discussion	22
6 Conclusion	24
Support Information	26
References	37

List of Figures

1	The result of the coherence measurement	14
2	Mean similarity and stability between models	15
3	Interactive visualisation for topic distributions	16
4	Probability distributions of all topics from all documents	17
5	Eight wordclouds of top 25 stemmed words	18
6	Heatmap between topic and agencies	19
7	Heatmaps between topics and 50 intervals of discrete funding amount	20
8	The accuracy of prediction of funding amount using LightGBM	22
9	Visualisation of topic 1 by pyLDAvis	30
10	Visualisation of topic 2 by pyLDAvis	30
11	Visualisation of topic 3 by pyLDAvis	31
12	Visualisation of topic 4 by pyLDAvis	31
13	Visualisation of topic 5 by pyLDAvis	32
14	Visualisation of topic 6 by pyLDAvis	32
15	Visualisation of topic 7 by pyLDAvis	33
16	Visualisation of topic 8 by pyLDAvis	33
17	The probability distribution of topic probabilities after cut small probabilities . .	34
18	t-SNE clustering for visualisation of eight topics	34
19	Word count and importance of topic keywords	35
20	Sentences coloring of dominant topic in part of documents	36

¹ 1 Abstract

² Academic funding plays an essential role in research projects. Due to the finite amount of
³ academic funding and an increasing number of emerging research fields, funding agencies could
⁴ provide more funding to some specific topics, resulting in a topic bias in funding agencies. Com-
⁵ bining topic modelling with a supervised machine-learning algorithm, I examined the topic bias
⁶ from the information of funded projects in the National Institutes of Health (NIH), European
⁷ Research Council (ERC), National Science Foundation (NSF) and UK Research and Innova-
⁸ tion (UKRI). Then, I used a supervised classifier to predict the amount intervals of funding,
⁹ with topic distributions generated from topic models and project duration as predictors. I col-
¹⁰ lected 145,787 funded project data between 2015 and 2021 from 4 academic funding agencies
¹¹ and extracted topic distributions using latent Dirichlet allocation. Then, I evaluated the topic
¹² distributions in each agency and each amount interval of funding per day in project duration.
¹³ Afterwards, I created Light Gradient Boosting Machine to classify three intervals of funding
¹⁴ amounts (i.e low-, medium- and high-amount classes) and made comparisons to models with
¹⁵ different predictors in different agencies. The results show that (1) topic modelling is helpful
¹⁶ for topic assignment of project abstracts; (2) The topic “Computer Systems and Commercial
¹⁷ Application” dominated NSF, ERC and UKRI agencies, whereas the topic “Cancer Treatment
¹⁸ and Immunology” dominated NIH agency; (3) With combining the topics “Computer Systems
¹⁹ and Commercial Application” with “Cancer Treatment and Immunology” (e.g. computational
²⁰ biology or medical engineering), a project could have a higher probability to obtain more fund-
²¹ ing in all agencies; (4) Topic distributions are found to improve the accuracy of prediction of
²² the funding amount by 5% – 20% once I combined them with project duration as predictors
²³ together, even though the model fails to have sufficient prediction power when the model solely
²⁴ considered topic distributions as predictors.

²⁵ **Keywords:** Text Mining, Bibliometrics, latent Dirichlet allocation, Multiclassification, Fund-
²⁶ ing agencies, Academic funding, National Institutes of Health, European Research Council,
²⁷ National Science Foundation, UK Research and Innovation, LightGBM, Machine Learning

28 The Analysis of Topic Bias in Four Academic Funding Agencies
29 using Machine-Learning Approaches

30 Jingkai Sun (CID: 01991822)

31 August 26, 2021

32 **2 Introduction**

33 The advancement of academic research is not only dependent on the advancement of technology,
34 but also on funding support. Academic funding can provide essential support for scientists,
35 especially for medicine and natural science studies [1]. However, with an increase in the number
36 of researchers, scientists and emerging research fields all over the world [2], the competition of the
37 academic market is increasingly fierce so that funding fail to cover all excellent research projects.
38 This indicates a decrease in opportunities for scientists to obtain enough funding to complete
39 their research projects, many scientists worry about their research careers [3]. Therefore, it is
40 important to understand determinants of funding success for research projects and how they
41 will be biased due to different reasons.

42 In 2011, Allesina *et al.* studied whether nepotism can affect scientific success in Italian Academia
43 using conventional statistical analysis [4]. They leveraged the Monte Carlo approach to show
44 there are fewer different last names of principal investigators (PIs) compared with the expected
45 situation at random for each discipline in Italian research institutions. This revealed the exis-
46 tence of nepotism in most disciplines, even though the results are underestimated due to the
47 exclusion of cases of nepotism of the mother-child type. Furthermore, they also used logistic
48 regression models to examine how the probability of sharing the last name changes with factors
49 such as geography, institutions, sub-disciplines and latitude, concluding a significant latitudinal
50 influence on nepotism.

51 However, at that time, in most of the research that studies the factors of scientific success such
52 as what Allesina *et al.* did, they used traditional statistical analysis to uncover the factors
53 that determine the funding success. In recent years, with the development of machine-learning
54 techniques and the computational power of computers, big data are more influential for us in

recent decades than in the past [5]. It is possible to find the relationship between things and hence predict what will happen in the future through big data in many research fields. For example, Mahajan *et al.* utilised the Naïve Bayes model with a large amount of data in the stock market to predict the movement of the stock market and conclude that the model is useful for predicting the future prices of shares [6].

Therefore, more researchers utilised machine-learning and text mining techniques to extract information and patterns from numerous text data such as scientific papers. In 2014, Emre *et al.* analysed 100,000 publications from Computer Science to examine the centrality in the coauthorship network and use it with a random forest algorithm to predict whether a paper can be highly cited five years after publication [7]. The final outcome showed a high accuracy of prediction. Additionally, Michael *et al.* applied several machine-learning classifiers to predict the funding success of projects from a crowdfunding website – Kickstarter. They used the information of projects such as project duration, the number of Twitter followers and the number of Facebook connections *etc.* as features, obtaining 68% of the accuracy of prediction [8]. Instead of using the basic information of projects exclusively, Tanushree and Eric (2014) extracted phrases in the project description by text-mining techniques from Kickstarter and explored whether they are the factors determining the funding success [9].

To explore the impact factors of academic funding success using text-mining approaches, Magua *et al.* (2017) studied the problem of gender bias in the National Institutes of Health (NIH), an academic funding agency. They evaluated the impact on the funding process caused by gender differences in NIH by combining Latent Dirichlet Allocation (LDA) and qualitative analysis [1]. They first used the regression models to examine the sex differences by assigned scores for male (PIs) and female PIs. The outcome suggested, compared with female PIs, that male PIs' applications have higher priority, approach and significance scores with having the same experience level and research quality. With latent Dirichlet allocation (LDA), they also found that funded male PIs are usually described as “pioneers” in specific fields with “highly innovative” and “significant research”, whereas funded female PIs are described as “expertise” or “excellent”. The outcome revealed a distinct criterion that may lead to different funding success rates of funding applications between male and female PIs.

The difference of research topics of the projects is also another one of the most important impact factors of funding success [10]. By finding the ‘hot’ topics that can be more funded, researchers can study their own fields by combining these ‘hot’ topics to increase the probability of obtaining more funding. However, using big data approaches to study the relationship of topics and academic funding are not found yet.

Thus, the objective of this project is to explore the relationship between topics and academic

90 funding in the National Institutes of Health (NIH), National Science Foundation (NSF), Eu-
91 ropean Research Council (ERC) and UK Research and Innovation (UKRI), which are four
92 well-known academic funding agencies in the world. At the end of the project, I try to answer
93 the following questions:

- 94 • Is text-mining techniques useful for finding ‘hot’ topics in different agencies?
95 • Whether different topics could affect the funding amount in four different academic funding
96 agencies?
97 • What are some topics that could improve the probability of obtaining more funding in
98 four agencies?
99 • Are academic funding can be predicted using Machine Learning Approach?
100 • Can topics distributions improve the accuracy of prediction of academic funding amount?

101 **3 Methods**

102 **3.1 Experimental Data**

103 Data was collected from the websites of UKRI, NIH, NSF and ERC. After removing projects of
104 which abstracts were unavailable, data from 145,787 funded projects awarded between 2015 and
105 2021 were collected, including 27,035 project abstracts from the UKRI funding agency collected
106 by the GtR API program [11], 41,796 abstracts from the NSF database of funded projects [12],
107 71,008 abstracts from NIH collected using RePORTER database and finally 5,948 abstracts
108 from ERC collected directly using its website database [13, 14].

109 To obtain more accurate results, project abstracts need to be preprocessed in ways that are
110 required by any text mining technique. The first step was to remove all HTML character
111 references (e.g. & or <), which is a way to display special characters on a website page
112 since some of the abstracts in my dataset are scraped from websites. Then, I removed all
113 punctuations, non-English words and English stop words such as “You”, “I”, “necessary” and
114 “on” etc. Some of the stop words are from the NLTK package in Python and some are identified
115 manually (see Support Information). In addition, to improve the accuracy of the LDA model,
116 bigram and trigram models were utilised to identify two words or three words that have a high
117 probability to occur together [15]. For example, if the word set {chimeric, antigen, receptor}
118 could have a higher probability to show up together in the genetic- or protein-related topics
119 than physics-related topics, then the trigram will treat them as a single word with an underscore
120 (e.g. chimeric_antigen_receptor).

121 Also, these co-occurring words are helpful for describing topics produced by the LDA model,
122 even though a few co-occurring words are not ideal such as ‘erc_start’ [16]. Afterwards, all
123 remaining words were stemmed, meaning that all inflected words were cut back to a common
124 root [17]. For instance, all words like “damage”, “damaged” and “damaging” were transformed
125 into “damag”. Next, all high-frequency (occur in more than 80% of abstracts) and low-frequency
126 (occur in less than 0.1% of abstracts) words were removed. Finally, 7060 unique words are left
127 to fit the topic models.

128 The final dataset consists of project title, project abstracts, funding amount of projects, the
129 project duration of projects. Because four funding agencies are not founded in the same country,
130 meaning the currency are distinct from each other. For consistency, the currency of funding
131 amount in four agencies is converted to U.S. Dollars (\$) with real-time exchange rates using the
132 Fixer API program [18].

133 3.2 Latent Dirichlet Allocation

134 Latent Dirichlet Allocation (LDA) is one of the most popular models in topic analysis, which is
135 an unsupervised machine learning technique. First, I introduce D as the number of documents,
136 V as the number of words in corpus and K as the number of topics. LDA treats D documents
137 in the dataset as a probability distribution of K topics. LDA is used to infer K latent topics
138 that can summarise V words in text documents [16]. From this perspective, the LDA model can
139 be also regarded as a dimensionality reduction tool, transforming V words into K clusters (i.e.
140 topics). Two matrices are inferred in the LDA model for obtaining topic distributions for D
141 documents [19]. The first matrix is a document-topic matrix with D rows of documents and K
142 columns of topics, the values in the matrix are probabilities of each document belonging to each
143 topic. Hence, the sum of probabilities of K topics for each document should be one. Similarly,
144 the topic-word matrix enables K topics to have words with different weights representing how
145 important each word is. For instance, if the high weights are generated by the LDA model for
146 some words such as ‘climate’, ‘emission’, ‘change’, ‘environment’ and ‘temperature’, then this
147 might be a topic related to “Climate Change” or “Environment Protection”.

148 In this project, the LDA model is implemented by collapsed Gibbs sampling (see Support
149 Information section). The collapsed Gibbs sampling method has the following basic process:
150 (1) Select a document and randomly assign topics for each word in the document. (2) Repeat
151 the same process for every document in the corpus and then summarise local and global (i.e. for
152 all documents) statistics for each topic. For example, if a document has words such as ‘epilepsy’,
153 ‘dynamic’, ‘bayesian’, ‘EEG’ and ‘model’, one first assign topics for each word in the document
154 and summarises the number of words each topic included in the document (local statistics) and

155 the number of times for each word grouped into each topic in all documents (global statistics).
 156 (3) Topics are reassigned randomly again and the probability of a new assignment is estimated.
 157 (4) The same process is iterated throughout the words and documents. The probability of a
 158 new assignment can be estimated by [20]:

$$P = \frac{n_{ik} + \alpha}{N_i - 1 + K\alpha} \frac{m_{t,k} + \beta}{\sum_{w \in V} m_{t,k} + V\beta} \quad (1)$$

159 where n_{ik} is the number of current assignments to topic k in doc i , N_i is the number of words
 160 in document i , α is a hyperparameter that controls the prior distribution over topic weights in
 161 each document, K represents the number of topics, $m_{t,k}$ means the number of assignments of
 162 given word t to topic k , V is the size of vocabulary and β is a hyperparameter for the prior
 163 distribution over word weights in each topic. Tuning two hyperparameters also affect inferences
 164 of the LDA model. Higher α can make document preference “smoother” over topics, enabling
 165 the model to generate more topics with higher probability; higher β also enables more terms to
 166 be important for topic description and vice versa, i.e. topics preferences are “smoother” over
 167 their keywords. Therefore, for simplicity, my research used default values of α and β in Gensim
 168 [21], which both are one over the number of topics. Additionally, because collapsed Gibbs
 169 sampling is based on Markov Chain Monte Carlo (MCMC) algorithm, leading to the stochastic
 170 output when the same LDA model with fixed parameters is run for multiple times, explaining
 171 why I need to conduct a *stability* test for multiple topic models in the following section [22, 23].

172 3.3 Coherence Measurement

173 The coherence measurement is an approach to examine how the words in the same topic support
 174 each other [24]. Many metrics have been developed to measure how coherent a set of words in
 175 each topic such as C_{uci} [25], U_{mass} [26], C_{npmi} [27], C_v [28] and C_p [29]. In this project, the C_v
 176 measure is applied since it has been tested that it outperformed all other metrics in coherence
 177 measurement [28]. C_v is based on a sliding window [30], which uses one-set segmentation (i.e.
 178 any single word within a set is paired with every other single word for comparison) between
 179 topic words as well as an indirect confirmation measure [24] to calculate normalized pointwise
 180 mutual information (NPMI) for each word pair [28]. Afterwards, the cosine similarity for each
 181 NPMI of word pairs is calculated [30]. The equation of NPMI can be expressed as:

$$NPMI(w_i, w_j)^\gamma = \left(\frac{PMI(w_i, w_j)}{-\log(P(w_i, w_j) + \epsilon)} \right)^\gamma \quad (2)$$

182 where

$$PMI(w_i, w_j) = \log \frac{P(w_i, w_j) + \epsilon}{P(w_i) \cdot P(w_j)} \quad (3)$$

183 where $P(w_i, w_j)$ is the probability of words w_i and w_j co-occur, ϵ is a parameter that can avoid
184 the logarithm of zero, γ is a parameter that can give more weights on higher NPMI values as γ
185 increases.

186 3.4 Model Selection

187 To obtain an efficient topic model, I first run topic models with a given number of topics k
188 with 10,000 iterations [31], where each model M^k are repeated by 10 times and each repetition
189 has a single C_v value. Therefore, I obtained 10 different coherence values for each model M^k .
190 Mathematically speaking, each model M^k has a vector of coherence values, expressed as:

$$\begin{bmatrix} Coh_1^k & Coh_2^k & \dots & Coh_n^k \end{bmatrix} \quad (4)$$

191 where $n = 10$ and Coh_i^k means the i -th coherence measurement (i.e. C_v value) for a k-topic
192 model (i.e. the topic model with k topics).

193 In this project, I tested 32 topic models whose the k starts from 2 to 64 with increment 2 (i.e.
194 $k \in \{2, 4, 6, \dots, 64\}$). Therefore, I gained 32 vectors of coherence values, expressed as a $(p \times n)$
195 coherence matrix:

$$\begin{bmatrix} Coh_{11}^{k=2} & Coh_{21}^{k=2} & \dots & Coh_{n1}^{k=2} \\ Coh_{12}^{k=4} & Coh_{22}^{k=4} & \dots & Coh_{n2}^{k=4} \\ \vdots & \vdots & \ddots & \vdots \\ Coh_{1p}^{k=64} & Coh_{2p}^{k=64} & \dots & Coh_{np}^{k=64} \end{bmatrix} \quad (5)$$

196 where n is the number of coherence values in each k-topic model, p is the number of k-topic
197 models.

198 Afterwards, I computed the median C_v (Cv_{med}^k) values for each row in the coherence matrix,

199 transforming the coherence matrix into a ($p \times 1$) vector, shown as follows:

$$\begin{bmatrix} Cv_{med}^{k=2} \\ Cv_{med}^{k=4} \\ \vdots \\ Cv_{med}^{k=64} \end{bmatrix} \quad (6)$$

200 where each element in the vector represents the median C_v value for a k-topic model M^k . Then
201 I choose a subset of M^k s with relatively high Cv_{med}^k and make a comparison with each other
202 using *stability* and mean *similarity* approach (see Support Information) to find the final topic
203 model.

204 3.5 Prediction of Funding Amount

205 To find out whether the topic distributions I obtained using the LDA model have prediction
206 power or improved prediction compared to the project duration, I used a supervised machine-
207 learning tool to forecast the discrete funding amount of projects. First, I divided the four
208 agencies into two groups: the NIH group (the health-specialist agency) and all the other agencies
209 including NSF, ERC and UKRI. Two targets were predicted, which are the total funding amount
210 for each project and the funding amount per day, called daily funding amount in this project,
211 which is calculated by:

$$A_{daily} = \frac{A_{total}}{Duration} \quad (7)$$

212 where A_{daily} is the daily funding amount, A_{total} is the total funding amount, $Duration$ is the
213 project duration in the unit of days.

214 Then, I treat the prediction as a multiclassification problem. Thus, all funding amounts are dis-
215 cretized to three different classes (i.e. low-amount zone, medium-amount zone and high-amount
216 zone). For the multiclassification problem, I utilised a popular gradient boosting framework
217 known as Light Gradient Boosting Machine (LightGBM) [32]. LightGBM is based on the gra-
218 dient boosting decision tree (GBDT), which is a prevalent machine-learning algorithm [33] that
219 can be used to solve multiclassification problems [34] and learning to rank [35]. However, differ-
220 ent from GBDT, LightGBM leverages approaches of Gradient-based One-Side Sampling (GOSS)
221 and Exclusive Feature Bundling (EFB). GOSS method retains only data instances with larger
222 gradients. It can provide more information gain, obtaining more accurate outcomes. Meanwhile,
223 the EFB approach can optimize the efficiency of feature selection, enabling the LightGBM algo-

224 rithm to handle data with less computational complexity [32]. The LightGBM has been proved
225 that it can outperform XGBoost and SGB algorithms concerning the efficiency of computation
226 of memory consumption [32]. Therefore, the prediction was done by solely LightGBM.

227 I created separate LightGBM models according to datasets of two groups (NIH and all other
228 agencies). For each group, I built three models with distinct predicting variables. The first one is
229 the LightGBM model with all topic distributions as predicting variables (defined as M_{topic}), the
230 second model is created with only project duration as a feature, defined as $M_{duration}$. The third
231 model has all topic distributions and project duration as predicting variables, defined as M_{all} .
232 All predicting variables used for the three models are normalised by using the “StandardScaler”
233 function in the “sklearn” package in Python [36]. Finally, these two groups comprised 70,842 and
234 74,141 observations respectively, with a balanced number of classes and six models in total (i.e.
235 three models for each group). 70% of observations are used as the training set, the remaining
236 observations are used as the test set for prediction in both two groups.

237 There are three steps to make a prediction: (1) I fitted all models described above by using
238 normalised predictors. (2) I chose parameters for each model, where the learning rate is 0.02,
239 boosting type is GBDT (i.e Gradient Boosting Decision Tree), the maximum depth of the tree
240 is 80, the number of leaves is 200 and the number of iterations is 150. The parameters of all
241 models are the same. (3) The probabilities of three classes for each data point were gained. I
242 then chose a class with the highest probability as a label for each observation. The accuracy of
243 each model was calculated by the expression:

$$Accuracy = \frac{tp}{tp + fp} \quad (8)$$

244 where tp is the number of true positives and fp is the number of false positives, indicating the
245 ability of the model not to assign positives to the data that is negative.

246 3.6 Computing Tools

247 In this project, all scraped data are stored in a MySQL database. ALL analytical techniques
248 are conducted using Python 3, where the “pymysql” package is used for extracting data from
249 MySQL database, packages “numpy”, “pandas”, “NLTK”, “pyinflect” and “sklearn” are used
250 for data preprocessing, the LDA model used in the project is implemented by the “Gensim”
251 package, the “pyLDAvis” package is used for model visualisation.

252 **3.7 Data Availability**

253 All scripts of data preprocessing, text-mining and machine-learning modelling are available at
254 https://github.com/Kyle-J-Sun/CMEE_Final_Project

255 **4 Results**

256 **4.1 Topic Selection**

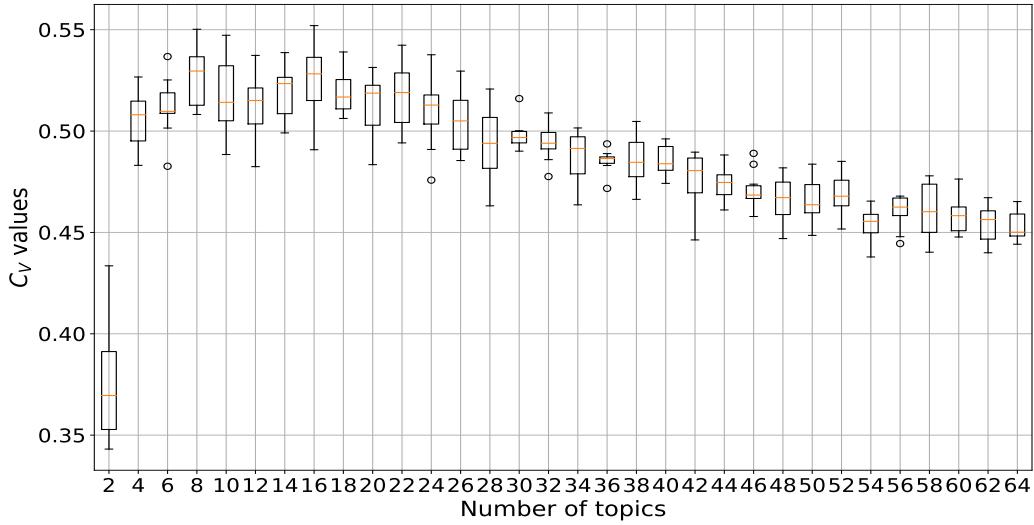


Figure 1: **Coherence measurement:** The figure shows the result of the coherence measure (C_v) for topic models with 2 to 64 topics with increment by 2, every observation is the boxplot of coherence vector of each k-topic model. The orange line of every boxplot represents the median C_v value of each 10-run LDA model with given topics. Higher median C_v values are preferred.

257 Fig 1 shows an increase in C_v values from 2 and 8 topics, while a decreasing trend occurred from
258 22 to 64 topics. Therefore, the topic models with 8 to 22 topics (8 topic models in total) were
259 selected for the *stability* test since their median C_v values are significantly higher than others,
260 where the 8-topic models have maximum median C_v value (Median $C_v = 0.53$) and 10-topic
261 models have minimum median C_v value (Median $C_v = 0.51$).

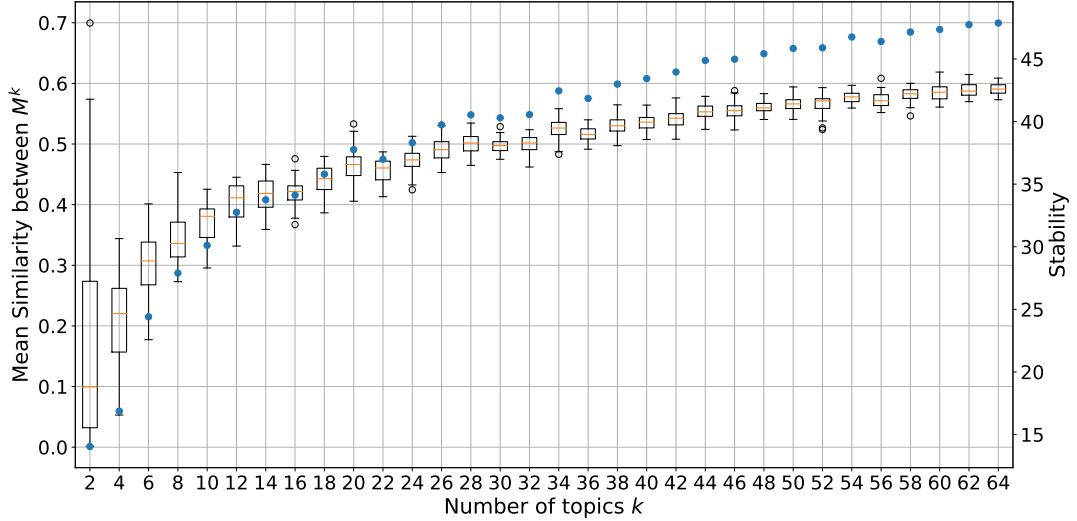


Figure 2: Mean similarity and stability between models: The main y-axis is the averaged *similarity*. Each boxplot is a vector of averaged *similarity* values. The secondary y-axis represents *stability* values, which are blue points in the figure. The best model should have the lowest mean *similarity* and *stability* values (see Support Information section) in the subset of k -topic models selected by C_v metrics.

Furthermore, the averaged *similarity* and model *stability* both show an increasing trend from topic 2 to topic 64 in general. The 8-topic model has the lowest averaged *similarity* ($Sim(M^k) = 0.34$) and *stability* ($Stab(M^k) = 27.9$) between eight models selected by coherence measurement. In conclusion, by figures above, the 8-topic LDA model has the highest median C_v value while its median mean *similarity* and *stability* values are also lowest compared to other selected topic models. Therefore, the 8-topic LDA model is considered as the final model throughout the analysis.

269 4.2 LDA Model Fitting and Topic Visualisation

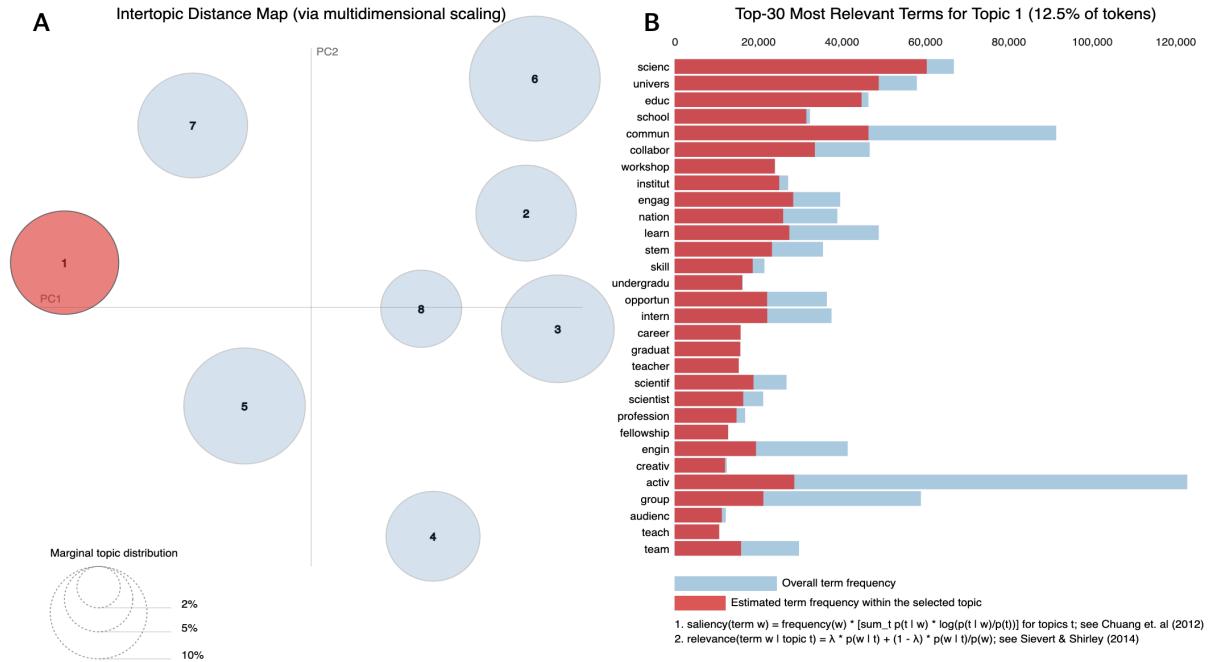


Figure 3: **Topic Visualisation:** An example of interactive visualisation for topic 1 produced from the 8-topic LDA model. Figure A projects eight topics onto a two-dimensional plane by Principal Component Analysis (PCA), where centres of topic circles represent the distance between topics. The areas of circles represent the frequency of topics. Figure B provides a bar chart, where red bars are the frequency of the most relevant stemmed words in topic 1, light blue bars show the corpus-wide word frequency of the individual topic words.

270 For the interpretability of topics, λ is adjusted to 0.61 (see Support Information section). From
 271 Fig.3 A, there is no overlap between the areas of 8 topics, indicating that 8 topics are well-
 272 distributed on the two-dimensional plane, supporting that the topic model with 8 topics could
 273 be the best-fit model in my dataset. For Fig.3 B, the most five relevant words are ‘science’,
 274 ‘university’, ‘education’, ‘school’ and ‘collaboration’ respectively. Thus, the first topic might
 275 be relevant with “Higher Education”. I used the same logic to assign labels for all eight top-
 276 ics. The labels of eight topics are: “Higher Education”, “Brain Science”, “Protein Chemistry
 277 and Microbiology”, “Material Science and Physics”, “Computer System and Commercial Ap-
 278 plication”, “Cancer Treatment and Immunology”, “Social Policy and Public Healthcare” and
 279 “Climate Change and Environment”.

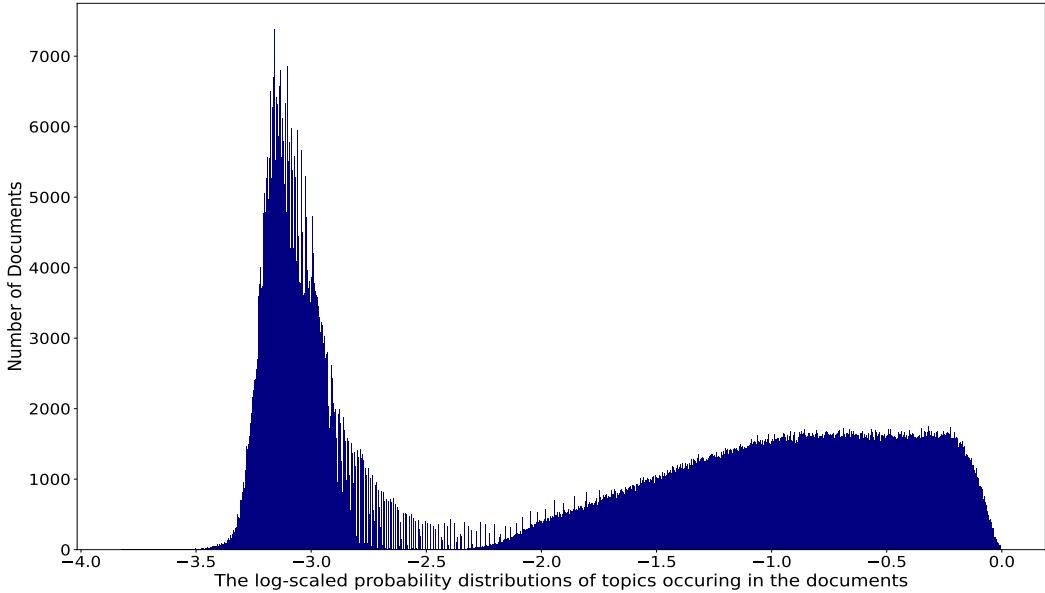


Figure 4: **Topic Distribution:** The figure shows log-scaled (with base 10) probability distributions of all topics from all documents.

280 The probability distribution of topics in each document might not always be helpful for topic
 281 identification since all topics can have non-zero probabilities occurring in each document [19],
 282 indicating that the probabilities for some topics are too small to be informative. Therefore, I
 283 plotted the topic probabilities throughout documents based on base 10 logarithm. Two peaks
 284 can be seen in Fig. 3, where one occurred as the probabilities are between $10^{-3.5}$ and $10^{-2.5}$,
 285 another peak is distributed from approximately 10^{-1} to $10^{-0.2}$. Hence, in order to remove all
 286 uninformative topics for each topic, any topics of which the probability is less than the minima
 287 between two peaks, at $10^{-2.3}$ (i.e. 0.005), are ignored.

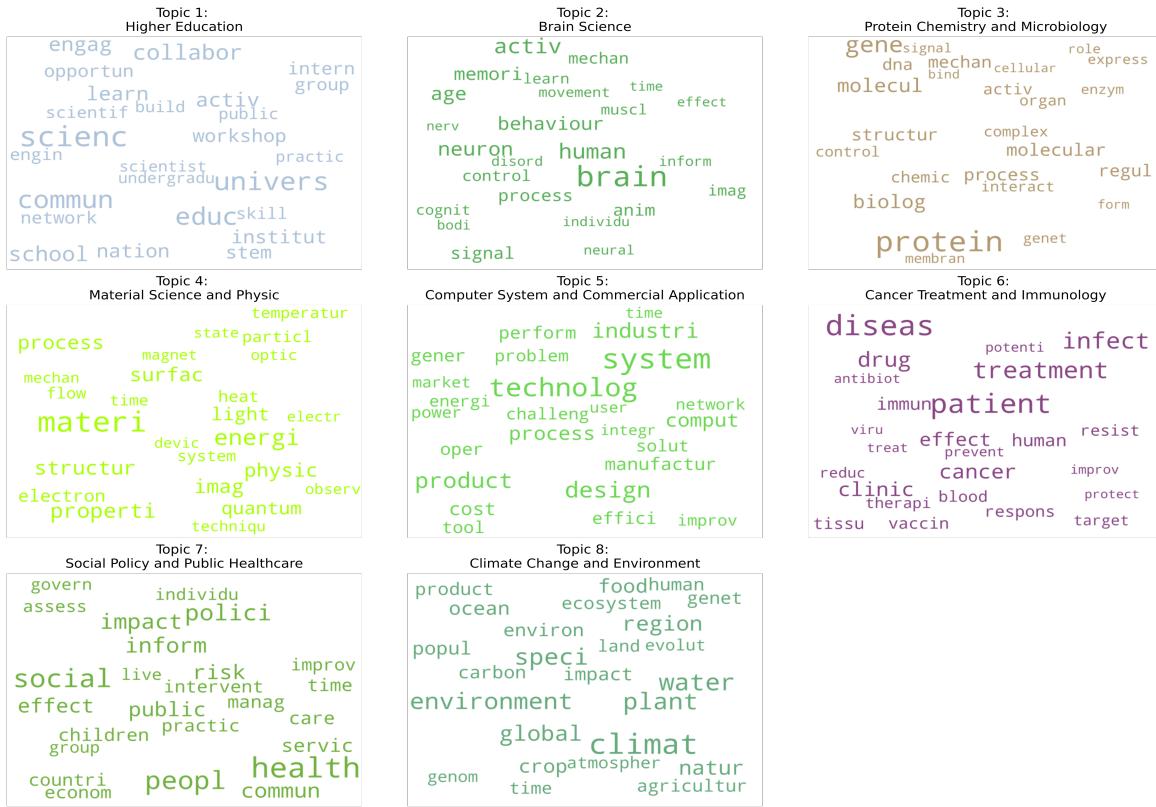


Figure 5: **Topic Wordclouds:** The figure shows the wordclouds including top 25 stemmed words of each topic derived from the LDA model, titles for the wordclouds are topic labels assigned manually.

288 4.3 Topic Analysis

- 289 In general, excluding the “Higher Education” Topic in ERC, the topic occurrence probabilities of
 290 topic 5 (Computer System and Commercial Application), topic 4 (Material Science and Physics)
 291 and topic 1 (Higher Education) are relatively high in NSF, ERC and UKRI, whereas other topics
 292 occurred relatively more rarely. For NIH agency, “Cancer Treatment and Immunology” (topic
 293 6), “Protein Chemistry and Microbiology” (topic 3), “Brain Science” (topic 2), and “Social
 294 Policy and Public Healthcare” (topic 7) have higher mean topic probabilities in total, which
 295 is not surprising as NIH is a funding agency mainly focusing on medical and health-related
 296 research, supporting that my LDA model can obtain topic distributions for projects accurately.
 297 For the group of all other agencies (Fig. 6), they are all dominated by “Computer System and
 298 Commercial Application” with the mean probability approximately $10^{-0.56}$, $10^{-0.61}$ and $10^{-0.54}$
 299 in NSF, ERC and UKRI respectively. For the agency NIH, it is dominated by “Cancer Treatment
 300 and Immunology” whose mean probability is $10^{-0.53}$. Moreover, the “Climate Change and
 301 Environment” and “Brain Science” topics are prevalent (i.e. larger than or equal to 10^{-1})
 302 solely in NSF and NIH respectively, whereas other topics are prevalent in 2 agencies at least.

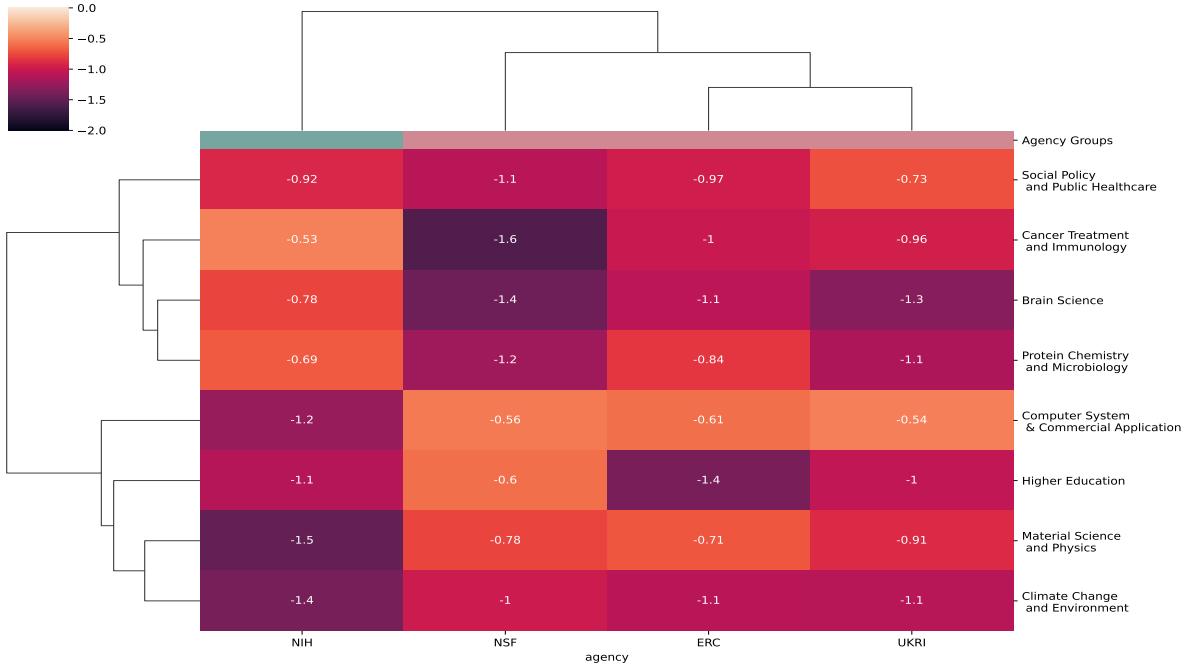


Figure 6: **Topics-Agencies Heatmap:** The heatmap between topics and four agencies. The colours inside the heatmap represent the base 10 log-scaled mean occurrence probabilities of topics. The darker is the colour, the fewer occurrence probabilities for the topic. The column and row dendograms represent the results of hierarchical clustering on agencies and topics respectively. The two colours at the top heatmap represent different amount areas of daily funding amount: dark-green is the NIH group, dark-pink is the group of all other agencies.

303 For the group of all other agencies (Fig. 7A), the mean probabilities of topic “Brain Science” that
 304 occur at the low-amount and medium-amount areas are low, but becoming larger at the high-
 305 amount area of the daily funding amount. For the topics “Protein Chemistry and Microbiology”
 306 and “Cancer Treatment and Immunology”, they are also less probable of occurring at the low-
 307 amount area and increasingly likely to be at the areas of daily medium- and high-amount of
 308 funding, even though they have the lowest occurrence probabilities occurring in the group of
 309 all other agencies. For the “Social Policy and Public Healthcare” topic, it is more likely to be
 310 funded at the range from 0 to 69.6 dollars and the range from 344.5 to 388.3 dollars, comparing
 311 with other intervals of the funding amount. The “Climate Change and Environment” topic has
 312 a relatively high probability of occurrence between 193.7 and 203.1 dollars, at the highest ranges
 313 of the low-amount area. Except for the interval between 0 to 36.5 dollars, the probability of
 314 appearance of the topic “Computer System & Commercial Application” is evenly distributed.
 315 The funded projects related to “Higher Education” are more probable to be funded at low-
 316 amount funding area than at the medium- and high-amount areas. For the topic “Material
 317 Science and Physics”, the projects in the group of all other agencies are more likely to be
 318 funded at the medium-amount area, compared to be at the low- and high-amount areas.

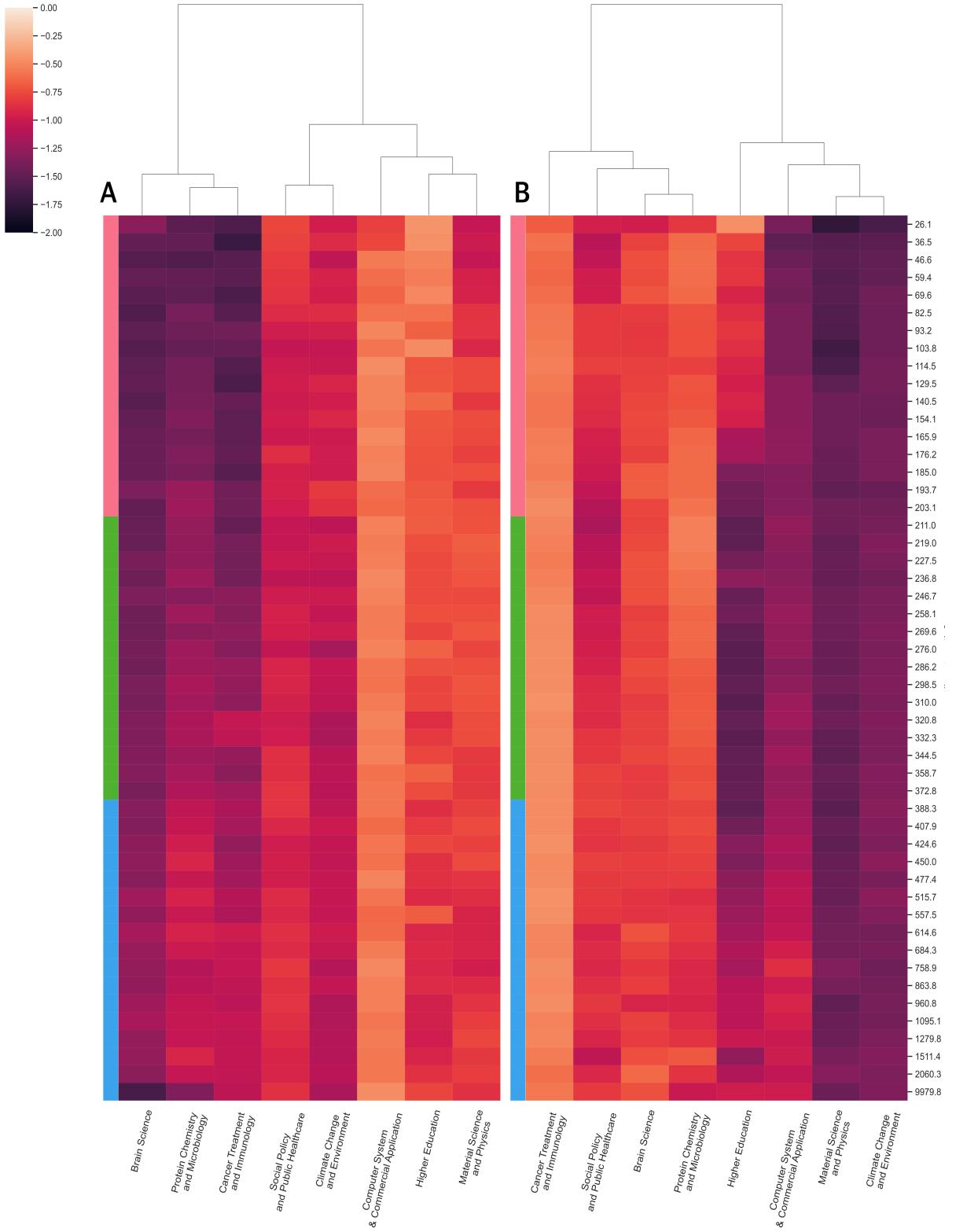


Figure 7: **Topics-Amount Heatmap:** The heatmap A and B gives the base 10 log-scaled mean occurrence probabilities of eight topics between 50 intervals of daily funding amount (see section 2.5) in NSF, ERC and UKRI agencies and in the NIH agency respectively. The vertical axis represents the discrete daily funding amount, the horizontal axis represents the agencies. Three colours on the left side of A and B describe three areas of the daily funding amount: Pink is the low-amount area, Green is the medium-amount area and Blue is the high-amount area. The column dendrogram gives the result of hierarchical clustering on topics.

319 For the NIH, a health-focused agency (Fig. 7B), topics “Climate Change and Environment”

320 and “Material Science and Physics” have equal distributions but the lowest occurrence proba-
321 bilities over three areas of daily funding agencies. The probability of occurrence of the “Cancer
322 Treatment and Immunology” topic is evenly distributed over the intervals of daily funding
323 amount and, compared with other topics, is highest in NIH agency. There is an opposite trend
324 that can be identified between the topics “Computer System & Commercial Application” and
325 “Higher Education”. The projects related to the former topic is more probable to be funded at
326 a high-amount funding area between 684.3 and 9979.8 dollars each day, whereas projects in NIH
327 agency that are relevant to the latter topic are more likely to be funded at the low-amount area
328 between 0 dollar and 154.1 dollars per day. Furthermore, the NIH funded projects related to
329 “Brain Science” have a higher probability of obtaining funding at ranges of 176.2–193.7 dollars
330 and 1279.8–2060.3 dollars per day. The “Protein Chemistry and Microbiology” appears with a
331 higher probability at the range from 26.1 and 69.6 dollars and the range from \$203.1 to \$227.5
332 per day. Additionally, the topic “Social Policy and Public Healthcare” is more likely to appear
333 in some low intervals of daily funding amount at the low-amount and high-amount areas.

334 In general, the occurrence probabilities of two dominated topics (“Computer System & Com-
335 mercial Application” and “Cancer Treatment and Immunology”) in the NIH agency group and
336 the group of all other agencies are both distributed evenly at the three areas of the daily funding
337 amount. The probability of occurrence of these two topics increases over the intervals of daily
338 funding amount if the two topics exchange groups with each other. In addition, in the group of
339 all other agencies, with the unchanged distribution of the probability that the topic “Computer
340 System & Commercial Application” occurs, the occurrence probabilities of “Protein Chemistry
341 and Microbiology” and “Brain Science” topics rise with increasing daily funding amount.

342 **4.4 Prediction of Funding Amount**

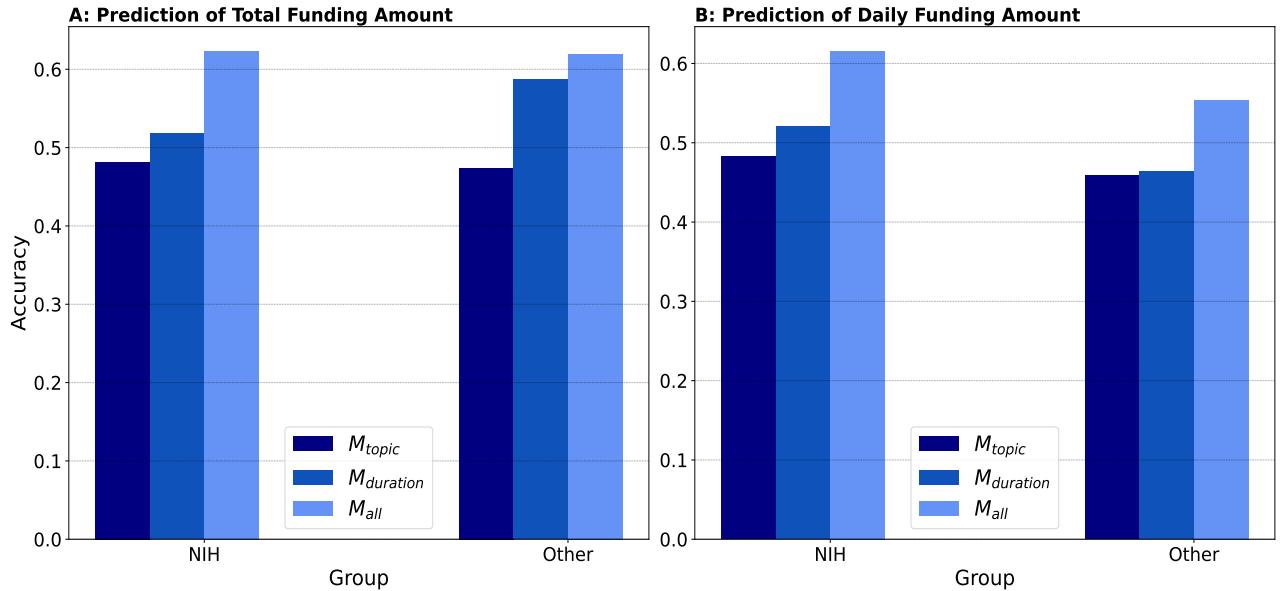


Figure 8: **Prediction Accuracy:** The accuracy of prediction of the total and daily funding amount using LightGBM in NIH group and the group of all other agencies.

343 As predicting the total funding amount, the M_{topic} in both the NIH group and all other agencies
 344 obtained the lowest accuracy, at approximately 0.482 and 0.474. Having only project duration
 345 as a feature, the accuracy of $M_{duration}$ in all other agencies is 0.587, 13.3% higher than that of
 346 $M_{duration}$ in the NIH group. As combining all topic distributions and project durations together
 347 as predictors, both the accuracy of M_{all} are beyond 0.6, boosting 20.4% in the NIH group and
 348 5.5% in all other agencies compared with the accuracy of M_{topic} . When the prediction target is
 349 the daily funding amount, a different situation is seen in the group of all other agencies. The
 350 accuracy of $M_{duration}$ and M_{all} in all other agencies decreases by 10.5% and 8.5% respectively,
 351 compared with that of corresponding models in the NIH group. However, regardless of whether
 352 the prediction target is either total or daily funding amount, the accuracy of prediction is
 353 improved in all groups with the topic distributions as predictors.

354 **5 Discussion**

355 By using combined machine-learning approaches, I found that topic modelling techniques are
 356 helpful for topic assignment in these four funding agencies because the LDA model correctly
 357 assigned medicine- and health-related projects into NIH agency. Also, I looked into the part
 358 of project abstracts in the other three agencies and the topics also can be found interpretable.
 359 Moreover, the LightGBM might be useful for the prediction of funding amount since the accu-
 360 racy of models M_{all} can beyond 0.6 as I used both the topic distributions and project duration

361 as prediction variables, even though the accuracy of M_{all} is poorer as predicting daily funding
362 amount than total funding amount in the group of all other agencies. However, the accuracy is
363 below 0.5 if I use only topic distributions as predictors to forecast either total or daily funding
364 amount, indicating that the topic distributions generated from the LDA model don't show the
365 prediction power.

366 From heatmap analysis on topics, a significant topic bias can be seen in four prevalent academic
367 agencies. Except for NIH, which is a health-focused funding agency, the funding of the other
368 three agencies is occupied by topics related to Computer Systems, Commercial Applications,
369 Material Science and Physics. In addition, protein chemistry and microbiology are also prevalent
370 topics in the ERC agency. Thus, NSF, ERC and UKRI are biased by STEM-related topics.
371 The topic "Climate Change, Environment" are ignored significantly by all four popular funding
372 agencies. Only two dominated topics in NIH ("Computer System and Commercial Application")
373 and in all other agencies ("Cancer Treatment and Immunology") are equally distributed in
374 almost 50 intervals of the daily funding amount.

375 However, a positive or negative relationship with daily funding amount also can be seen in
376 two groups. For example, "Cancer Treatment and Immunology", dominating the group of all
377 other agencies, increases its occurrence probability with funding amount in the NIH group, the
378 occurrence probability of "Computer System & Commercial Application", dominating the NIH
379 group, also rises with daily funding amount in the group of all other agencies. This suggests
380 that a project that is mainly relevant to "Cancer Treatment and Immunology" and "Computer
381 System and Commercial Application" (e.g. Medical Engineering) will have a higher probability
382 of obtaining a higher daily funding amount in any of four agencies. Similarly, a project will be
383 more likely to be funded with a high amount in the group of all other agencies if it includes a
384 large portion of "Computer System & Commercial Application" with "Cancer Treatment and
385 Immunology" or "Protein Chemistry and Microbiology" or "Brain Science". (e.g. possibly
386 Computational Biology or Computational Neuroscience).

387 So far, the quantitative analysis combining the LDA model and LightGBM is first introduced in
388 this field. Thus, there are some limitations to the project and possible solutions to be completed
389 in the future. In this project, the dataset of four agencies is obtained in different ways, resulting
390 in the difficulty of data collection. Therefore, one can develop a program that can automatically
391 scrape up-to-date data from the websites of the agencies without the breach of data privacy
392 statement, allowing the models used in the project to update Regularly.

393 Moreover, I used an 8-topic model by comparing models that are run multiple times with
394 different k . However, because of time and computational restriction, only 32 topic models
395 are included and each model is repeated only 10 times, leading the final topics to the lack of

396 specificity. For example, for the topic “social policy and public healthcare”, I can not give a
397 precise conclusion that the prevalence of the topic is due to “Social Policy” or because of “Public
398 Healthcare”. The problem can be solved by testing a wider range of k and more repetition times
399 for each model if time and computational power allow.

400 Second, the parameters of all LightGBM models are all set to be the same on a rule-of-thumb
401 basis. However, a better approach to parameter optimisation could be implemented in order
402 to mitigate the risk of overfitting. For example, the Bayesian hyper-parameter optimization
403 algorithm was utilised by Wang *et al.* to improve the performance of LightGBM [37]. Next, I
404 used $M_{duration}$ and M_{all} to study the improvement of the model accuracy by topic distributions.
405 However, the only project duration can not always provide sufficient information of the data,
406 leading to poor accuracy. This might be the reason why $M_{duration}$ and M_{all} in the group of all
407 other agencies give worse performance as predicting the daily funding amount than predicting
408 the total funding amount. Therefore, one can add more factors that could affect funding amount
409 to improve the accuracy of the models such as the h-index of PI and the rank of PI’s university
410 *etc.*

411 Then, to assess the performance of the LightGBM models, only the accuracy metric is used
412 (see section 2.5). However, the accuracy (or precision) score can be reasonable only if the cost
413 of false positives is high. Therefore, more approaches can be applied to assess the model per-
414 formance such as recall, F1 Score and Receiver Operating Characteristic Curve. Furthermore,
415 prediction on funding amount has been treated as a multiclassification problem. So, other pre-
416 diction methods are worthy of consideration in order to acquire more specific outcomes such as
417 regression on funding amount.

418 Last but not the least, due to the impossible collection of unfunded project data, the relationship
419 between topic distributions and funding failure can not be explored yet. So, if the information
420 of unfunded projects is available, one can use the same approaches in the research to forecast
421 funding success and failure. For example, instead of predicting academic funding amount, a
422 similar method was developed by Mitra *et al.* (2014) for predicting crowdfunding success and
423 failure by penalized logistic regression with phrases produced by unigram, bigram and trigram
424 models as features [9].

425 6 Conclusion

426 In conclusion, I answer the five questions mentioned in the introduction section. (1) With
427 topic modelling approach, it enables us to obtain the interpretable topics from four academic
428 funding agencies. (2) Except for the topics “Cancer Treatment and Immunology” in NIH agency

429 and “Computer System and Commercial Application” in all other agencies, the eight topics
430 have different occurrence probabilities over discrete daily funding amount, meaning that the
431 existence of topic bias in different intervals of daily funding amount and they may provide some
432 information for prediction of the funding amount. (3) The research projects that are relevant
433 to “Computer System and Commercial Application” and “Cancer Treatment and Immunology”
434 possibly like computational epidemiology or medical engineering *etc.* have a higher probability
435 of obtaining more funding in all funding agencies. However, some projects that are related
436 to Environment, Ecosystem and Environment with Education or Social Policy are more likely
437 to be funded with the less daily amount. (4) I found it useful for LightGBM to classify the
438 three classes in both of total and daily funding amount, even though the accuracy is poorer as
439 predicting the discrete daily funding amount in the group of all other agencies. (5) Using topic
440 distributions as predictors is able to boost the accuracy of prediction by approximately 5% to
441 20%.

442 Support Information

443 All stop words used in this project

444 {"'november', 'shouldn', 'necessary', 'inc', 'twice', 'whether', 'from', "you'll", 'aspect', 'pro-
445 posed', 'measure', 'as', 'significant', 'took', 'w', 't', 'his', 'thru', 'almost', 'proposal', 'vi', 'out',
446 'say', 'someone', 'd', 'www', "didn't", 'mostly', 'gone', 'edu', 'f', 'never', 'h', 'chinese', 'october',
447 'down', 'hither', 'those', 'serious', 'specific', 'our', 'is', 'are', "you've", 'poorly', 'himself', 'pro-
448 gram', 'i', 'become', 'awfully', 'o', 'something', 'near', 'come', 'detected', 'pursue', 'less', 'febru-
449 ary', 'able', 'consider', 'just', 'gotten', 'until', 'going', 'changes', 'nine', 'away', 'three', 'unless',
450 'somebody', 'zero', 'progress', 'doi', 'like', 'upon', 'although', 'seemed', 'together', "it's", 'en-
451 able', 'behind', 'recent', 'e', 'itself', 'august', 'because', "isn't", 'iv', 'far', 'considering', 'often',
452 'made', 'therein', 'indeed', 'need', "shan't", 'x', 'inasmuch', 'all', 'so', 'same', 'rather', 'obvi-
453 ously', 'others', 'no', 'theirs', 'associated', 'l', 'thoroughly', "that'll", 'four', 'beside', 'thereby',
454 'needed', 'herself', 'training', 'latter', 'let', 'knows', 'main', 'exactly', 'aren', 'third', 'follow-
455 ing', 'linkage', 're', 'beyond', 'enough', 'thats', 'not', 'll', 'application', 'got', 'lately', 'herein',
456 'respectively', 'cover', 'lest', "couldn't", 'but', 'value', 'focus', 'before', 'large', 'he', 'seem', 'inso-
457 far', 'truly', 'english', 'do', 'yourself', 'nor', 'beforehand', 'vii', 'toward', 'certainly', 'containing',
458 'second', 'been', 'sure', 'academic', 'q', 'ought', 'applicant', 'willing', 'example', 'about', 'becom-
459 ing', 'currently', 'between', 'everything', 'liked', 'very', 'many', 'study', 'one', 'anyhow', 'across',
460 'relate', 'etc', 'forth', 'seeming', 'hopefully', 'aim', 'km', 'next', 'specifying', 'tell', 'during', 'in-
461 crease', 'since', 'fifth', 'still', 'plus', 'against', 'though', 'small', 'quite', 'durham', 'oh', 'cer-
462 tain', 'approach', 'baseoxford', 'too', 'name', 'causes', 'theres', 'such', 'after', 'course', 'through',
463 'count', 'might', 'indicates', 'neither', 'future', 'examine', 'face', 'five', 'either', 'k', 'bed', 'no-
464 body', 'wherever', 'had', 'cannot', 'kept', 'concerning', 'formerly', 'here', 'go', 'once', 'way',
465 'understand', 'lead', 'available', 'however', 'uses', 'alone', 'specify', 'well', 'whither', 'off', 'eas-
466 ily', 'last', 'thanx', 'whereby', 'below', 'immediate', 'please', 'how', 'at', 's', 'taken', "needn't",
467 'a', 'without', 'then', 'only', 'induce', 'seen', 'ourselves', 'six', 'http', 'take', 'further', 'give',
468 'moreover', 'africa', 'reasonably', 'wasn', 'least', 'okay', 'unlikely', 'everyone', 'known', 'myself',
469 'afterwards', 'elsewhere', 'seven', 'therefore', 'contains', "don't", 'needn', 'later', 'local', 'accord-
470 ing', 'she', 'good', 'most', 'little', 'which', 'closely', 'highly', 'see', 'sup', 'comes', 'much', '_', 'im-
471 portant', 'anyone', 'japanese', 'getting', 'recently', 'always', 'p', 'david', 'ok', 'using', 'et', 'some-
472 what', 'entirely', 'ma', 'became', "you'd", 'dataset', 'continued', 'c', 'field', 'they', 'great', 'some-
473 time', 'trying', 'essential', 'work', 'won', 'mustn', 'relatively', 'look', 'know', 'along', 'amongst',
474 'apart', 'west', 'now', 'possible', 'somewhere', 'seeing', 'definitely', 'descriptionabstract', 'viz',
475 'have', 'within', 'tried', 'overall project summary', 'throughout', 'any', 'nothing', 'test', 'per-
476 haps', 'programme', 'my', 'march', 'under', 'better', 'given', 'whole', 'sheet', 'thereafter', 'goes',

477 'shan', 'sensible', 'channel', 'and', 'in', 'consequently', 'asking', 'each', 'greetings', 'description
478 provided by applicant:', 'december', 'july', 'anything', 'onto', 'isn', 'extend', 'whereupon', 'hav-
479 ing', 'to', 'vast', 'lt', 'likely', 'while', 'inward', 'meanwhile', 'april', 'hereby', 'usually', 'you', 'co',
480 'than', 'survey', 'why', 'be', 'outside', 'even', 'different', 'we', 'develop', 'also', 'knowledge', 'how-
481 beit', 'successful', 'provides', 'indicated', 'y', 'nowhere', 'normally', 'went', 'whenever', 'whereas',
482 'investigate', 'un', 'hers', 'soon', 'affect', 'range', 'first', 'think', 'believe', 'award', 'seriously', 'lat-
483 terly', 'particular', 'questions', 've', 'them', 'instead', 'secondly', 'whereafter', 'th', 'other', 'was',
484 'project description', 'india', 'hello', 'mainly', 'used', 'by', 'anybody', 'whatever', 'done', 'needs',
485 'help', 'whoever', 'various', 'this', 'couldn', 'above', 'r', 'did', "doesn't", 'cell', 'does', 'gets', 'de-
486 cide', 'follows', 'somehow', 'namely', 'there', 'unto', 'haven', 'another', 'include', 'september',
487 'description', 'hi', 'line', 'more', 'maybe', "hasn't", 'came', 'says', 'based', 'besides', 'neverthe-
488 less', 'cause', 'g', 'thus', 'don', 'with', 'placed', 'should', 'due', 'doing', 'specified', 'own', 'yes',
489 'among', 'mightn', 'themselves', 'ignored', 'sorry', 'again', "shouldn't", 'nice', 'saw', 'anyway',
490 "hadn't", "mustn't", 'june', 'show', 'successfully', 'yours', 'mean', 'right', 'contain', 'students',
491 'especially', 'shall', 'underly', 'china', 'its', 'can', 'looks', 'everybody', "weren't", 'ms', 'clearly',
492 'us', 'hence', 'anywhere', 'hasn', 'selves', 'brief', 'ltd', 'tries', 'inner', 'up', 'appear', 'sent', 're-
493 mains', 'big', 'provide', 'an', 'use', 'around', "wasn't", 'except', 'regarding', 'for', 'probably',
494 'receive', 'uk', 'appropriate', 'sector', "you're", 'their', 'method', 'appreciate', 'south', 'presum-
495 ably', 'could', 'want', 'b', 'greatly', 'get', 'ii', 'few', "should've", 'u', 'the', 'really', 'didn', 'apply',
496 'hadn', 'research', 'some', 'project', 'seems', 'her', 'looking', 'tends', 'becomes', 'furthermore',
497 'nearly', 'despite', 'goal', 'may', 'towards', 'funded', 'ex', 'unfortunately', 'support', 'thanks',
498 'cant', 'try', 'none', 'whom', 'into', 'former', 'me', 'thereupon', 'ask', 'happens', 'keep', 'said',
499 "haven't", 'thence', 'hardly', 'doesn', 'wherein', 'downwards', 'merely', 'call', 'british', 'two', 'm',
500 'whence', 'seminar', 'on', 'z', 'sometimes', 'yet', 'these', 'com', 'where', 'several', 'london', 'your-
501 selves', 'subject', 'nd', 'qv', 'both', 'corresponding', 'org', 'input', 'rd', 'new', 'dr', 'otherwise',
502 'purpose', 'your', 'would', 'everywhere', 'best', 'particularly', 'thorough', 'make', 'supports',
503 'regards', "she's", 'sub', 'predict', 'area', 'actually', 'wonder', 'being', 'easy', 'lack', 'propose',
504 'must', 'ie', 'easet', 'data', 'him', 'will', 'thank', 'hereupon', 'wouldn', 'hereafter', 'has', 'ab-
505 stract', 'of', 'https', 'indicate', 'per', 'worldwide', 'regardless', 'que', 'shown', "aren't", 'that',
506 'north', "wouldn't", 'vs', 'over', 'whose', 'model', 'areas', 'novel', 'weren', 'uucp', 'when', 'use-
507 ful', 'gives', 'every', 'ones', 'january', 'were', 'high', 'described', 'v', "mightn't", 'accordingly',
508 'unique', 'n', 'old', 'am', 'ours', 'if', 'what', 'already', 'or', 'iii', 'else', 'objective', 'function', 'eg',
509 'anyways', 'followed', 'summary', "won't", 'it', 'estimate', 'j', 'allow', 'summaryabstract', 'say-
510 ing', 'require', 'lot', 'non', 'ain', 'funding', 'overall', 'keeps', 'amountdetail', 'contact', 'wants',
511 'run', 'identify', 'ever', 'viii', 'self', 'experience', 'who', 'wish', 'via', 'eight', 'noone', 'aside',
512 'welcome', 'participation' }

513 **Topic Similarity and Model Stability**

514 To explore how the two topics in two sub-models at the same k-topic model differ from each
515 other. I calculated Jensen-Shannon Divergence (JSD) between best-matched topic pairs since
516 it has been proven to outperform other divergence-based similarity metrics [38]. Due to the
517 stochasticity of the LDA model, topics with similar words won't be in the same order be-
518 tween sub-models with different random states. I applied the Kuhn-Munkres algorithm in the
519 'scipy' package in Python [39], also known as the Hungarian algorithm, in order to solve this
520 minimum weight bipartite matching problem [40]. Afterwards, the JSD distance between two
521 best-matched topics from two different sub-models of a k-topic model (see Section 2.5 Model
522 Selection) can be calculated by using the 'Gensim' package in Python [21]. Then, from these
523 calculated distances between topics, the averaged *similarity* between two sub-models in the
524 same k-topic model can be expressed as follows:

$$Sim(M_u^k, M_v^k) = \frac{1}{k} \sum_{x=1}^k JSD(Z_{ux}, \phi(Z_{vx})) \quad (9)$$

525 where M_u^k and M_v^k are $u - th$ and $v - th$ sub-models in the model set with the topic indicator
526 k ; $JSD(Z_{ux}, \phi(Z_{ux}))$ are the JSD measure between the best-matched topics Z_{ix} and $\phi(Z_{ix})$,
527 where Z_{ix} is the $x - th$ topic in M_u^k and $\phi(Z_{ux})$ is the most similar topic (i.e. best-matched)
528 with Z_{ux} in M_v^k .

529 Afterwards, the *stability* between two k-topic models is expressed as:

$$Stab(M^k) = \frac{2}{n \times (n-1)} \sum_{u,v,u < v}^r Sim(M_u^k, M_v^k) \quad (10)$$

530 where n is the number of sub-models in a single k-topic model, which is 10 in my project.
531 The topic model with the lowest *stability* and median *similarity* of pair-wise sub-models are
532 considered as my final topic model.

533 **Hierarchical Clustering**

534 To explore how similar topics change with different funding agencies and funding amounts. The
535 hierarchical clustering with the Ward linkage method, known as the minimal increase of sum-of-
536 squares (MISSQ), is applied for clustering similar topics. Ward linkage are more explicable way
537 to cluster noisy data [41]. The distance between two clusters can be calculated by the sum of
538 squares of the deviations from mean or the centroid (i.e. the error of sum squares or ESS)[42].

539 The Ward approach tries to find the minimal increase of ESS at every iteration. The distance
540 between clusters using the Ward approach can be expressed as [41]:

$$Dist(X, Y) = ESS(X \cup Y) - [ESS(X) + ESS(Y)] \quad (11)$$

541 where X, Y are two different clusters, $X \cup Y$ is a combined cluster of X and Y.

542 Topic Label Assignment

543 To assign more interpretable labels for topics derived from the LDA model, pyLDAvis is utilised,
544 which is a Python package providing a method of interactive visualisation for the LDA model
545 [43]. In the approach, the *relevance* measure of words of topics is defined, allowing to rank
546 topic terms by not only the probability under a topic but also its lift. The *relevance* can be
547 described as [43]:

$$r(w, k|\lambda) = \lambda \log(\phi_{kw}) + (1 + \lambda) \log\left(\frac{\phi_{kw}}{p_w}\right) \quad (12)$$

548 where ϕ_{kw} is the probability of a word w for topic k , p_w is defined as the marginal probability of
549 term w , λ ($0 \leq \lambda \leq 1$) is a parameter that can tune the weight between logged probability (ϕ_{kw})
550 and its logged lift ($\frac{\phi_{kw}}{p_w}$). If topic terms are ranked by only logged probability estimated by the
551 LDA model, it could affect interpretability of topics due to globally frequent tokens; if the topic
552 terms are ranked by solely the *lift*, which could weaken the impact of globally frequent tokens
553 but also could result in uninterpretable rare words [44]. Hence, a suitable λ value can give a
554 solution for the problem of the non-interpretability of topics brought by either globally frequent
555 words or field-specialised terminologies.

556 Topic Inference of the LDA Model

557 Two common approaches are utilised to infer the coefficients for document-topic and topic-
558 word matrices [20]: variational Bayesian approach and collapsed Gibbs sampling approach.
559 The former method is a variational Bayesian method with an Expectation-Maximization (EM)
560 algorithm proposed by Blei et al. (2003) [19] in their original paper of the LDA model. By the
561 variational approach, the computational process is faster but may produce inaccurate inferences
562 [45]. The latter approach is to use collapsed Gibbs sampling to infer the coefficients of matrices,
563 which is introduced by Griffiths and Steyvers(2004). For this method, theoretically, the accuracy
564 of inference can be improved but with the cost of more intensive computation, which could

565 therefore affect the efficiency of the model as the size of datasets are larger [20].

566 Topic space and top-30 words of eight topics by pyLDAvis

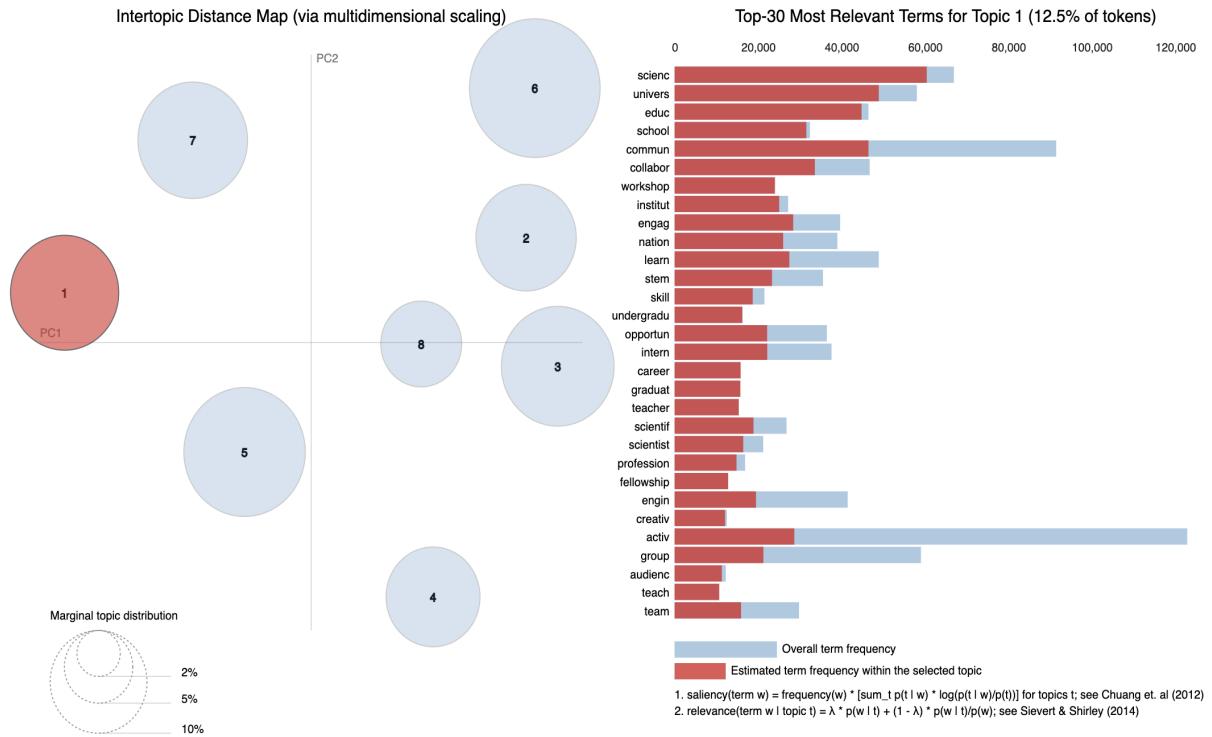


Figure 9: Visualisation of topic 1 by pyLDAvis

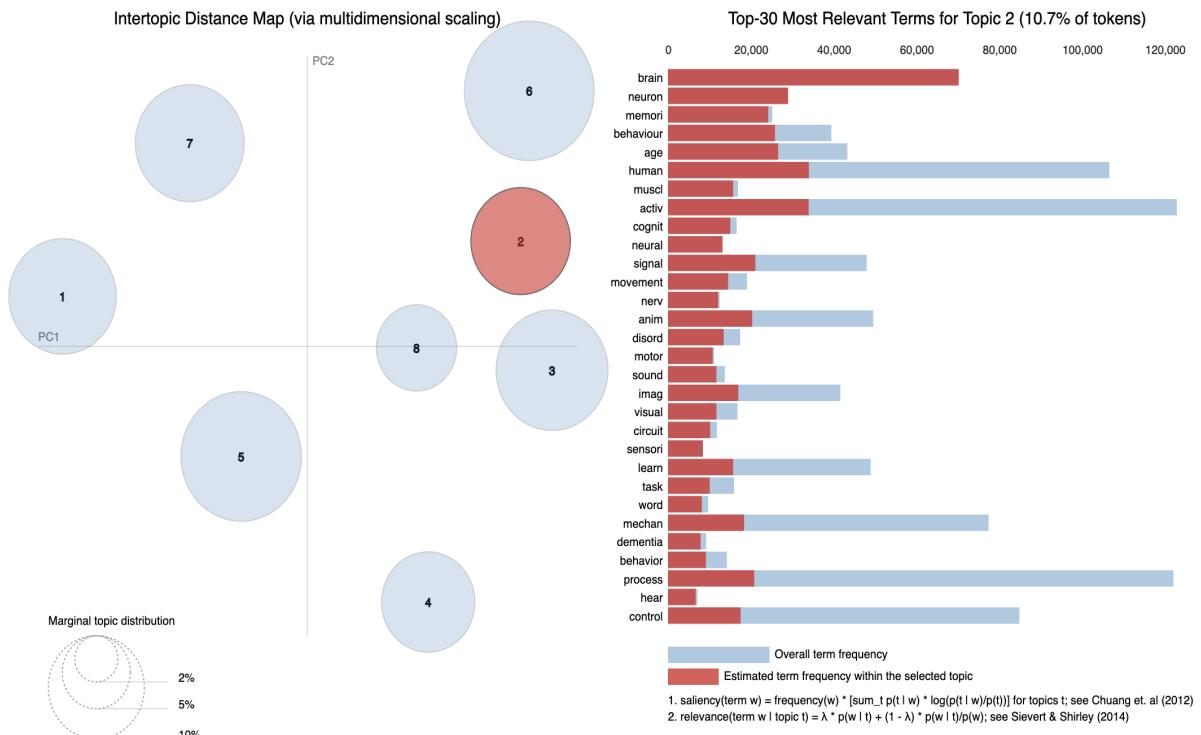


Figure 10: Visualisation of topic 2 by pyLDAvis

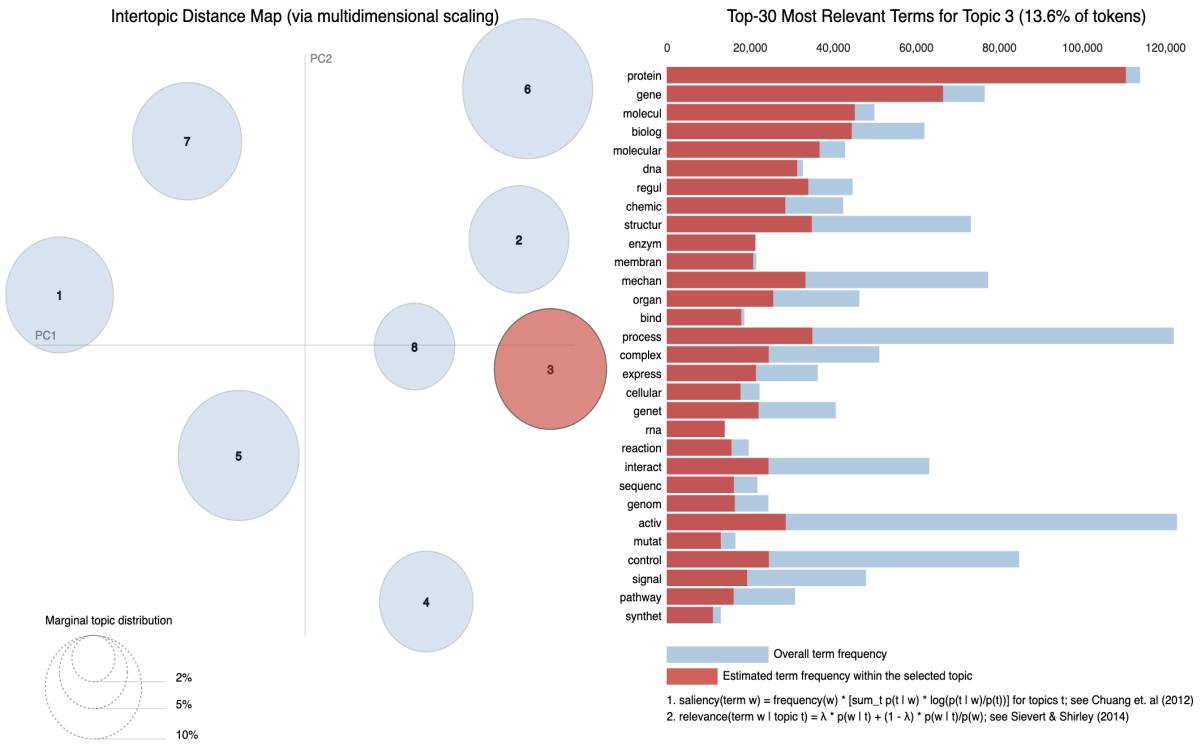


Figure 11: Visualisation of topic 3 by pyLDAvis

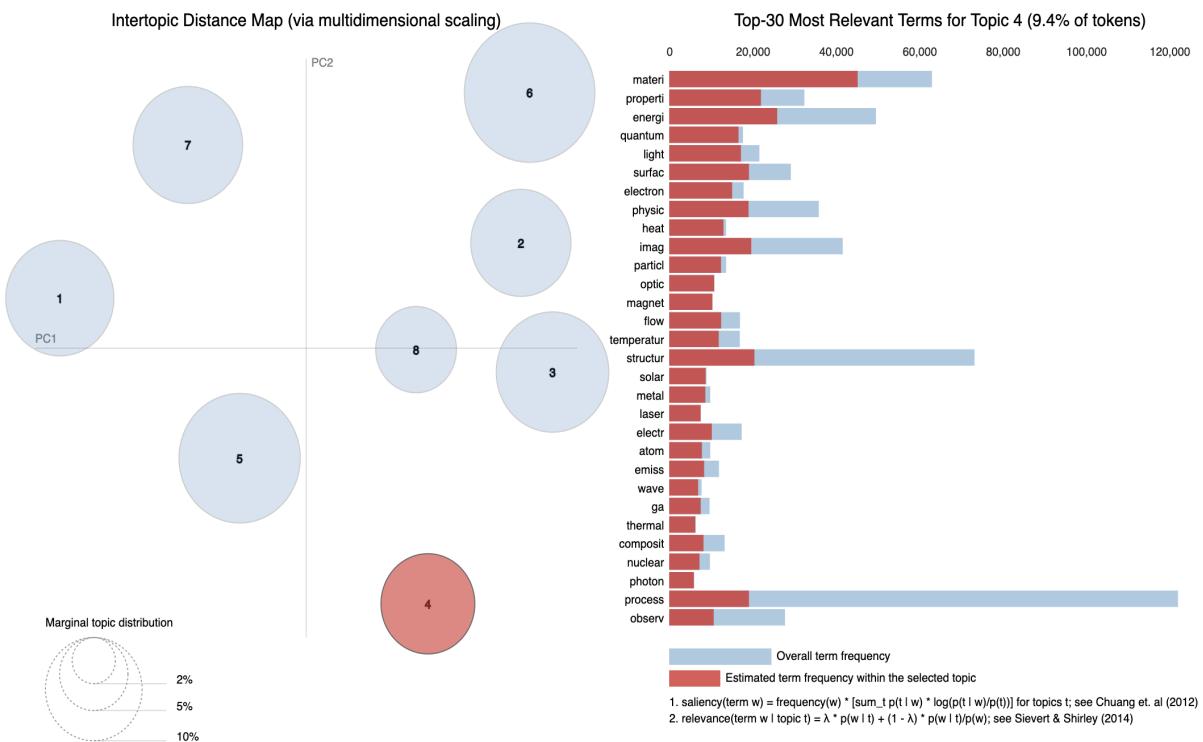


Figure 12: Visualisation of topic 4 by pyLDAvis

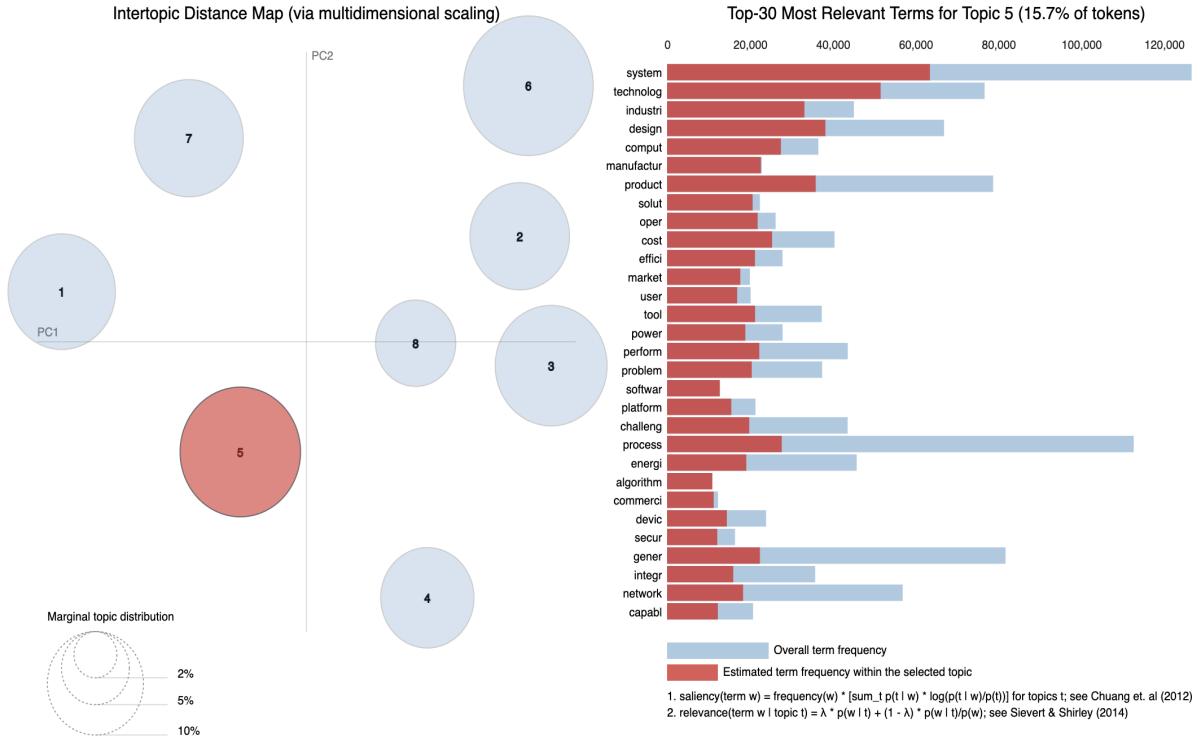


Figure 13: Visualisation of topic 5 by pyLDAvis

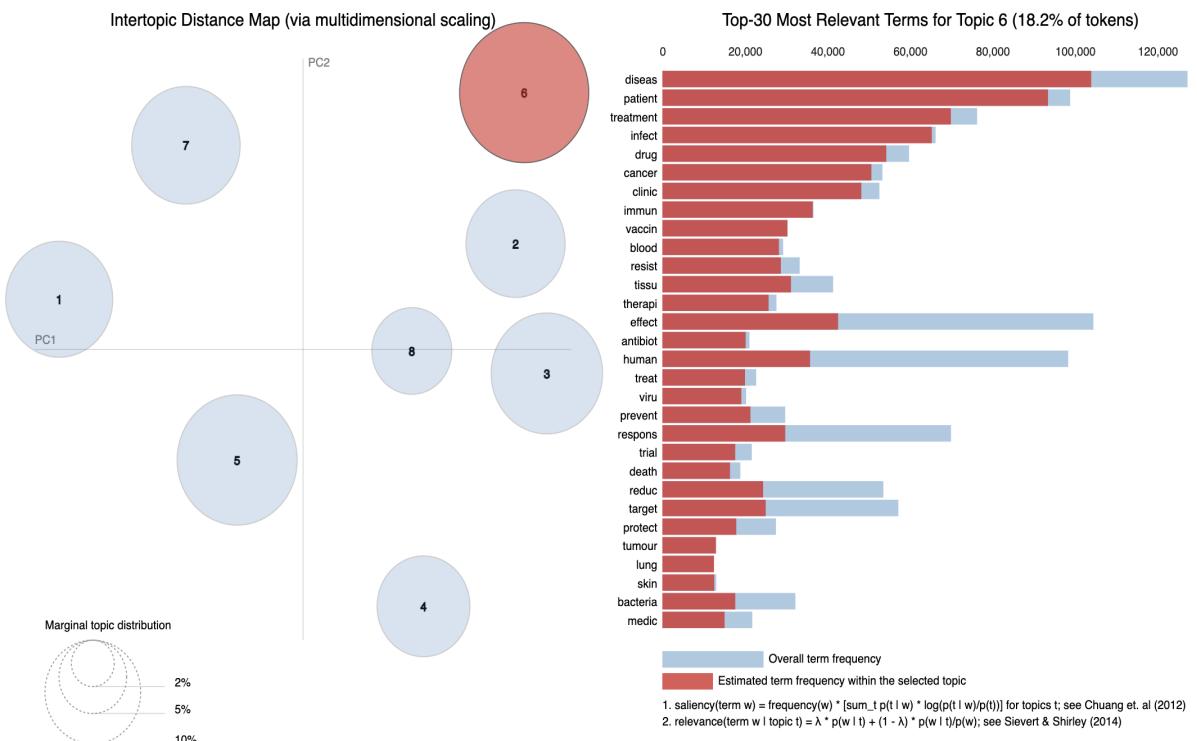


Figure 14: Visualisation of topic 6 by pyLDAvis

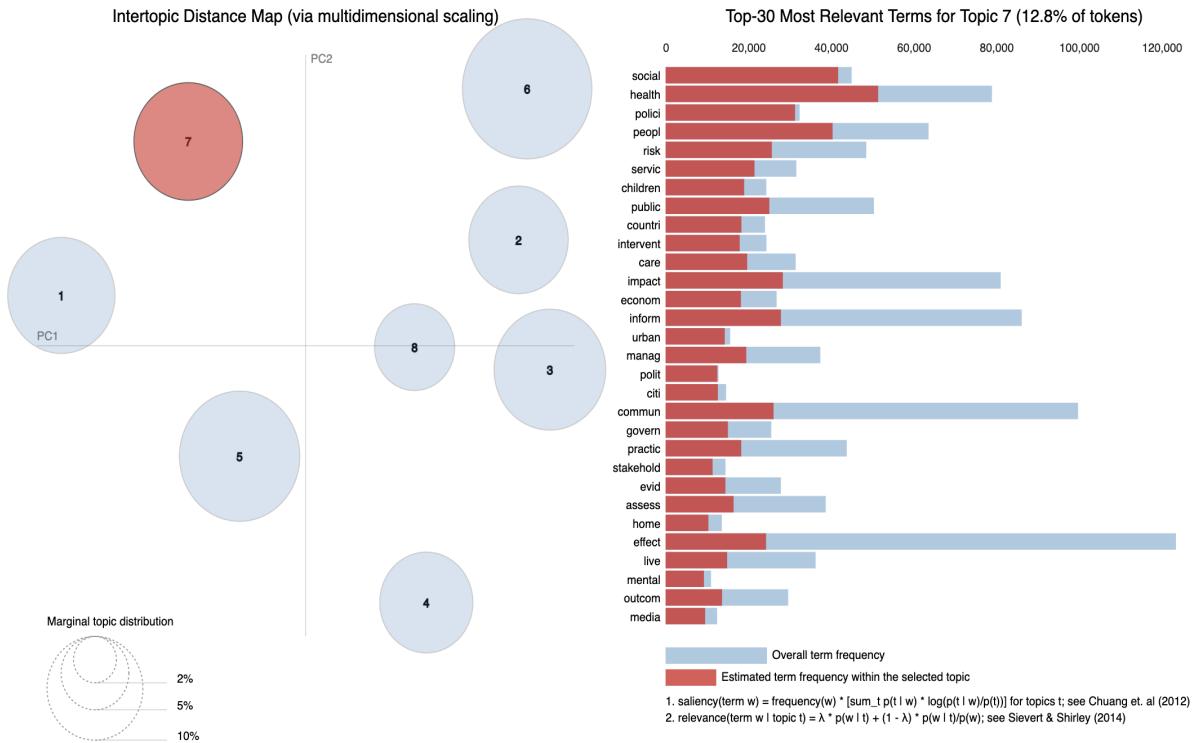


Figure 15: Visualisation of topic 7 by pyLDAvis

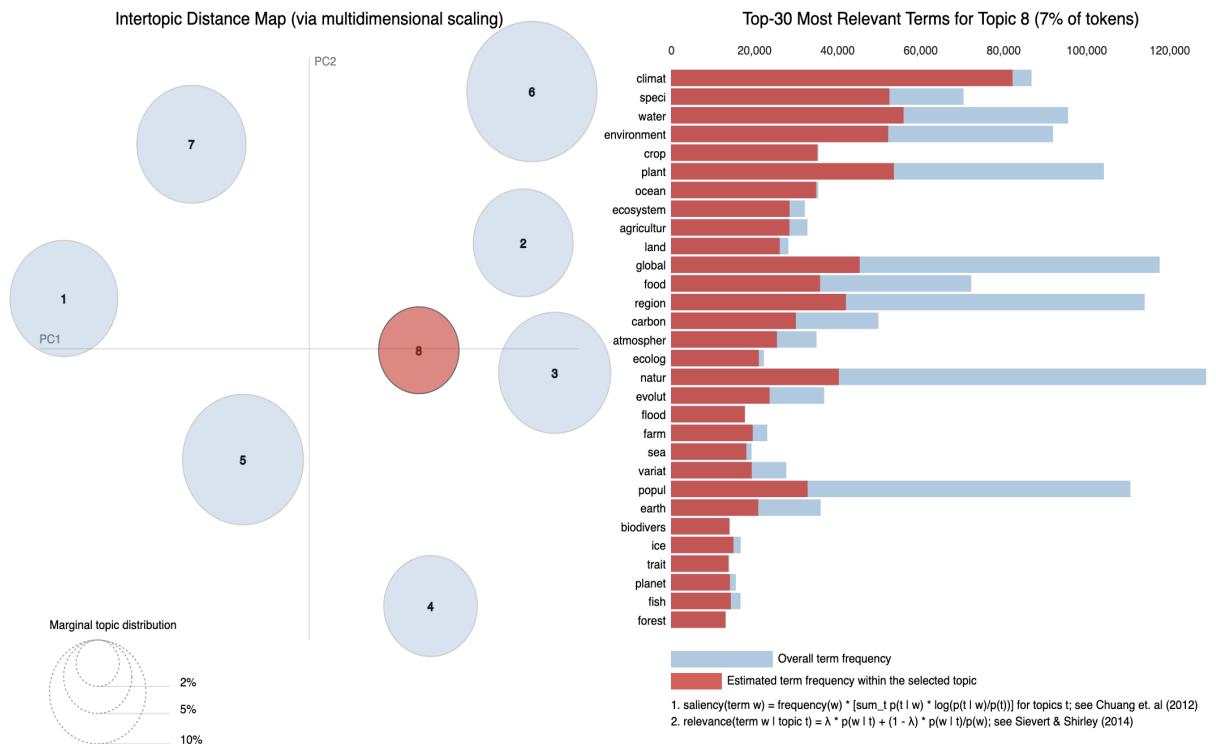


Figure 16: Visualisation of topic 8 by pyLDAvis

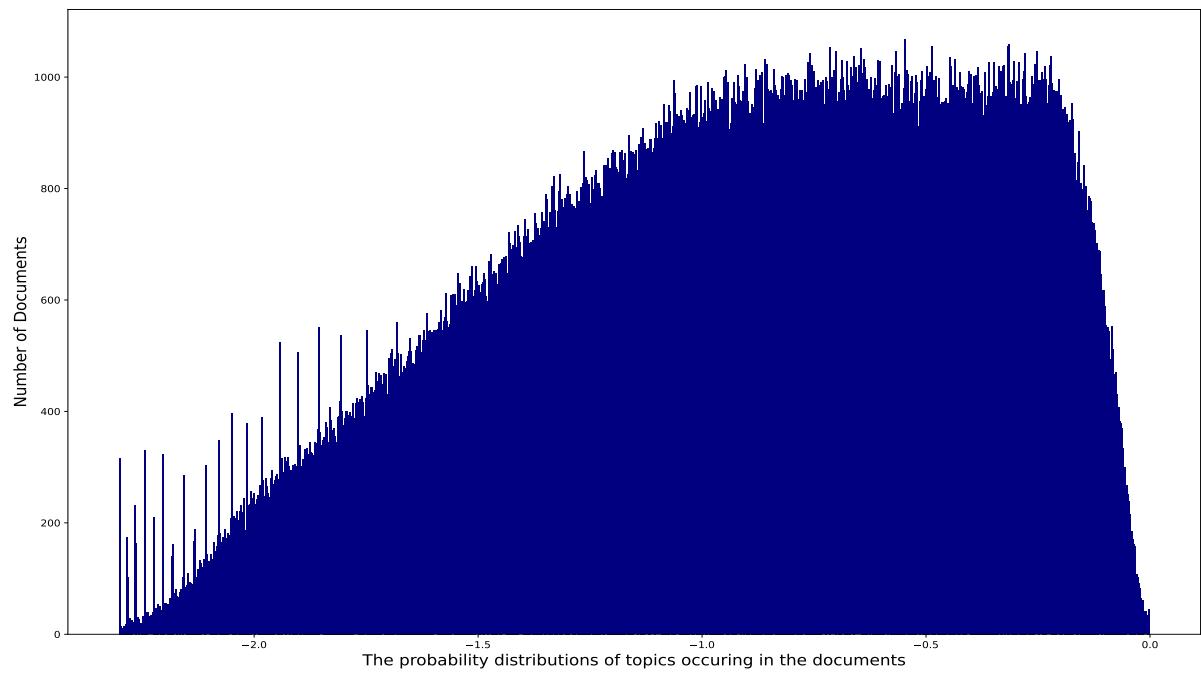


Figure 17: The probability distribution of topic probabilities that are more than 0.005 in all documents

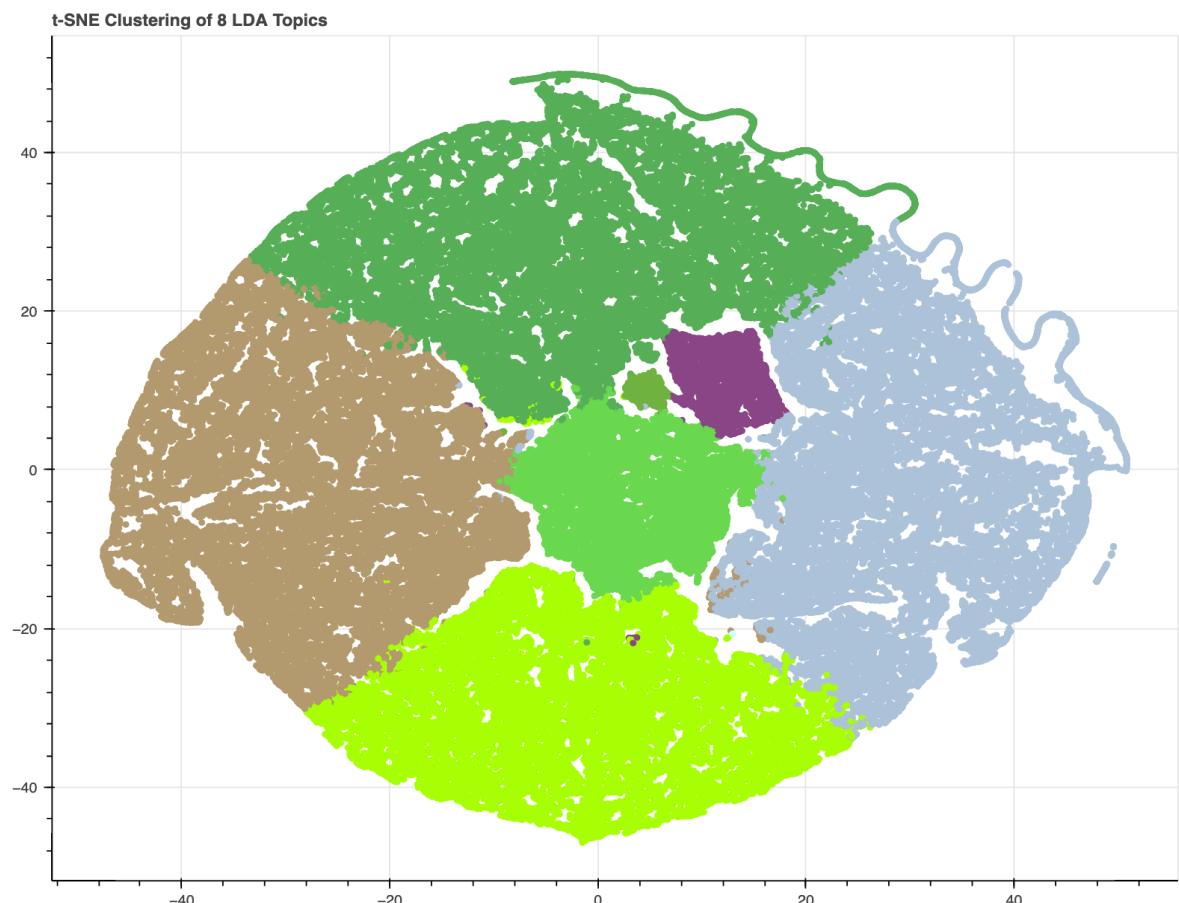


Figure 18: t-SNE clustering for visualisation of eight topics



Figure 19: Word count and importance of topic keywords

Doc 0: complex devic engin mechan membran properti system time character coupl fundamant interact isol long_term particl singl type
Protein Chemistry and Microbiology
Doc 1: challeng decreas effect find control reli role modul target urgent independ seek suggest gener strong treatment respond
Brain Science
Doc 2: design effect evalu mechan strategi establish experiment outcom result limit replac block demonstr influenc process promis role
Brain Science
Doc 3: key mechan properti strategi system circuit collect exhibit long_term reduc util accomplish role highlight initi modul previou
Brain Science
Doc 4: build integr life mechan system team establish find result set character interact reduc state collabor demonstr emerg
Brain Science
Doc 5: decreas institut reduc replac state major discoveri oper art advanc plan improv instrument addit microscop user footprint
Higher Education
Doc 6: challeng effect energi product technolog analyz establish result time coupl photon reduc singl util control demonstr role
Brain Science
Doc 7: build combin mechan natur power comput featur time collect interact phenomena util core influenc process distinct societ
Brain Science
Doc 8: design effect strategi directli establish experiment circuit interact limit tool block extern potenti role modul target central
Brain Science
Doc 9: product topic limit reservoir extern role speed creat host oper order urgent decad direct maintain critic therapi
Cancer Treatment and Immunology
Doc 10: effici effort manag product technolog approxim time character long_term state type origin class distinct genet initi previou
Climate Change and Environment
Doc 11: mechan rout structur character encod major potenti process role creat modul fact gap suggest activ brain diseas
Protein Chemistry and Microbiology
Doc 12: combin conduct effect evalu lower natur perform system directli languag result scienc set social error eventu
Computer System & Commercial Application
Doc 13: composit decreas evalu perform featur result character encod exhibit isol level reduc singl type address analysi core
Brain Science
Doc 14: mechan power structur system technolog establish exhibit implement resourc tool address basi genet contribut imag loss critic
Brain Science
Doc 15: mechan product system establish result level limit manipul reduc potenti cellular order target vivo independ proven brain
Cancer Treatment and Immunology
Doc 16: decreas directli fast find result coupl level reduc block demonstr potenti role suppress modul previou vivo gap
Cancer Treatment and Immunology
Doc 17: build complex design lower environ fundamant interact level limit thousand tool crucial demonstr major open discoveri genet
Climate Change and Environment
Doc 18: effect evalu integr mechan strategi team experiment framework result time act inform long_term probe protect singl state
Brain Science
Doc 19: build effect innov strategi system establish experiment find practic set reduc analysi commun demonstr potenti deliv target
Social Policy and Public Healthcare
Doc 20: build challeng complex construct integr system comput experiment featur scale social character circuit collect environ interact manipul
Brain Science
Doc 21: composit construct decreas geometri solut comput experiment result occupi reduc reson type influenc major simul arrang chemic
Protein Chemistry and Microbiology
Doc 22: challeng composit mechan system real time circuit fundamant inform intens level long_term tool address commun control influenc
Brain Science
Doc 23: mechan technolog character isol level accomplish analysi control demonstr promis cellular protein target vivo yield connect direct
Brain Science
Doc 24: build design effect key manufactur materi mechan robust directli featur find obtain outcom practic result act environ
Cancer Treatment and Immunology
Doc 25: design innov key rout strategi algorithm choic comput framework outcom result time circuit encod inform long_term singl
Brain Science
Doc 26: challeng conduct mechan robust structur system scienc time collect environ inform reduc analysi commun basi creat target
Social Policy and Public Healthcare
Doc 27: build combin conduct effect effici effort evalu innov strategi sustain establish find scienc set social implement resourc
Social Policy and Public Healthcare
Doc 28: conduct effect effort innov life featur outcom problem scienc social character environ level protect util commun potenti
Social Policy and Public Healthcare
Doc 29: complex construct effect lower power outcom principl time inform intrins level limit state advantag control process stabl
Brain Science
Doc 30: decreas effect mechan product establish reduc block major role cellular class protein contribut gener respons mild produc
Brain Science

Figure 20: Sentences coloring of dominant topic in part of documents

567 **References**

- 568 [1] W. Magua, X. Zhu, A. Bhattacharya, A. Filut, A. Potvien, R. Leatherberry, Y.-G. Lee,
569 M. Jens, D. Malikireddy, M. Carnes, and A. Kaatz, “Are female applicants disadvantaged
570 in national institutes of health peer review? combining algorithmic text mining and qual-
571 itative methods to detect evaluative differences in r01 reviewers’ critiques,” *Journal of*
572 *Women’s Health*, vol. 26, 03 2017.
- 573 [2] D. van Dijk, O. Manor, and L. B. Carey, “Publication metrics and success on the
574 academic job market,” *Current Biology*, vol. 24, no. 11, pp. R516–R517, 2014. [Online].
575 Available: <https://www.sciencedirect.com/science/article/pii/S0960982214004771>
- 576 [3] D. Acuna, S. Allesina, and K. Kording, “Future impact: Predicting scientific success,”
577 *Nature*, vol. 489, pp. 201–2, 09 2012.
- 578 [4] S. Allesina, “Measuring nepotism through shared last names: The case of italian
579 academia,” *PLOS ONE*, vol. 6, no. 8, pp. 1–6, 08 2011. [Online]. Available:
580 <https://doi.org/10.1371/journal.pone.0021160>
- 581 [5] M. Hilbert, “Big data for development: A review of promises and challenges,”
582 *Development Policy Review*, vol. 34, no. 1, pp. 135–174, 2016. [Online]. Available:
583 <https://onlinelibrary.wiley.com/doi/abs/10.1111/dpr.12142>
- 584 [6] K. Kannan, P. Sekar, M. Sathik, and A. Assaf, “Stock market prediction and analysis using
585 naïve bayes,” 2016.
- 586 [7] E. Sarigöl, R. Pfitzner, I. Scholtes, A. Garas, and F. Schweitzer, “Predicting scientific
587 success based on coauthorship networks,” *EPJ Data Science*, 02 2014.
- 588 [8] M. Greenberg, B. Pardo, K. Hariharan, and E. Gerber, “Crowdfunding support tools:
589 Predicting success & failure,” 04 2013, pp. 1815–1820.
- 590 [9] T. Mitra and E. Gilbert, “The language that gets people to give: Phrases that predict
591 success on kickstarter,” 02 2014, pp. 49–61.
- 592 [10] B. Alberts, “Impact factor distortions,” *Science*, vol. 340, no. 6134, pp. 787–787, 2013.
593 [Online]. Available: <https://science.sciencemag.org/content/340/6134/787>
- 594 [11] *UK Research and Innovation GtR Database*, https://gtr.ukri.org/search/project?term=*&,
595 2021, [Accessed: August 20, 2021].
- 596 [12] *National Science Foundation Funded Project Database*, <https://www.nsf.gov/awardsearch/download.jsp>, 2021, [Accessed: August 20, 2021].

- 598 [13] *National Institutes of Health RePORTER Database*, <https://reporter.nih.gov/search/JbRgAtchAUKHOoW2X1vObQ/projects>, 2021, [Accessed: August 20, 2021].
- 600 [14] *European Research Council Funded Project Database*, <https://erc.europa.eu/projects-figures/project-database>, 2021, [Accessed: August 20, 2021].
- 602 [15] F. Ahmed and A. Nürnberg, “Evaluation of n-gram conflation approaches
603 for arabic text retrieval,” *Journal of the American Society for Information
604 Science and Technology*, vol. 60, no. 7, pp. 1448–1465, 2009. [Online]. Available:
605 <https://onlinelibrary.wiley.com/doi/abs/10.1002/asi.21063>
- 606 [16] S. W. Thomas, B. Adams, A. Hassan, and D. Blostein, “Studying software evolution using
607 topic models,” *Sci. Comput. Program.*, vol. 80, pp. 457–479, 2014.
- 608 [17] R. E. Hintzen, M. Papadopoulou, R. Mounce, C. Banks-Leite, R. D. Holt,
609 M. Mills, A. T. Knight, A. M. Leroi, and J. Rosindell, “Relationship between
610 conservation biology and ecology shown through machine reading of 32,000 articles,”
611 *Conservation Biology*, vol. 34, no. 3, pp. 721–732, 2020. [Online]. Available:
612 <https://onlinelibrary.wiley.com/doi/abs/10.1111/cobi.13435>
- 613 [18] *Foreign exchange rates and currency conversion JSON API*, <https://fixer.io/documentation>, 2021, [Accessed: August 20, 2021].
- 615 [19] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent dirichlet allocation,” *J. Mach. Learn.
616 Res.*, vol. 3, no. null, p. 993–1022, Mar. 2003.
- 617 [20] Z. Qiu, B. Wu, B. Wang, C. Shi, and L. Yu, “Collapsed gibbs sampling for latent dirichlet
618 allocation on spark,” in *Proceedings of the 3rd International Conference on Big Data,
619 Streams and Heterogeneous Source Mining: Algorithms, Systems, Programming Models
620 and Applications - Volume 36*, ser. BIGMINE’14. JMLR.org, 2014, p. 17–28.
- 621 [21] *Gensim: Topic Modelling for Humans*, https://radimrehurek.com/gensim/auto_examples/index.html#documentation, 2009, [Accessed: August 20, 2021].
- 623 [22] J. Chuang, M. E. Roberts, B. M. Stewart, R. Weiss, D. Tingley, J. Grimmer, and
624 J. Heer, “TopicCheck: Interactive alignment for assessing topic model stability,” in
625 *Proceedings of the 2015 Conference of the North American Chapter of the Association
626 for Computational Linguistics: Human Language Technologies*. Denver, Colorado:
627 Association for Computational Linguistics, May–Jun. 2015, pp. 175–184. [Online].
628 Available: <https://aclanthology.org/N15-1018>
- 629 [23] X. Wang, A. Fang, I. Ounis, and C. MacDonald, “Evaluating similarity metrics for latent
630 twitter topics,” in *ECIR*, 2019.

- 631 [24] N. Aletras and M. Stevenson, “Evaluating topic coherence using distributional semantics,”
632 03 2013, pp. 13–22.
- 633 [25] J. H. Lau, D. Newman, and T. Baldwin, “Machine reading tea leaves: Automatically
634 evaluating topic coherence and topic model quality,” in *Proceedings of the 14th Conference*
635 *of the European Chapter of the Association for Computational Linguistics*. Gothenburg,
636 Sweden: Association for Computational Linguistics, Apr. 2014, pp. 530–539. [Online].
637 Available: <https://aclanthology.org/E14-1056>
- 638 [26] D. Mimno, H. M. Wallach, E. Talley, M. Leenders, and A. McCallum, “Optimizing semantic
639 coherence in topic models,” in *Proceedings of the Conference on Empirical Methods in*
640 *Natural Language Processing*, ser. EMNLP ’11. USA: Association for Computational
641 Linguistics, 2011, p. 262–272.
- 642 [27] G. Bouma, “Normalized (pointwise) mutual information in collocation extraction,” 2009.
- 643 [28] M. Röder, A. Both, and A. Hinneburg, “Exploring the space of topic coherence measures,”
644 in *Proceedings of the Eighth ACM International Conference on Web Search and Data*
645 *Mining*, ser. WSDM ’15. New York, NY, USA: Association for Computing Machinery,
646 2015, p. 399–408. [Online]. Available: <https://doi.org/10.1145/2684822.2685324>
- 647 [29] B. Fitelson, “A probabilistic theory of coherence,” *Analysis*, vol. 63, no. 3, pp. 194–199,
648 2003. [Online]. Available: <http://www.jstor.org/stable/3329309>
- 649 [30] D. Newman, J. H. Lau, K. Grieser, and T. Baldwin, “Automatic evaluation of topic co-
650 herence,” in *Human Language Technologies: The 2010 Annual Conference of the North*
651 *American Chapter of the Association for Computational Linguistics*, ser. HLT ’10. USA:
652 Association for Computational Linguistics, 2010, p. 100–108.
- 653 [31] M. Papadopoulou, “The cultural evolution of evolution,” 2017.
- 654 [32] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu, “Lightgbm:
655 A highly efficient gradient boosting decision tree,” in *Proceedings of the 31st International*
656 *Conference on Neural Information Processing Systems*, ser. NIPS’17. Red Hook, NY,
657 USA: Curran Associates Inc., 2017, p. 3149–3157.
- 658 [33] J. H. Friedman, “Greedy function approximation: A gradient boosting machine.”
659 *The Annals of Statistics*, vol. 29, no. 5, pp. 1189 – 1232, 2001. [Online]. Available:
660 <https://doi.org/10.1214/aos/1013203451>
- 661 [34] P. Li, “Robust logitboost and adaptive base class (abc) logitboost,” in *Proceedings of the*
662 *Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, ser. UAI’10. Arlington,
663 Virginia, USA: AUAI Press, 2010, p. 302–311.

- 664 [35] C. J. Burges, “From ranknet to lambdarank to lambdamart: An overview,” Tech. Rep.
665 MSR-TR-2010-82, June 2010. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/from-ranknet-to-lambdarank-to-lambdamart-an-overview/>
- 667 [36] *scikit-learn: Machine Learning in Python.*, <https://scikit-learn.org/stable/about.html>,
668 2007, version 0.24.2.
- 669 [37] Y. Wang and T. Wang, “Application of improved lightgbm model in blood glucose prediction,” *Applied Sciences*, vol. 10, p. 3227, 05 2020.
- 671 [38] D. Kim and A. Oh, “Topic chains for understanding a news corpus,” 02 2011, pp. 163–176.
- 672 [39] *Hungarian algorithm to solve the linear sum assignment problem*, <https://docs.scipy.org/doc/scipy/reference/tutorial/optimize.html#assignment-problems>, 2021, version 1.7.1.
- 674 [40] H. W. Kuhn, “The hungarian method for the assignment problem,” *Naval Research Logistics Quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/nav.3800020109>
- 677 [41] S. Miyamoto, R. Abe, Y. Endo, and J.-i. Takeshita, “Ward method of hierarchical clustering for non-euclidean similarity measures,” in *2015 7th International Conference of Soft Computing and Pattern Recognition (SoCPaR)*, 2015, pp. 60–63.
- 680 [42] J. H. W. Jr., “Hierarchical grouping to optimize an objective function,” *Journal of the American Statistical Association*, vol. 58, no. 301, pp. 236–244, 1963. [Online]. Available: <https://www.tandfonline.com/doi/abs/10.1080/01621459.1963.10500845>
- 683 [43] C. Sievert and K. Shirley, “Ldavis: A method for visualizing and interpreting topics,” 06 2014.
- 685 [44] M. A. Taddy, “On estimation and selection for topic models,” 2011.
- 686 [45] M. Jordan, Z. Ghahramani, T. Jaakkola, and L. Saul, “An introduction to variational methods for graphical models,” *Machine Learning*, vol. 37, pp. 183–233, 01 1999.