

*Article*

# Forecasting Emerging Topics in Research and Development in Pharmaceutical Industry

MahParsa<sup>1,t,‡</sup>,<sup>1</sup> mahboobeh.parsapoor@mail.mcgill.ca

\* Correspondence: mahboobeh.parsapoor@mail.mcgill.ca

† Current address: McGill University

Academic Editor: name

Version March 23, 2019 submitted to Journal Not Specified

- <sup>1</sup> **Keywords:** Bibliometrics Tools, Clustering Algorithms, Machine Learning Tools, Text Analysis, Topic Modeling, Topic Trends
- 

## <sup>3</sup> 1. Abstract

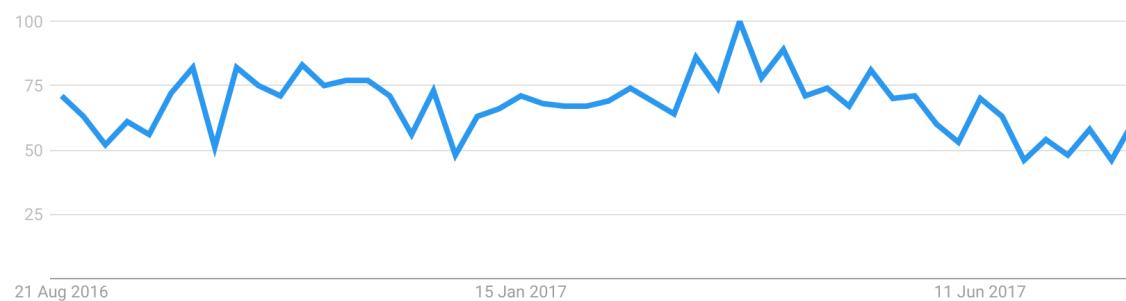
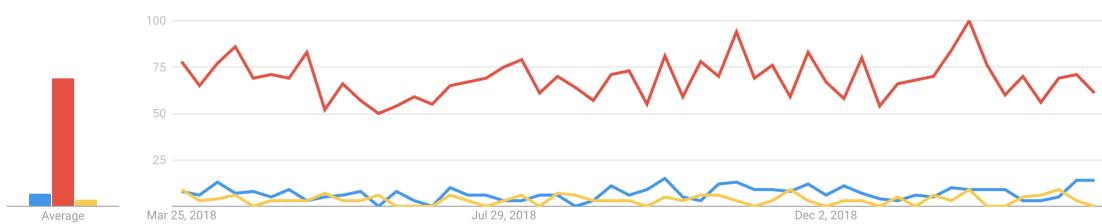
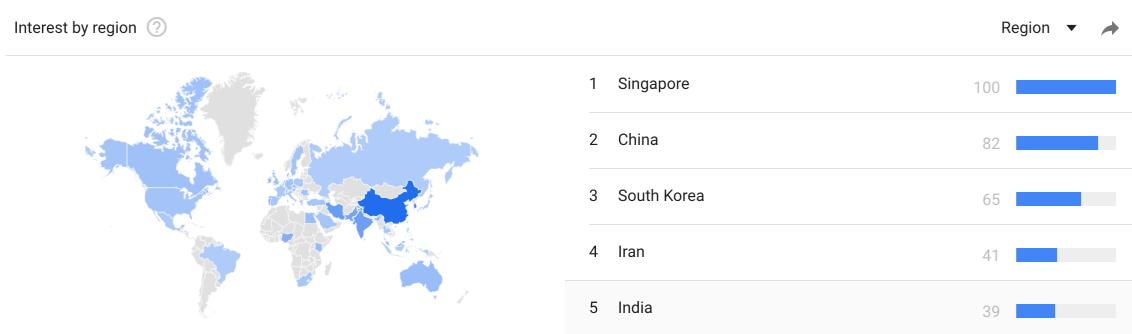
<sup>4</sup> Predicting the future trends in topics can impact on research projects in the academy by helping  
<sup>5</sup> grant agencies recognize emerging topics and assign future grants to scientists that are working on  
<sup>6</sup> those topics. Identifying future trends in a topic can also help companies and start ups invest in  
<sup>7</sup> emerging themes to change the game in their fields. This paper suggests a new strategy to identify  
<sup>8</sup> developing drifts.

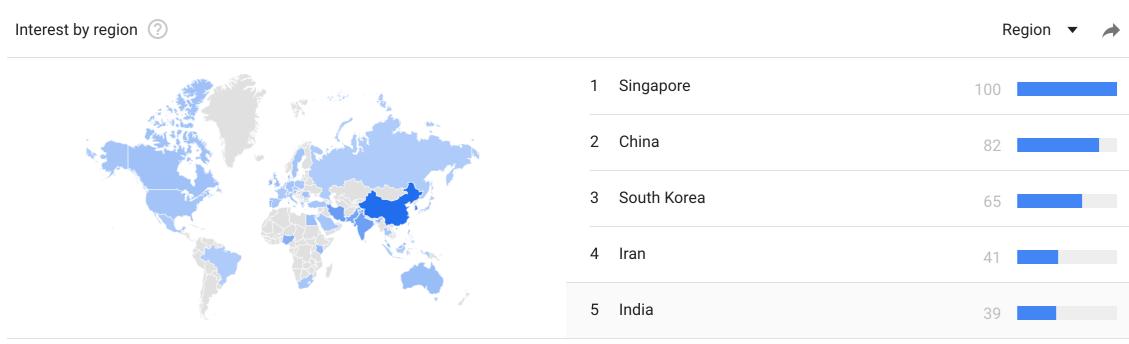
## <sup>9</sup> 2. Introduction

<sup>10</sup> Discovering research trends in pharmaceutical industry can directly influence governmental  
<sup>11</sup> policies, healthcare services, educational organizations, research institutes. Thus, it is entirely necessary  
<sup>12</sup> to predict what is the future trends in pharmaceutical industries. For example, discovering the  
<sup>13</sup> future trends in biosensors<sup>1</sup> not only assist pharmaceutical and biomedical researchers in grasping  
<sup>14</sup> information of emerging drifts in this field of bio electronics but also help institutions and grant  
<sup>15</sup> agencies to assign future grants to scientists that are working on emerging topics [2]. Online tools  
<sup>16</sup> such as "Google Trends" can be utilised to forecast the latest trending topics or determine (Figure  
<sup>17</sup> 1) the popularity of research themes (Figure 5) in the different fields of science and technology, e.g.,  
<sup>18</sup> biosensors. They can provide a comparison interest to different terms such as "Dendritic cells" (blue  
<sup>19</sup> line), "Biosensors" (yellow line) and "Cancer Treatment" (red line) Cell" see Figure 3. These tools cannot  
<sup>20</sup> detect emerging topics by just analysing scholarly articles. Thus, it is necessary to use a method to  
<sup>21</sup> analyze documents, discover topics and predict the future topics.

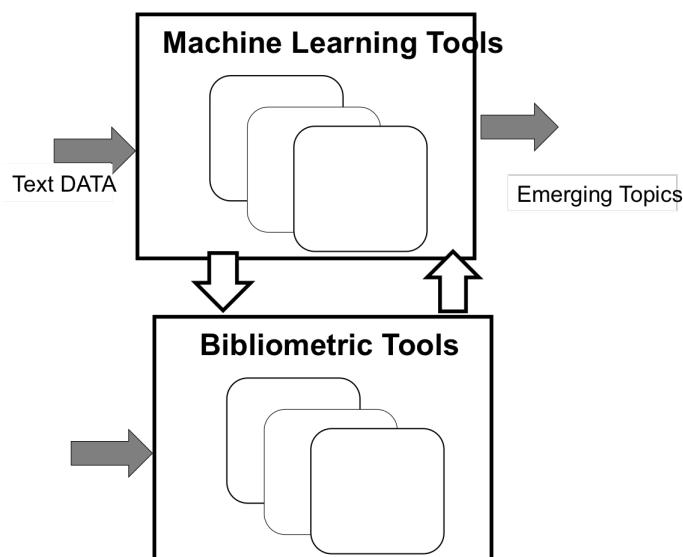
<sup>22</sup> Text analysis methods includes techniques that are useful for processing text that are obtained  
<sup>23</sup> from equipment logs, blogs, articles and social media such as Facebook and twitter. Using text  
<sup>24</sup> analysis we can extract useful information from text, remove unnecessary information from text and  
<sup>25</sup> represent text using numerical values. In other words, text analysis tools help us to develop statistical  
<sup>26</sup> tools. One of main class of text analysis is machine learning methods such as clustering algorithms,

<sup>1</sup> The term "biosensors", which firstly introduced by Griffin and Nelson, refers an analytical device which converts a biological response into an electrical signal [1]. A typical biosensor has two components, biological elements such as enzyme, antibody, cell, tissue and a physical transducer such as voltammetric, amperometric, conductometric, spectrophotometric

**Figure 1.** The number of interest of biosensors over time**Figure 2.** Different locations that the terms "biosensors" have been most noted during the specified time frame**Figure 3.** Different locations that the terms "Dendritic cells", "Biosensors" and "Cancer Treatment" Cell have been most noted during the specified time frame**Figure 4.** Different locations that the terms "biosensors" have been most noted during the specified time frame



**Figure 5.** Different locations that the terms "biosensors" have been most noted during the specified time frame



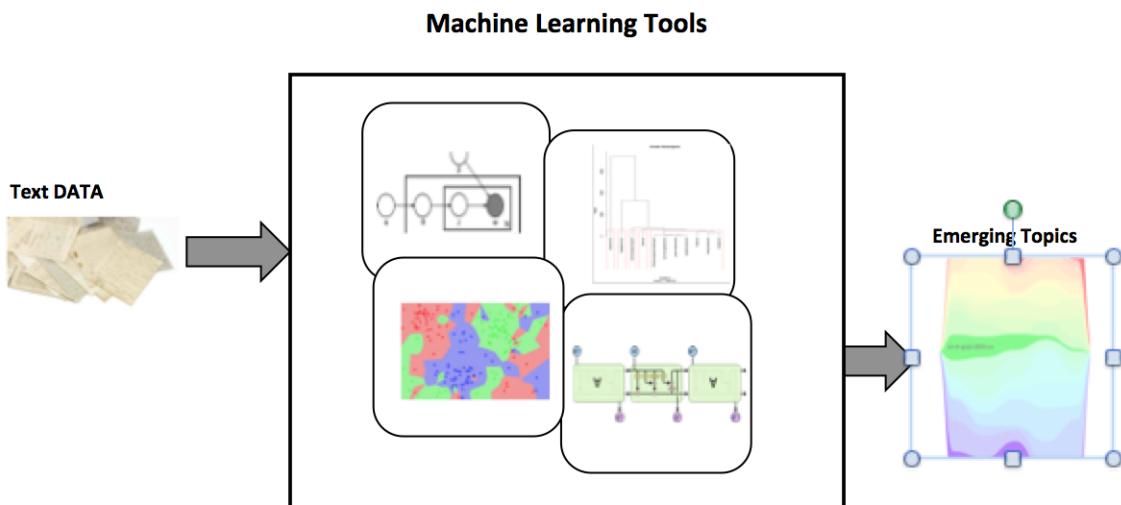
**Figure 6.** The general structure of framework to forecast emerging topics; The framework receives text documents (an abstract of publication) and discover emerging topics. Note that, MLT and BT modules utilize a bidirectional connections to exchange data and useful information

27 LDA and LSA. It should be noted that outputs of text analytical tools must be used with other machine  
 28 learning tools to provide meaning full results. This paper introduces a framework that merges machine  
 29 learning algorithms and bibliometrics to discover emerging topics and forecasts the latest trending of  
 30 emerging topics. The framework automatically detects emerging topics from a set of text documents  
 31 and estimate evolving topics. Figure 6 describes the framework with its two main modules that are  
 32 named "Machine Learning Tools" ("MLT") and "Bibliometrics Tools"("BT").

### 33 3. Framework

34 The MLT module in the framework integrates various statistical-based and bio-inspired machine  
 35 learning tools to do the following tasks.

- 36 1) Select useful terms
- 37 2) Find meaningful topics and assign text documents to highest frequency terms of topics
- 38 3) Categorise text documents
- 39 5) Predict bibilmetric indicators
- 40 6) Integrate data
- 41 7) Identify emerging topics
- 42 8) Predict emerging topic trend



**Figure 7.** An example of the MLT module

43 This paper is focusing on the three first steps of the above tasks. The first phase is to select useful  
 44 terms to gather documents from Blogs, Papers, websites. For example, we selected some keywords  
 45 such as "biosensors", "Continuous glucose monitoring", "glucose measurement", "Hypoglycemia",  
 46 "Insulin", "medical devices", "metabolic", "pharmacological target", "Protein", "real-time applications",  
 47 "wearable sensors" to gather documents related to biosensors. At the second phase, we can run natural  
 48 processing language tools and topic modeling tools (e.g., LDA) to find terms with highest frequency  
 49 and topics that are hidden in documents. The output of LDA is topics and terms that are related to each  
 50 topic. Thus, for each topic, we can determine one term with the highest probability of belonging to  
 51 topics. The third step categorises text documents to most important topics using a clustering algorithm.  
 52 The MLT module is an ensemble method<sup>2</sup>, that integrates multiple algorithms (such as topic modelling,  
 53 ensemble classifiers, clustering algorithms, multivariate time-series prediction models such as adaptive  
 54 neuro fuzzy inference system (ANFIS) and Brain Emotional Learning inspired Model (BELiMs) [3]  
 55 and [4], multivariate multiple regression) to analyse texts, detect emerging topics and determine the  
 56 evolution of emerging topics.

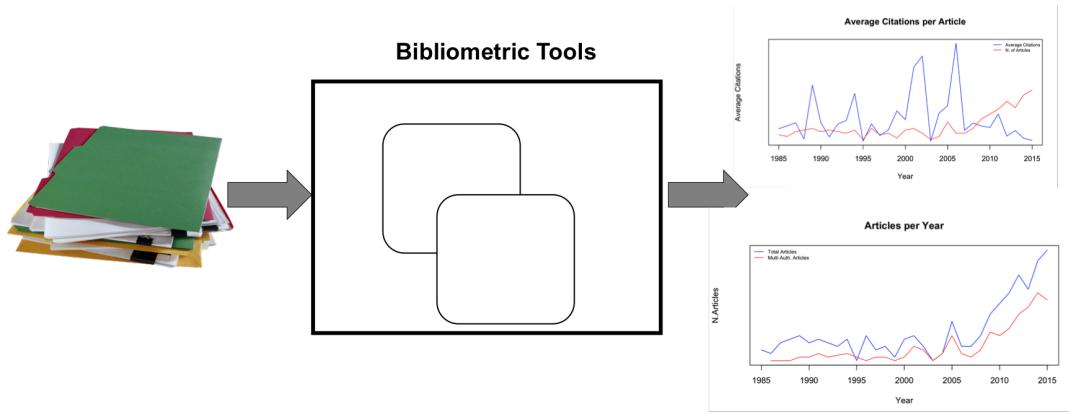
57 For example, an MLT module (see Figure 7) can encompass latent dirichlet allocation (LDA)  
 58 as a topic modelling algorithm, hierarchical clustering algorithm (HCA) as a clustering algorithm,  
 59 K nearest neighbour (KNN) as a classifier, and long short-term memory (LSTM)<sup>3</sup> as a time-series  
 60 prediction model) [5].

61 Another part of this framework is the BT (Bibliometrics Tools) module to analyse metadata of  
 62 academic literature related to most important topics. For example, "bibliometrix" (an R package for  
 63 bibliometric study) can interpret meta-data of scholarly articles related to topics to indicators (e.g., the  
 64 number of publications) that would be sent to the MLT module to construct multivariate time series  
 65 data sets.

66 Bibliometric tools (see Figure 8) implement quantitative analysis of academic literature to discover  
 67 the impact of research fields, the impression of a set of researchers, or the impact of particular articles.

<sup>2</sup> Ensemble learning is the process to combine multiple learning models such as classifiers, prediction to solve a complex problem. The main advantage of an ensemble approach is to increase the accuracy

<sup>3</sup> LSTM has been inspired by human memory and learning systems and can learn about long-term dependencies of values of a variable in the past to predict future values of variable



**Figure 8.** The BT module that encompasses bibliometrix

68 The performance and effectiveness of proposed framework would be evaluated on a variety of  
 69 text datasets (e.g., Scopus); the obtained results would be compared with other methodologies that  
 70 most are based on bibliometric tools.

#### 71 4. Latent Dirichlet Allocation (LDA) of MLTs

72 This section explains Latent Dirichlet Allocation (LDA) as a well-known topic modelling. Topic  
 73 modelling algorithms mostly are statistical methods and belongs to a group of algorithms that are  
 74 named generative probabilistic modeling methods. Generative probabilistic modelling methods are are  
 75 based on a probabilistic generative procedure to find hidden variables underlying the observed  
 76 data. Then they utilize a joint probability distribution over both the observed and hidden variables  
 77 and compute the posterior distribution (conditional distribution) of the hidden variables given the  
 78 observed variables. aim to determine a probabilistic model of a corpus that can assign high probability  
 79 to members of the corpus (i.e unseen documents).

80 One well-known generative probabilistic topic modeling algorithm is Latent Dirichlet Allocation  
 81 (LDA) which have been excellent performance in finding topics in a collection of documremts. In LDA,  
 82 Latent variables are topics, while apparent variables are words. LDA receives words as an input vector  
 83 and generates topics which are probability distribution over words based on a generative process.  
 84 LDA uses a joint probability distribution over both the observed and hidden random variables and  
 85 compute the posterior distribution (conditional distribution) of the hidden variables given the observed  
 86 variables. The fundamental assumption of LDA is that documents can be assigned to multiple topics.  
 87 Another assumption is that topics are hidden variables and words in documents are apparent variables.  
 88 Thus, LDA performs a generative process by receiving words as an input vector to provide topics  
 89 which are probability distribution over words. In addition to the above assumption, the following  
 90 assumption should be considered:

- 91 1. The number of topics,  $k$ , is known.
  - 92 2. The size of vocabulary,  $V$ , is known.
  - 93 3. Each topic can be defined as a Dirichlet distribution over words.
  - 94 4. Each document has two probabilistic distributions, the first one is a Dirichlet distribution over  
 95 topics to choose a topic proportion. The second one is mutlinomial disterbution that be used to  
 96 choose a word for the document.
  - 97 5. The number of words,  $N$  is choosen using a poisson distribution as  $\sim (\xi)$ .
  - 98 6. Choose  $\theta \sim (\alpha)$ .
  - 99 7. For each of the  $N$  words  $w_n$ :
- 100 (a) Choose a topic  $z_n \sim (\theta)$ .

101 (b) Choose a word  $w_n$  from  $(w_n z_n, \beta)$ , a multinomial probability conditioned on the topic  $z_n$ .

102 It should be noted that, LDA considers the dimensional of the topic variable  $z$  and the dimensional  
 103 of the Dirichlet distribution  $k$  are known; moreover, LDA uses a matrix as  $\beta$  with  $k \times V$  where  
 104  $\beta_{ij} = p(w^j = 1 | z^i = 1)$ , the matrix is named as word probability.

A  $k$ -dimensional Dirichlet random variable  $\theta$  can take values in the  $(k - 1)$ -simplex (a  $k$ -vector  $\theta$  lies in the  $(k - 1)$ -simplex if  $\theta_i \geq 0, \sum_{i=1}^k \theta_i = 1$ ), and has the following probability density on this simplex:

$$(1) \quad (\theta|\alpha) = \frac{\Gamma\left(\sum_{i=1}^k \alpha_i\right)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \dots \theta_k^{\alpha_k-1},$$

105 where the parameter  $\alpha$  is a  $k$ -vector with components  $\alpha_i > 0$ , and where  $\Gamma(x)$  is the Gamma  
 106 function. The Dirichlet is a convenient distribution on the simplex — it is in the exponential family,  
 107 has finite dimensional sufficient statistics, and is conjugate to the multinomial distribution. In  
 108 inference-and-estimation, these properties will facilitate the development of inference and parameter  
 109 estimation algorithms for LDA.

Given the parameters  $\alpha$  and  $\beta$ , the joint distribution of a topic mixture  $\theta$ , a set of  $N$  topics, and a set of  $N$  words is given by:

$$(2) \quad (\theta, \alpha, \beta) = (\theta|\alpha) \prod_{n=1}^N (z_n \theta)(w_n z_n, \beta),$$

where  $(z_n \theta)$  is simply  $\theta_i$  for the unique  $i$  such that  $z_n^i = 1$ . Integrating over  $\theta$  and summing over  $z$ , we obtain the marginal distribution of a document:

$$(3) \quad p(\alpha, \beta) = \int p(\theta|\alpha) \left( \prod_{n=1}^N \sum_{z_n} p(z_n \theta) p(w_n z_n, \beta) \right) d\theta.$$

Finally, taking the product of the marginal probabilities of single documents, we obtain the probability of a corpus:

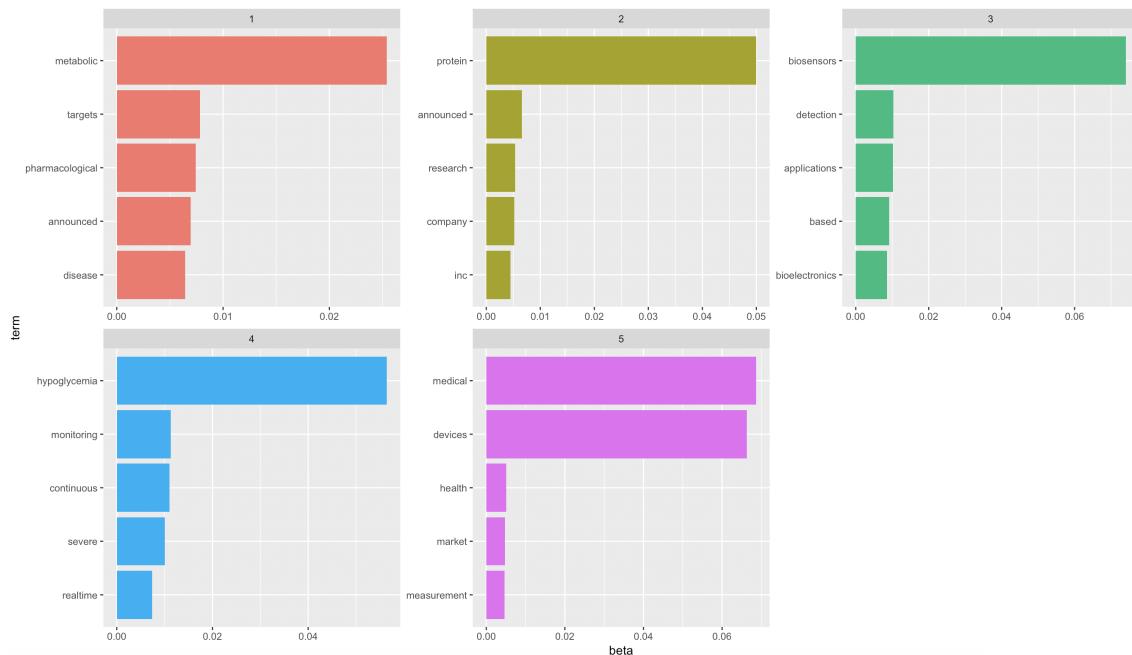
$$(4) \quad p(\alpha, \beta) = \prod_{d=1}^M \int p(\theta_d|\alpha) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} \theta_d) p(w_{dn} z_{dn}, \beta) \right) d\theta_d.$$

## 110 5. Clustering Algorithms of MLTs

111 The main aim of clustering methods is to categorize unstructured documents into a list of  
 112 meaningful groups [6]. Clustering algorithms can divide data object into separated groups [6]. Each  
 113 group encompasses similar objects; while objects between two groups must be different. In other words,  
 114 similar objects could be partitioned to a cluster, while patterns that are not alike to each other should be  
 115 distributed in separated clusters. The similarity measurement is defined using a clustering algorithm  
 116 A clustering algorithm needs a clear, implementable similarity rule. Defining different rules lead to  
 117 developing a different variation of clustering algorithms. Two main categories of clustering algorithms  
 118 are partitional clustering methods and hierarchical clustering methods. Partitional clustering algorithm  
 119 partitions data between separated groups which have not any common members, so each object only  
 120 belongs to one group. In contrast, hierarchical clustering provides a nested tree-like framework to  
 121 partition data. It should be noted that a member in lower clusters could be considered as a member  
 122 of upper clusters in the hierarchical structure. Clustering algorithms [7] as a group of unsupervised  
 123 learning algorithm have applications in text analysis and organizing documents.

### 124 5.0.1. K-means Clustering Algorithm

125 K-means clustering is an iterative unsupervised machine learning and has the goal to partition  
 126 objects into  $K$  separate groups. In the following, I have explained the steps of a K-means algorithm.



**Figure 9.** Generated five topics for a corpus with documents related to bio-sensors. The following terms, "Metabolic", "Protein", "Medical", "biosensors" and "Hypoglycemia" have highest probability for topics.

- 127 1. Determination Centroids which means with K cluster centres randomly selected.
- 128 2. Partitioning which means for each pattern cluster centre is found, and that pattern is considered from that cluster.
- 129 3. Replacement: Each centroid is replaced by average of data points in its partition
- 130 4. Iteration: The following step is repeated until convergence criteria is satisfied.

132 Step1:  $x_i = (x_{i1}, \dots, x_{ip})$ :

133 Step2: If centroids are  $m_1, m_2, \dots, m_k$ , and partitions are

134  $c_1, c_2, \dots, c_k$ , then one can show that K-means converges to a *local* minimum of

$$\sum_{k=1}^K \sum_{i \in c_k} \|x_i - m_k\|^2 \quad \text{Euclidean distance}$$

## 135 6. Experimental Results

136 The experimental results of this section are focusing on text analysis of two types of text data sets.  
 137 The first data set is realted to 6 documents that are gathered using the following keywords: "biosensors",  
 138 "Continuous glucose monitoring", "glucose measurement", "Hypoglycemia", "Insulin", "Medical  
 139 devices", "Metabolic", "Pharmacological target", "Protein", "Real-time Applications", "Wearable sensors".  
 140 The second data set is documents that have been extracted by using "Dendritic cells" as keywords.

### 141 6.1. Topic Modeling Results

### 142 6.2. First Data Set

143 Figures 9, 10 and 11 show different 5, 10 and 15 topics have been generated by using LDA models.  
 144 As it can be observed from the above figures, topic modeling algorithms provide us information  
 145 about words with highest probability. The information can be used to predict the future topics.

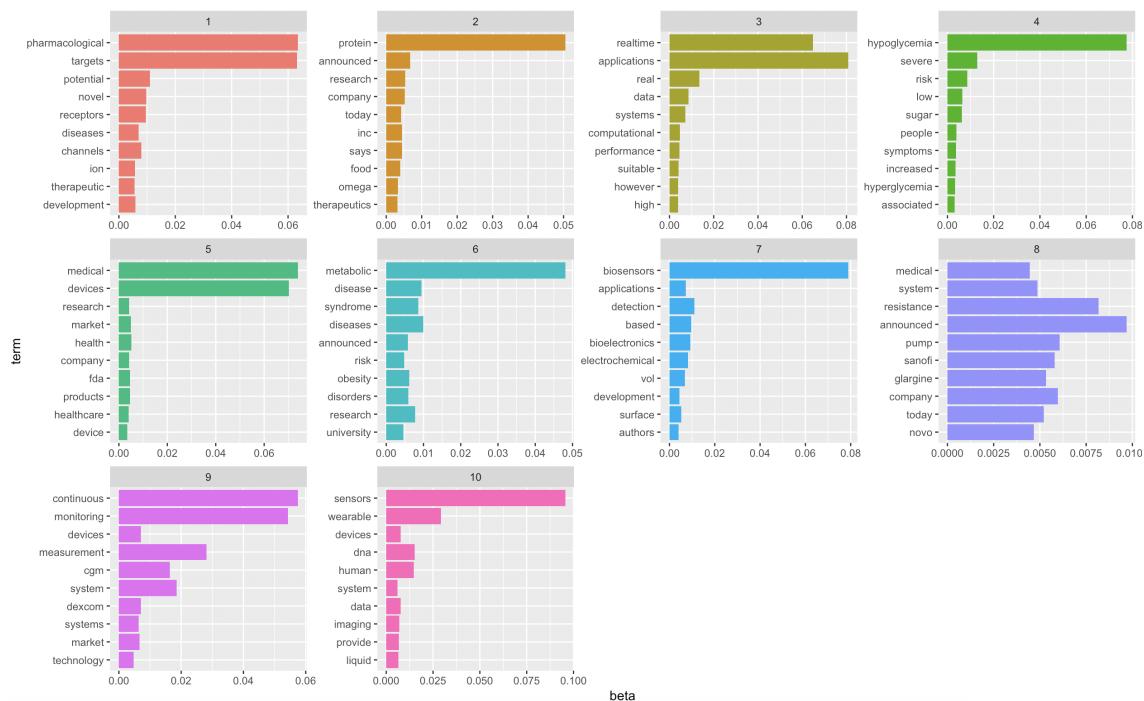


Figure 10. Generated ten topics for a corpus with documents related to bio-sensors

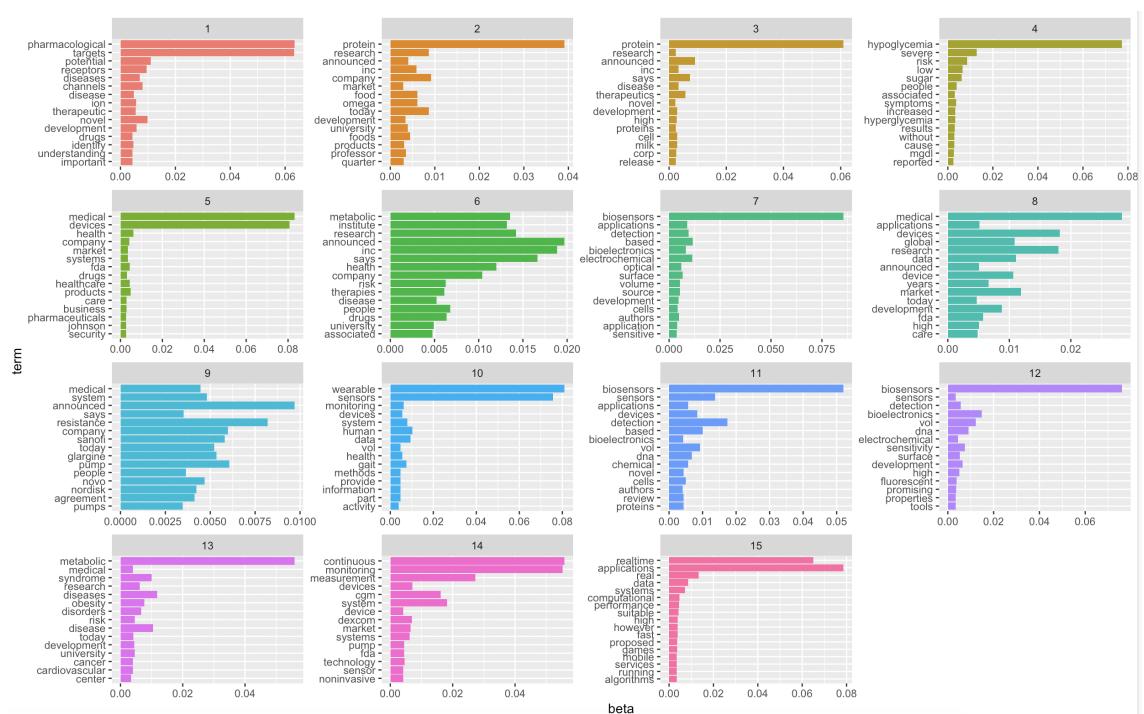
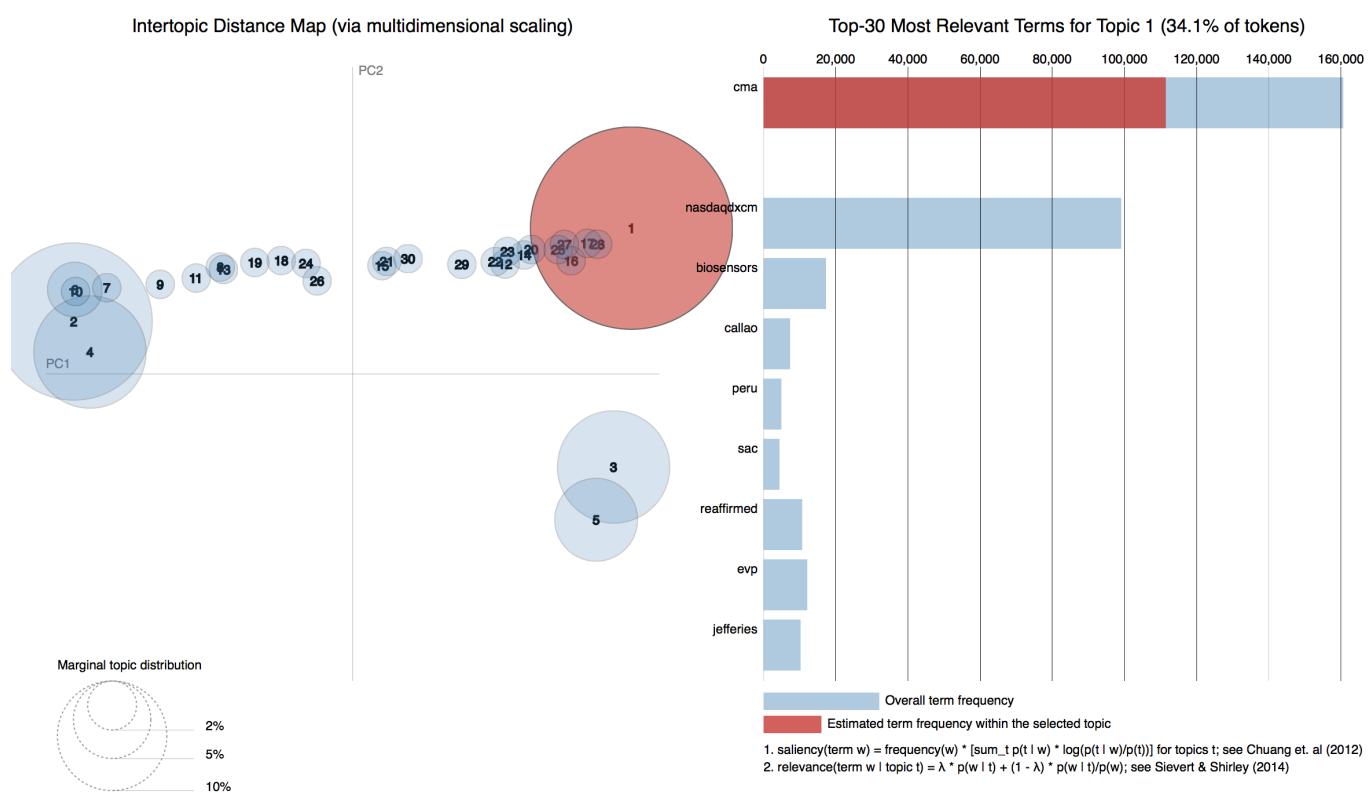
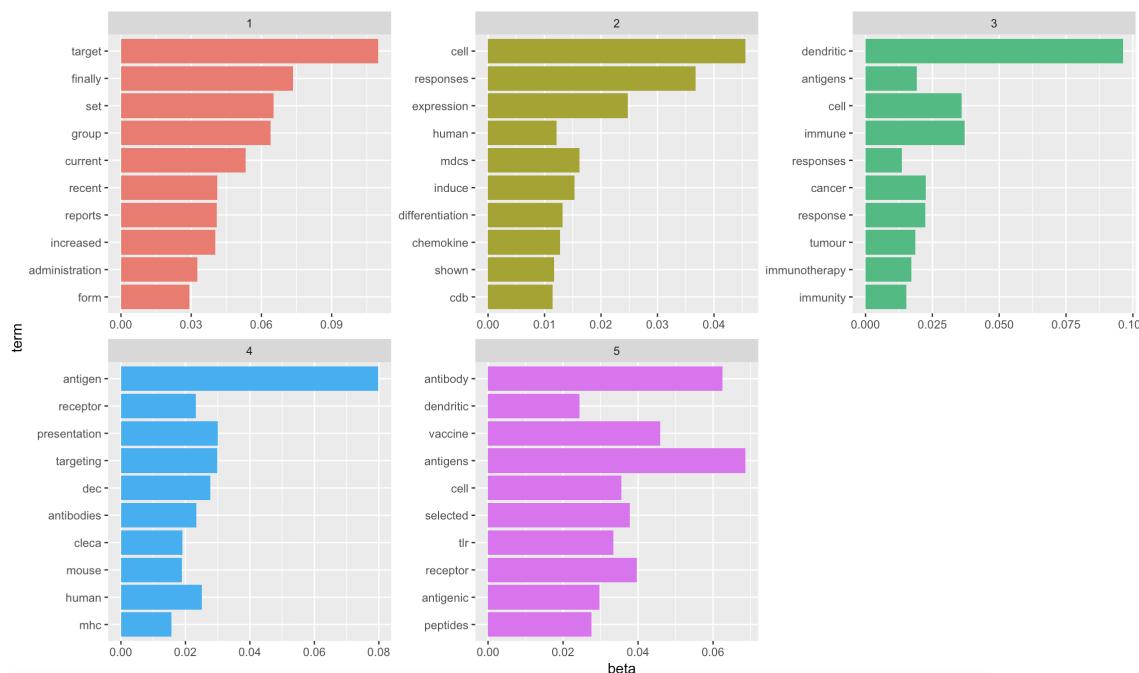


Figure 11. Generated fifteen topics for a corpus with documents related to bio-sensors

**Figure 12.** A visualization of 30 topics selected by LDA



**Figure 13.** Generated five topics for a corpus with documents related to Dendritic Cells

146      Figure 12 presents a nice visualization using "LDAvis" [8] of 30 topics that obtained from our  
 147      documents.

### 148      6.3. second Data set

149      The second example of this part is related to documents that are focusing on Dendritic cells<sup>4</sup>.  
 150      Figures ??, ?? and ?? show different 5, 10 and 15 topics have been generated by using LDA models.

### 151      6.4. Clustering Results

152      This section show how clustering algorithms can partition documents to clusters. It should be  
 153      mentioned we apply clustering algorithms on the documents and also the obtained results of topic  
 154      modeling.

#### 155      6.4.1. The First Data Set

156      Figures 16, 18 and ?? show the hierarchical clusters obtained from applying hirachical clustering  
 157      algorithms on the documents that are related to biosensors.

158      Figures 19, 20, 21, 22 and 23 show the hierarchical clusters obtained from applying K-means  
 159      clustering algorithms on the documents that are related to biosensors.

#### 160      6.4.2. The Second Data Set

161      This section show how clustering algorithms can partition documents related to Dendritic cells.  
 162      Figures 24, 25 and 26 show the hierarchical clusters obtained from applying hierarchical clustering  
 163      algorithms on those documents.

<sup>4</sup> Dendritic cells (DCs) can be classified as antigen-presenting cells (accessory cells) of the mammalian immune system. These cells are responsible for processing antigen material and offering it on the cell surface to the T cells of the immune system. In other words, DCs can play a role to exchange message be

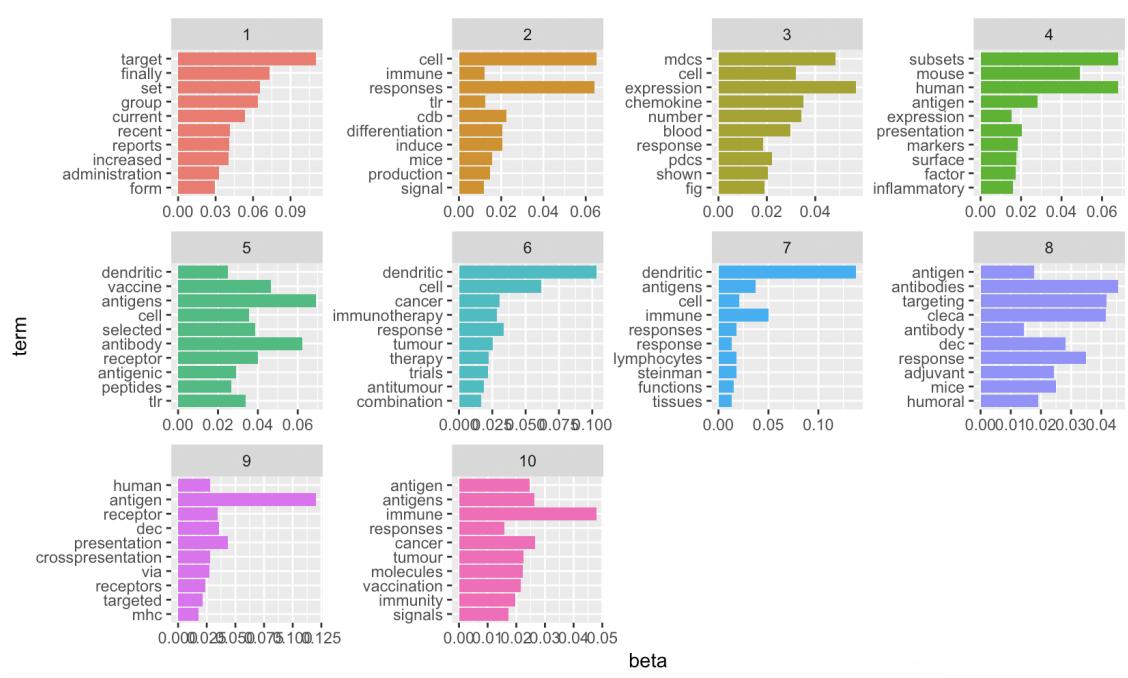


Figure 14. Generated ten topics for a corpus with documents related to Dendritic Cells

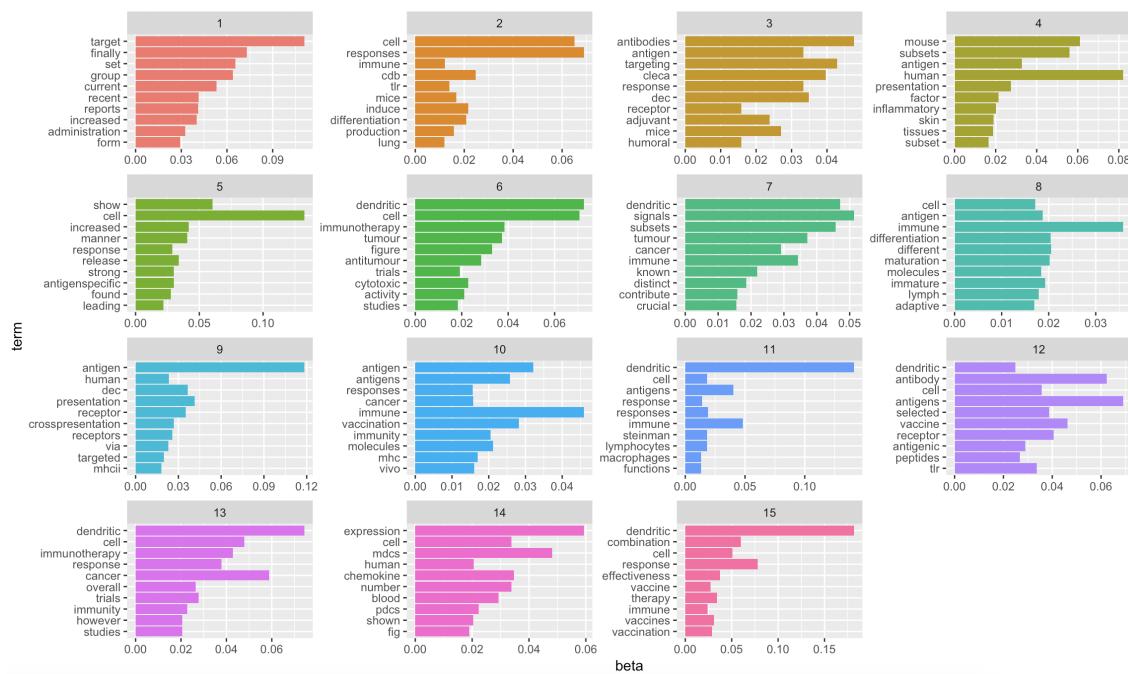
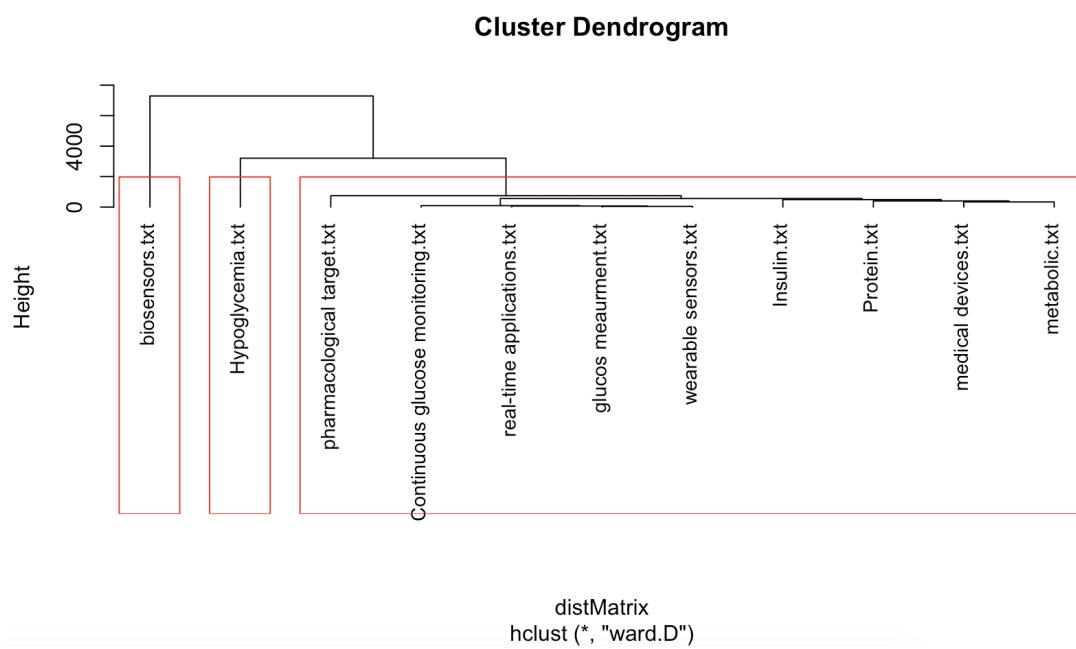
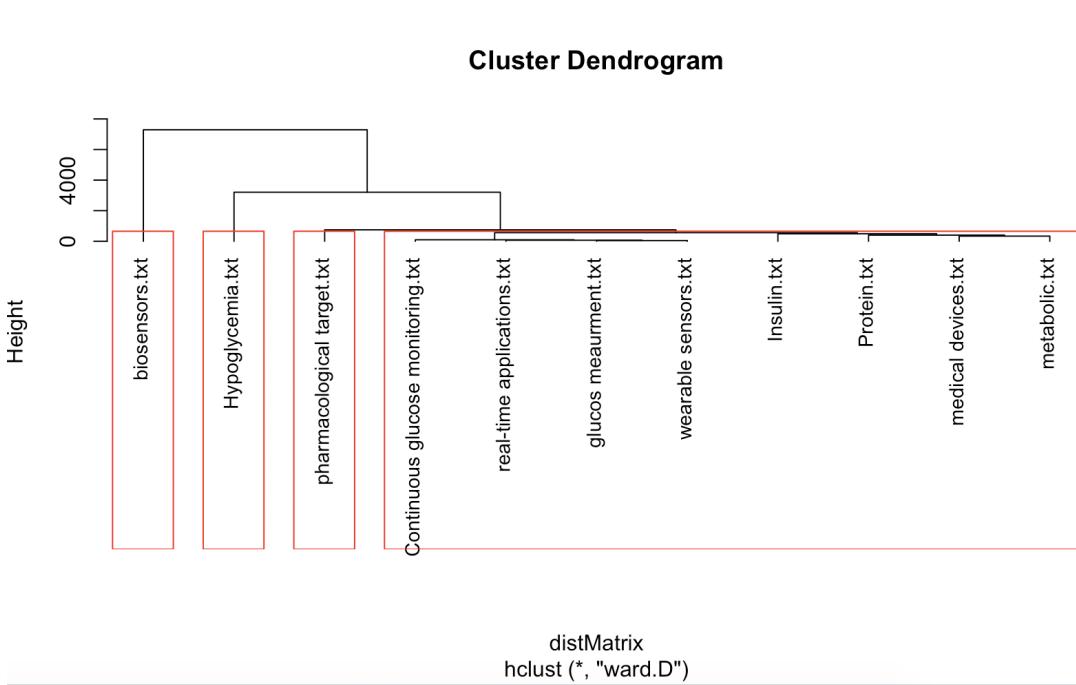


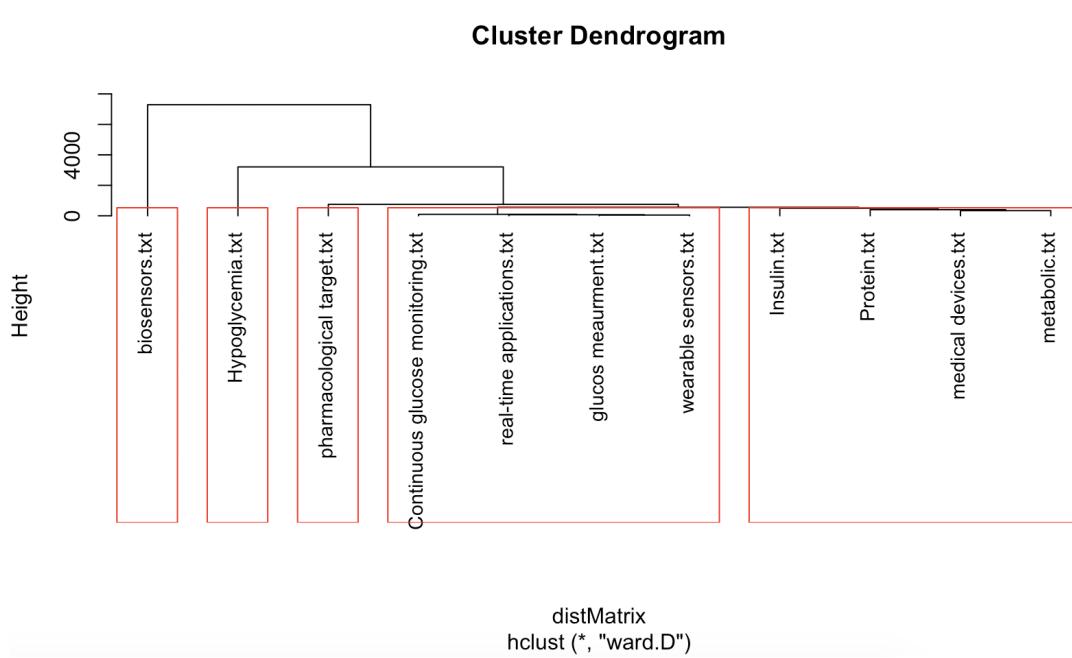
Figure 15. Generated 15 topics for a corpus with documents related to Dendritic Cells



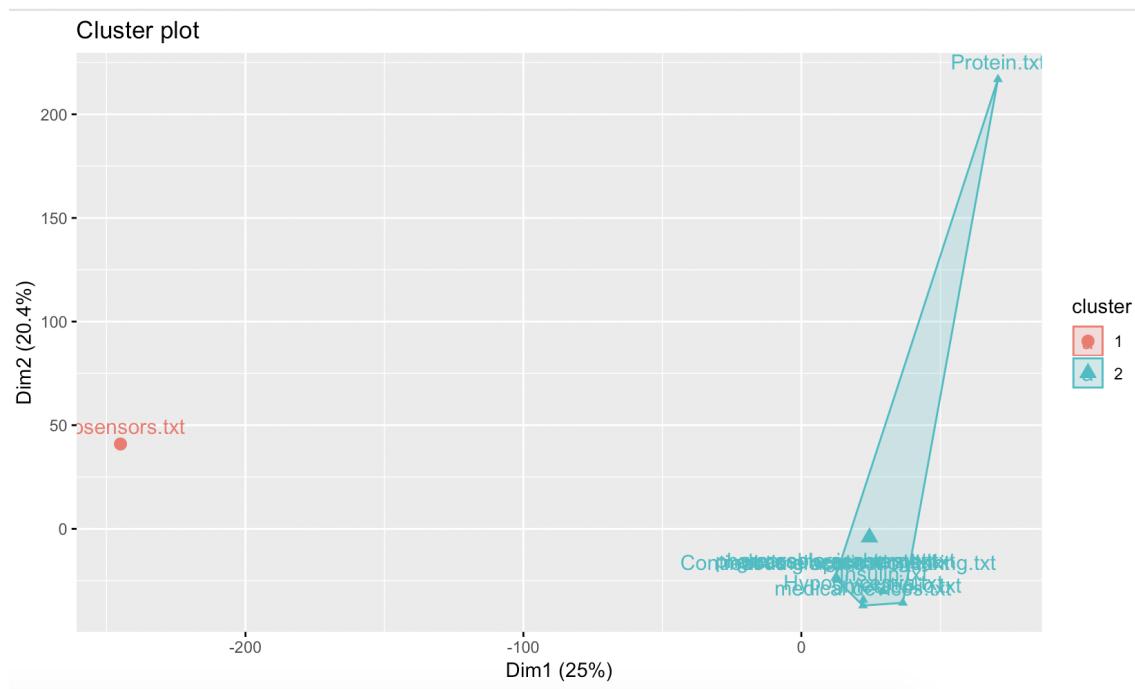
**Figure 16.** Using hierarchical clustering algorithm to have 3 cluster for 18 documents



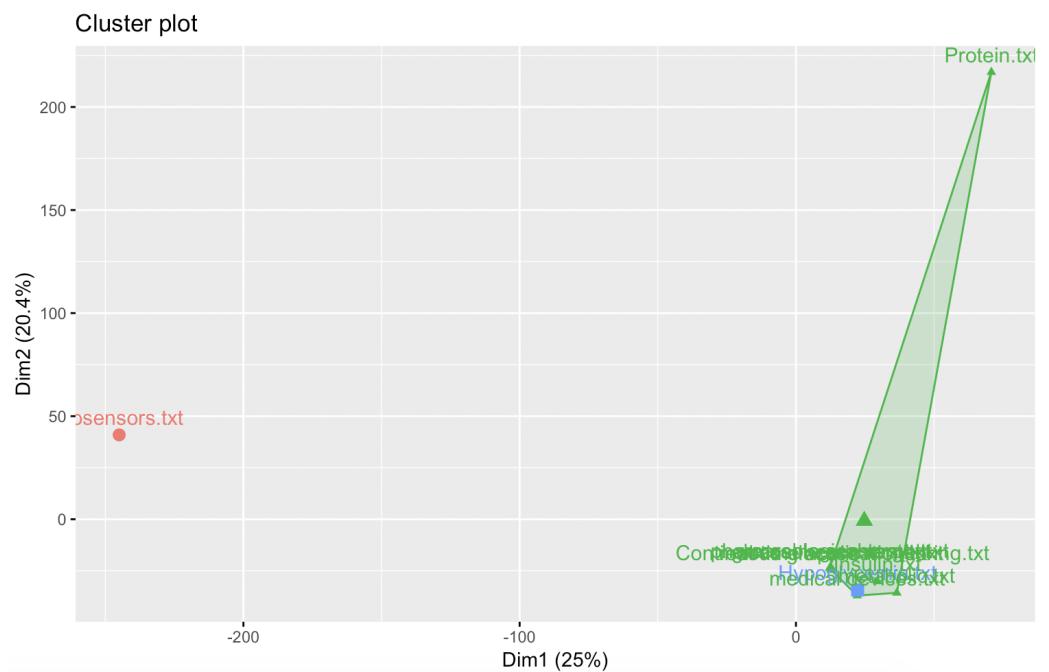
**Figure 17.** Using hierarchical clustering algorithm to have 4 clusters for 18 documents



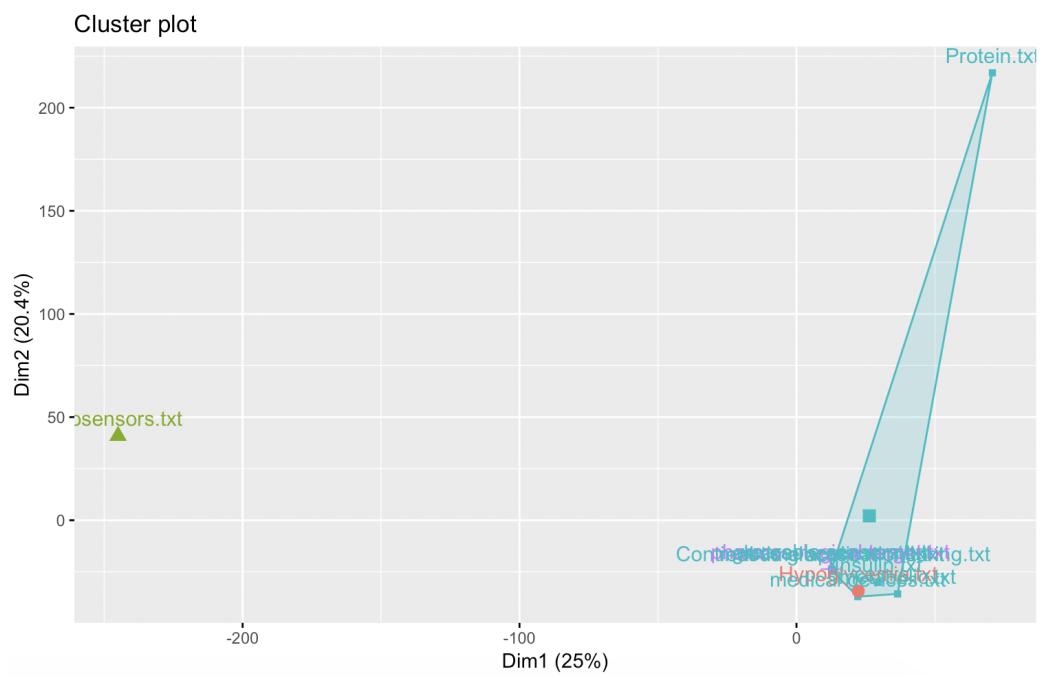
**Figure 18.** Using hierarchical clustering algorithm to have 4 clusters for 18 documents



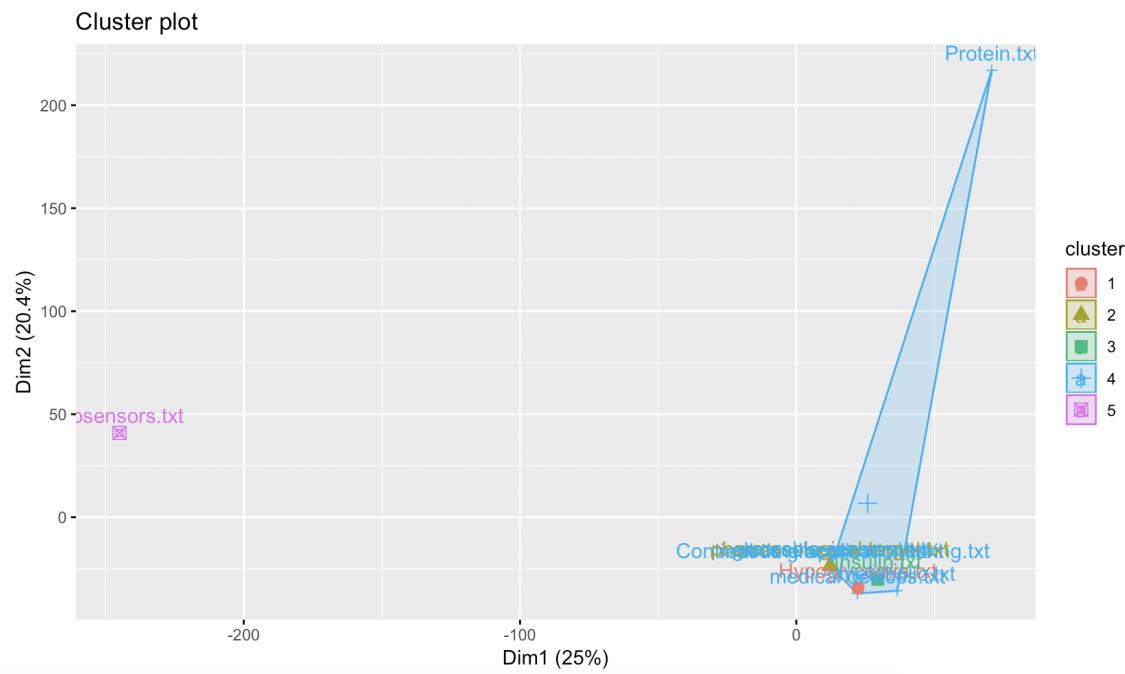
**Figure 19.** Using hierarchical clustering algorithm to have 2 cluster for 18 documents



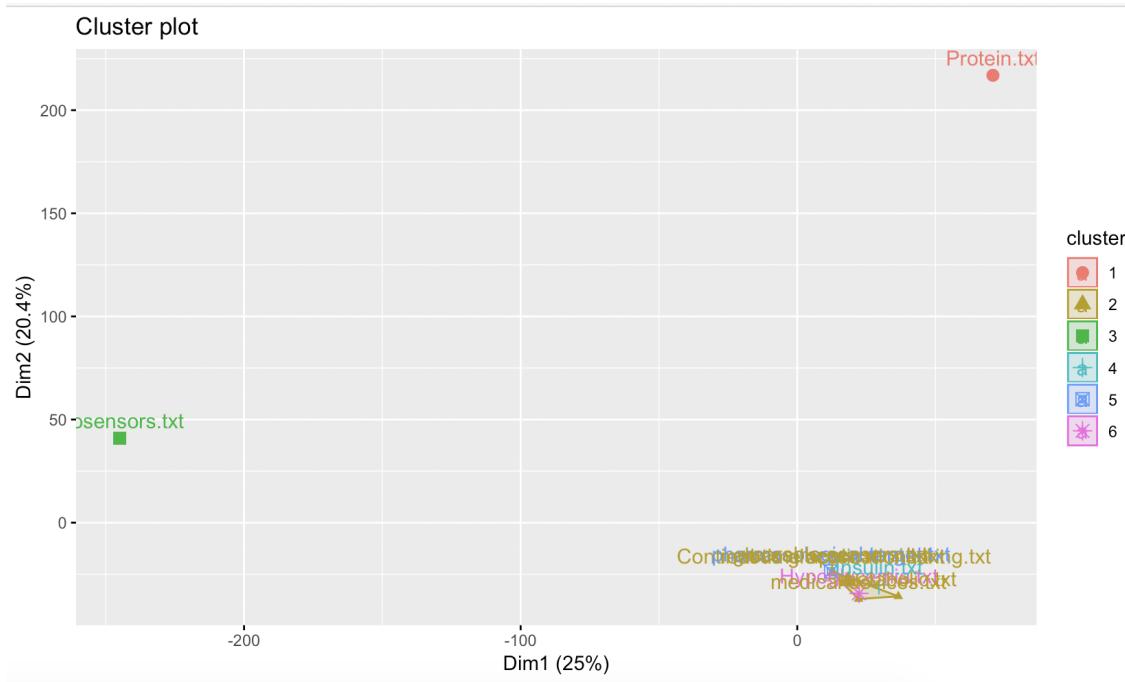
**Figure 20.** Using K-means clustering algorithm to have 3 clusters for 18 documents



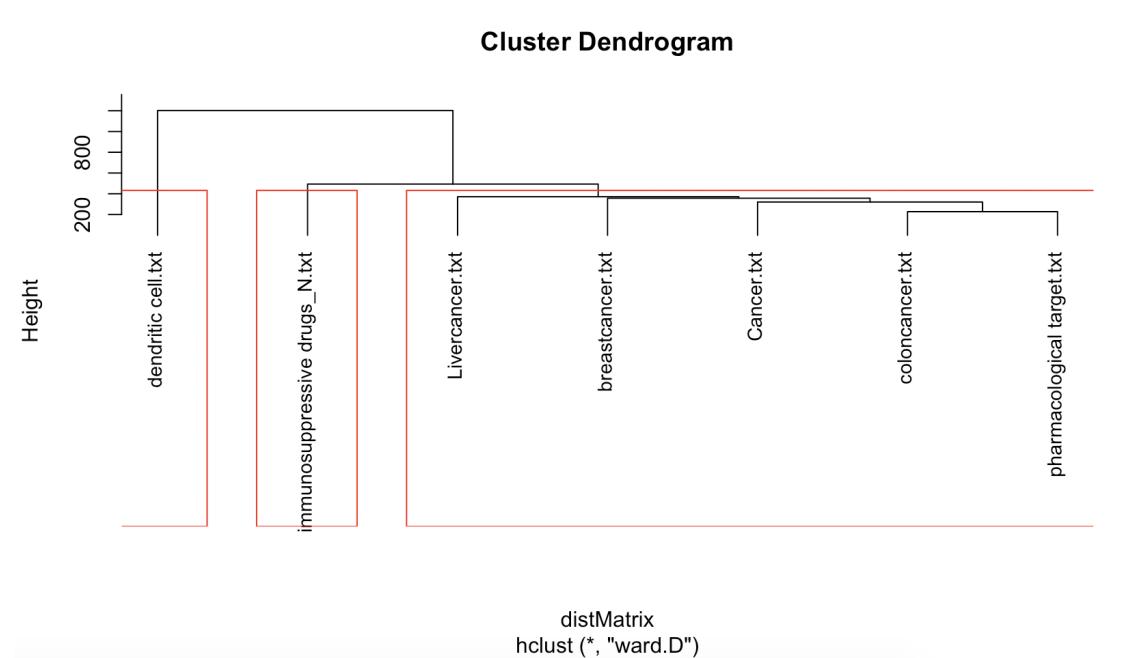
**Figure 21.** Using K-means clustering algorithm to have 4 clusters for 18 documents



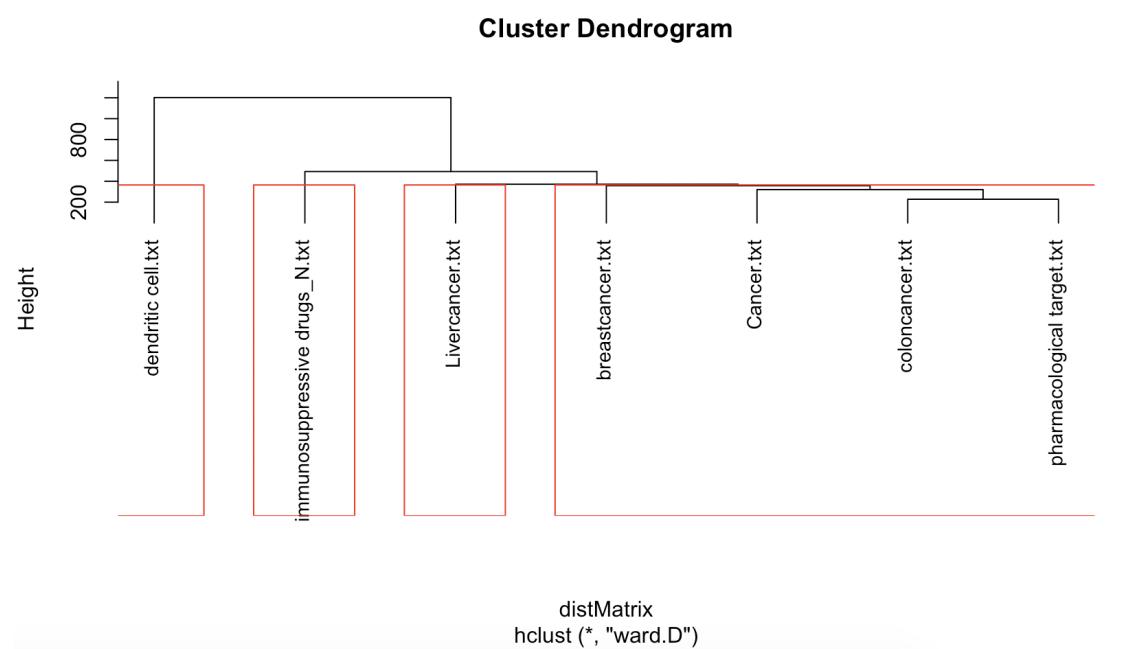
**Figure 22.** Using K-means clustering algorithm to have 5clusters for 18 documents



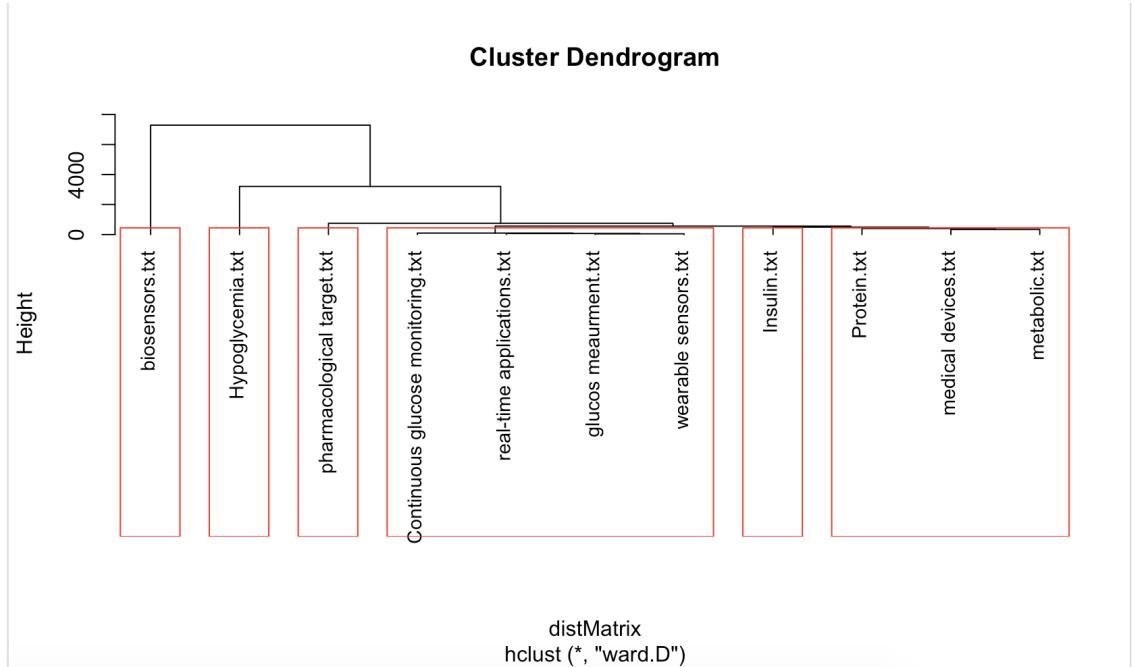
**Figure 23.** Using K-means clustering algorithm to have 6 clusters for 18 documents



**Figure 24.** Using hierarchical clustering algorithm to have 2 clusters for documents related to DCs



**Figure 25.** Using hierarchical clustering algorithm to have 3 clusters for documents related to DCs



**Figure 26.** Using hierarchical clustering algorithm to have 4 clusters for documents related to DCs

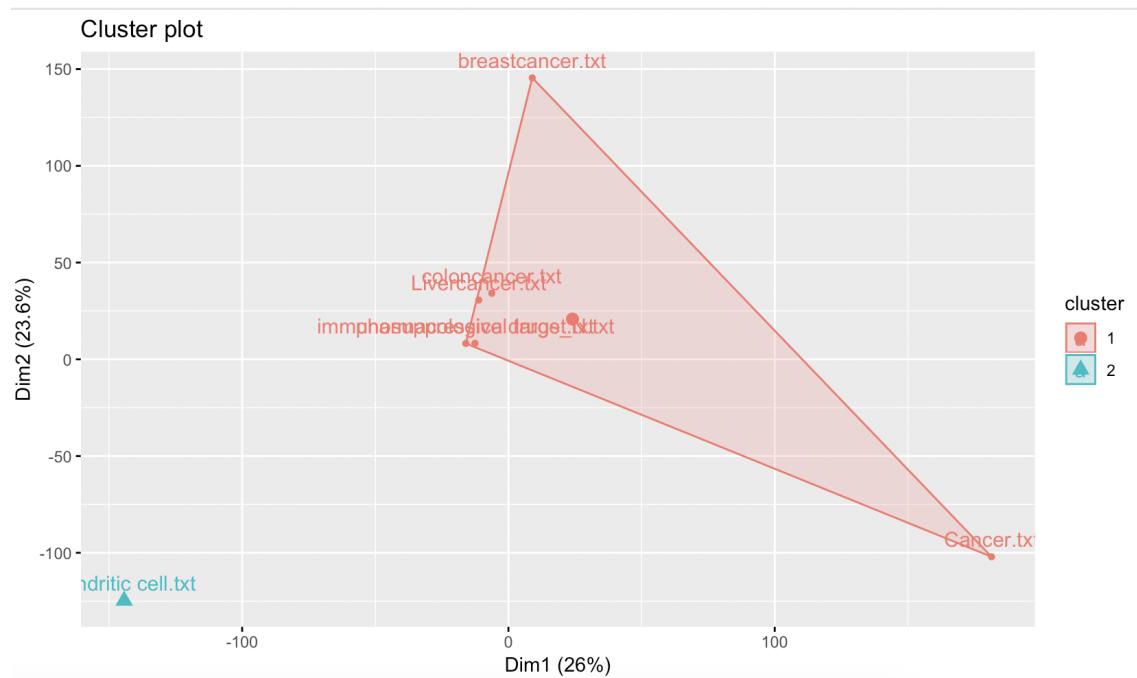
164 Figures 27, 28 and 29 show clusters obtained from applying K-means clustering algorithm on  
 165 those documents.

## 166 7. Conclusion

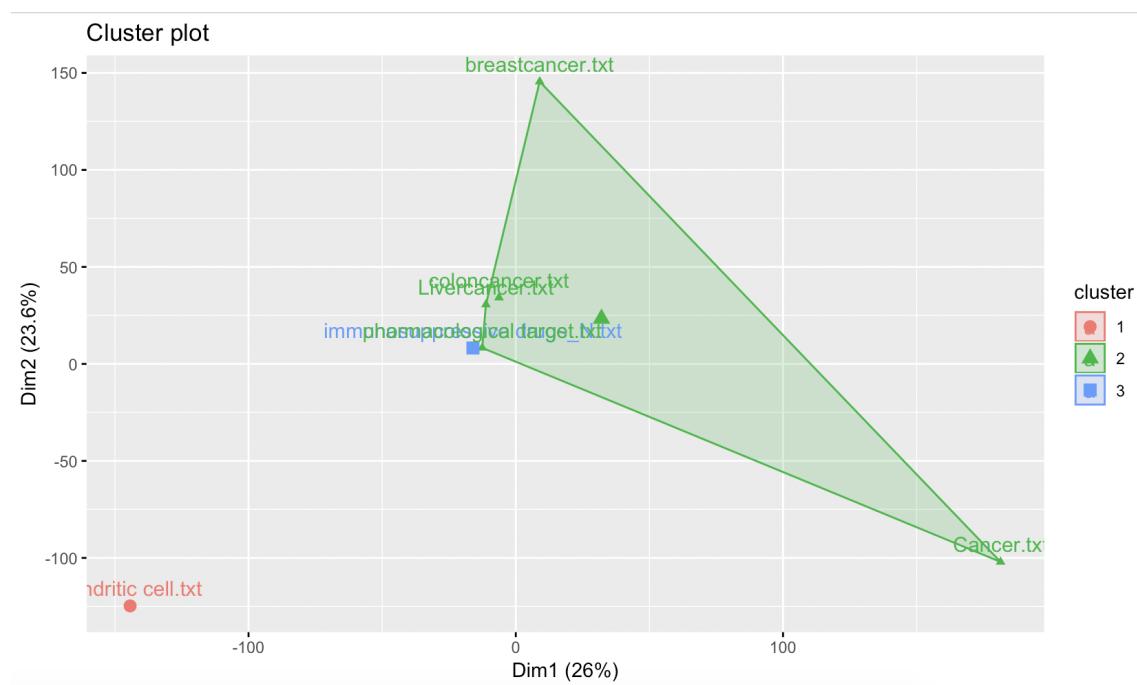
167 This paper explained how Multiple machine learning tools (MLT) and bibliometric tools can  
 168 be used to identify emerging topics and forecast trends of emerging topics. The MLT module takes  
 169 advantage of fusion algorithm to select important topics and determine relevant indices of those topics.  
 170 The BT module is responsible for constructing time-series data set of the indices. The motivation behind  
 171 presenting fusion algorithm is to develop a scalable data-driven framework with high performance.  
 172 The purpose of utilising bibliometrics is to provide multivariate time series data set. In this paper, we  
 173 showed how topic modeling and clustering algorithms are working and provided some results for our  
 174 next step. The performance and effectiveness of proposed framework would be evaluated on a variety  
 175 of text data sets (e.g., Scopus, ); the obtained results would be compared with other methodologies  
 176 that most are based on bibliometric tools. The fact that algorithm fusion of the MLT module integrate  
 177 time series prediction and identification system to elaborate topic evolution and evolving emerging  
 178 topic is one of the significant benefits of the framework.

## 179 References

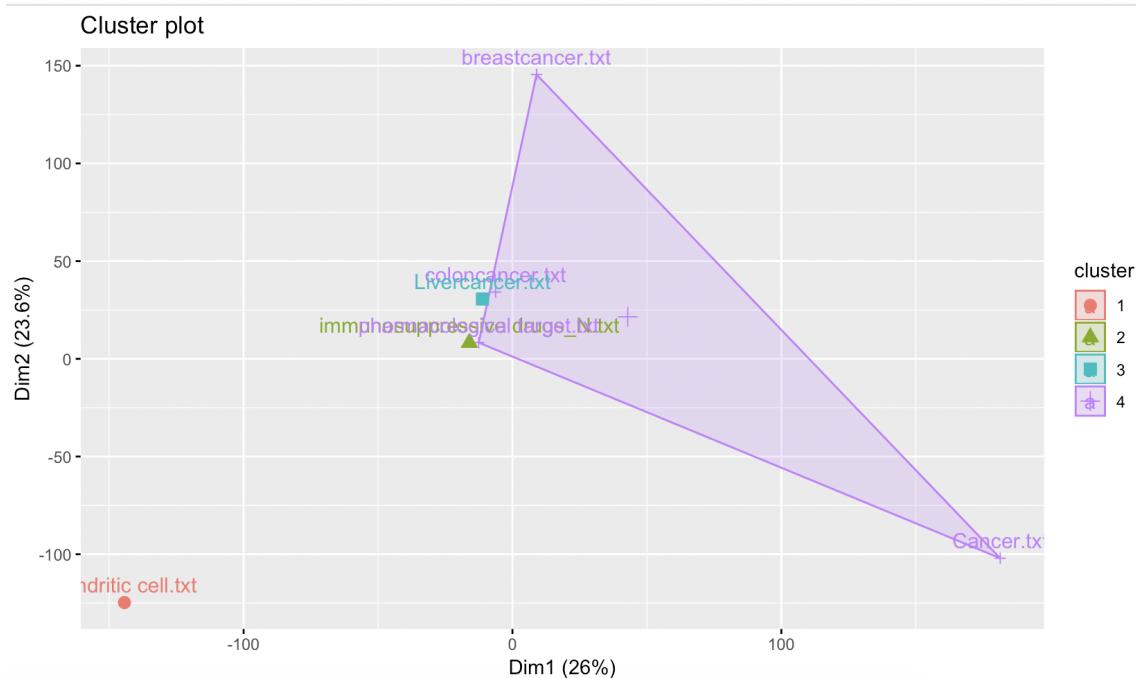
- 180 1. Mehrotra, P. Biosensors and their applications—A review. *Journal of oral biology and craniofacial research* **2016**,  
 181 6, 153–159.
- 182 2. Salatino, A. Early Detection and Forecasting of Research Trends **2015**.
- 183 3. Parsapoor, M. Brain Emotional Learning-Inspired Models **2014**.
- 184 4. Parsapoor, M. Towards Emotion-inspired Computational Intelligence (EiCI). PhD thesis, Halmstad  
 185 University, 2015.
- 186 5. Gers, F.A.; Eck, D.; Schmidhuber, J. Applying LSTM to time series predictable through time-window  
 187 approaches. In *Neural Nets WIRN Vietri-01*; Springer, 2002; pp. 193–200.
- 188 6. Xu, R.; Donald Wunsch, I. Survey of Clustering Algorithms. *IEEE TRANSACTIONS ON NEURAL  
 189 NETWORKS* **2005**, 16, 645.



**Figure 27.** Using K-means clustering algorithm to have 2 clusters for documents related to DCs



**Figure 28.** Using K-means clustering algorithm to have 3 clusters for documents related to DCs



**Figure 29.** Using K-means clustering algorithm to have 4 clusters for documents related to DCs

- 190 7. Allahyari, M.; Pouriyeh, S.; Assefi, M.; Safaei, S.; Trippe, E.D.; Gutierrez, J.B.; Kochut, K. A brief survey of  
191 text mining: Classification, clustering and extraction techniques. *arXiv preprint arXiv:1707.02919* **2017**.  
192 8. Sievert, C.; Shirley, K. LDAvis: A method for visualizing and interpreting topics. Proceedings of the  
193 workshop on interactive language learning, visualization, and interfaces, 2014, pp. 63–70.  
194 9. Blei, D.M. Probabilistic topic models. *Communications of the ACM* **2012**, 55, 77–84.  
195 10. Hurtado, J.L.; Agarwal, A.; Zhu, X. Topic discovery and future trend forecasting for texts. *Journal of Big  
196 Data* **2016**, 3, 7.

197 © 2019 by the authors. Submitted to *Journal Not Specified* for possible open access  
198 publication under the terms and conditions of the Creative Commons Attribution (CC BY) license  
199 (<http://creativecommons.org/licenses/by/4.0/>).