

# Assignment 4: Topic Modeling

Gurpreet Singh  
2018csb1092@iitrpr.ac.in

Indian Institute of Technology  
Ropar, Punjab, India

## Abstract

The main aim of this assignment is to give initial insight to topic modeling. It is a technique used in Natural Language Processing. This text mining tool is used to extract deeper information from a collection of texts. For this assignment, we have used the GENSIM library to apply topic modeling.

## 1 Dataset 1: Tasks 1, 2, 3 and 4

For the first three tasks, the dataset given is the State of Union addresses by the US presidents from 1790 to 2012. A few years contain more than one speech. Also, speeches for many years are missing from the .csv file. In total there are 155 speeches.

### 1.1 Data Preprocessing

Firstly, all the addresses are converted to lower case. After that, a list of special characters is prepared which contain new line character, full stop, comma, semi colon, question mark, exclamation mark and all types of parenthesis. All the occurrences of these characters in the dataset are replaced with spaces. Then, the text is tokenized by splitting the sentences about spaces. Those tokens which do not occur in the list 'stoplist' are kept for further analysis. 'Stoplist' is prepared by using suggestions given in the assignment itself.

Then a list which stores frequency of all the words is obtained. All those words which occur just once are then discarded. From this, we get the 'bag of words' representation using 'doc2bow' function.

### 1.2 Generation of tf-idf weighted vectors

Tf-idf stands for term frequency- inverse document frequency. It is used to give weights to words which are proportional to their frequency in a document, but inversely proportional to the number of documents which contain that word. Gensim library is used to do convert 'Bag of Words' representation to tf-idf weighted vectors. We obtain a list of 155 sub-lists which contain this vectorized form of all the words in 155 speeches.

### 1.3 LSI Topic Modeling

LSI is also known as LSA (Latent Semantic Analysis). The maximum number of topics that we can obtain is 155. Beyond that point, the number of topics generated remains 155. If we check coherence scores vs number of topics plot for LSI model, we find that the peak is in

**the range 130 to 140, slightly less than 135.** Also, the corresponding mark for 155 is also not very less. Both plots clearly depict peak in the same region.

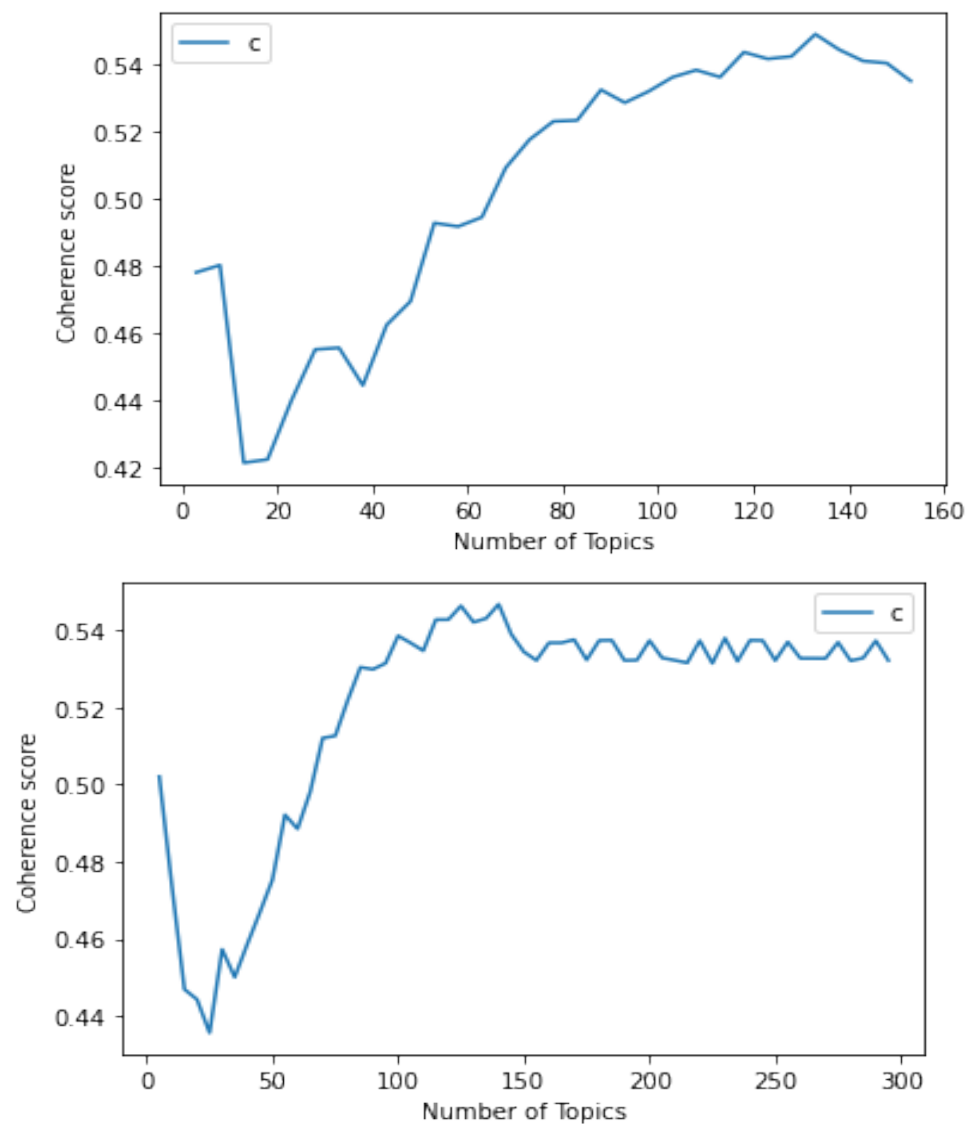


Fig-1.3(a) For LSA

From a rough glance at the topics generated we can get a feel for some events like World War 2, Saddam Hussein’s activities in Iraq, external affairs with Japan and Spain, actions against Al-Qaida and Taliban in Afghanistan, sense of togetherness and the idea of a new nation (during the first speech), etc. Let us examine any 10 out of 135 topics generated.

<pre>[112] lsi_model.show_topic(11)  [('japanese', -0.10165541557882866),  ('vietnam', 0.09343412980065383),  ('21st', -0.09015871719614697),  ('oil', 0.08271348541275785),  ('interstate', 0.08074860315658405),  ('1953', -0.07807700832064109),  ('college', -0.07748919058496871),  ('parents', -0.07678942464560579),  ('1942', -0.06856984987749444),  ('railways', 0.0610244584872311)]</pre>	<pre>[115] lsi_model.show_topic(47)  [('saddam', -0.0905326496506326),  ('1878', 0.06972488774696838),  ('award', -0.06501993505142797),  ('hussain', -0.064758694964178),  ('coinage', 0.0597686717366104),  ('kids', 0.05201148376265786),  ('japanese', 0.050663858625657294),  ('21st', -0.050512007491602035),  ('terrorists', 0.049422965522805506),  ('convenience', -0.0492512540955099)]</pre>	<pre>[118] lsi_model.show_topic(63)  [('vietnam', -0.062873429655418),  ('1964', 0.0674930631261896),  ('1993', -0.06102031679107241),  ('attitudes', 0.059771300464255),  ('-that', 0.0585248571872177),  ('steven', -0.0547319684031409),  ('communist', -0.053413218227278),  ('stress', 0.051693913015449115),  ('interdependence', 0.04964722755287733),  ('95th', 0.0486847571184441)]</pre>
<pre>[113] lsi_model.show_topic(129)  [('award', -0.0737906584905623),  ('axis', -0.06578801692459965),  ('hussain', -0.05872942763613019),  ('hats', -0.05408782344513443),  ('seventy-eight', -0.05408782344513443),  ('1878', -0.053517664557997426),  ('92d', -0.050533265283769046),  ('hitler', 0.049243149621794996),  ('1942', -0.0483853650576448),  ('1871', 0.04774975868612)]</pre>	<pre>[116] lsi_model.show_topic(92)  [('seventies', 0.06084417512635589),  ('damages', -0.05874158279721841),  ('blockades', -0.054962820826543864),  ('instruction', -0.05291873174798775),  ('grounds', -0.05062652274484741),  ('92d', -0.050212016258750845),  ('hunger', 0.045892051797787396),  ('multiplied', -0.04375020306256666),  ('sincerity', 0.0430525868157321),  ('lake', 0.04299905179862131)]</pre>	<pre>[119] lsi_model.show_topic(125)  [('award', -0.07170590384329778),  ('1827', -0.058610076328346884),  ('3rd', 0.05472961461876435),  ('1951', 0.05296158417179429),  ('productions', 0.052178802654546156),  ('00', 0.0485272013535477),  ('provinces', 0.04802264859795306),  ('1871', 0.0445502651414134),  ('articles', 0.0473134642681772),  ('strategic', 0.0440209111880496)]</pre>
<pre>[114] lsi_model.show_topic(26)  [('coinage', 0.10956647372510787),  ('democracy', 0.08998896956822168),  ('spain', 0.0813441751742563),  ('silver', 0.0710299312479884),  ('000', 0.06580638585784628),  ('1878', 0.06062497876301957),  ('coin', 0.0594878757219420),  ('navigation', -0.05744591948787535),  ('1827', 0.05508115098643661),  ('california', -0.05357947183368954)]</pre>	<pre>[117] lsi_model.show_topic(48)  [('eight-hour', 0.10293404811257947),  ('interstate', 0.09392455544048692),  ('democracy', -0.07774798308927448),  ('banks', 0.07186774387490655),  ('explanations', -0.062851519455877),  ('practically', 0.06015194014325645),  ('commission's', 0.0572251835270636),  ('1832', 0.05407616988377704),  ('spain', -0.054063980569322094),  ('1974', 0.051994736933416)]</pre>	<pre>[120] lsi_model.show_topic(118)  [('afghanistan', 0.07261966270119049),  ('al-qaida', -0.0546669503417413),  ('retreat', -0.05359808432137954),  ('exit', -0.05197322234029417),  ('crude', -0.05183744217945732),  ('1878', -0.05013565681266248),  ('challenge', -0.0487208878118028),  ('hostages', 0.04938698113513995),  ('17', -0.0484000072852628),  ('36', -0.0484000072852628)]</pre>
		<pre>[121] lsi_model.show_topic(18)  [('notes', 0.1513975665810703),  ('paper', 0.11874927419080015),  ('exchequer', 0.10730845797264134),  ('specie', 0.0912625034719027),  ('circulation', 0.0873569288387524),  ('insurrection', -0.0843005724134875),  ('silver', 0.0780738004455462),  ('minister', 0.07343338546106291),  ('currency', 0.07021053091083132),  ('rebellion', -0.0603479453130906)]</pre>

Fig-1.3(b), 1.3(c) and 1.3(d)

In 1.3(b), part 1 discusses about Vietnam and Japan. Years mentioned are 1942 and 1953. Its theme is possibly regarding US-Vietnam war. Part 2 mentions Hussein and Hitler which is a clear indication of its connection with Iraq and its dictator Saddam Hussein. Hitler is probably used here to give analogy. Part 3 discusses about democracy. 1827 and 1878 are mentioned. However, choosing one topic theme is difficult here.

In 1.3(c) part 1, Saddam Hussein is mentioned here again. He is the theme of this topic too. In part 2, hunger and damages are mentioned. The probable theme of this topic is regarding poverty. In part 3, the discussion is probably about economy and external affairs. But the conclusions here are not as concrete.

In 1.3(d), part 1 mentions Vietnam and communism. So, this seems to be a talk about the region of China and South-East Asia. It is very difficult to make sense out of Part 2. But still there is a mention of provinces and articles. This can be point towards the federal structure, power sharing, local and inter state affairs. Part 3 mentions Al-qaida, Afghanistan and hostages. This is definitely about some terrorist activity carried out by Al-Qaida. Part 4 seems very interesting. It mentions notes, paper, silver, currency, circulation and rebellion. Economy is the protagonist here for sure. But the word 'rebellion' adds a little twist. It might be indicative of the Civil War.

To sum up, analysis using LSI was very informative and useful.

## 1.4 LDA Topic Modeling

LDA stands for Latent Dirichlet Allocation. Unlike LSA, there is not any significant upper limit restriction on the number of topics here (The maximum I tried was around 1000). We again use the same method as LSA to check optimal number of topics in LDA.

From the coherence score plot, we get that the maximum score is occurring at number of topics equal to 35. **So, the optimal number is 35.**

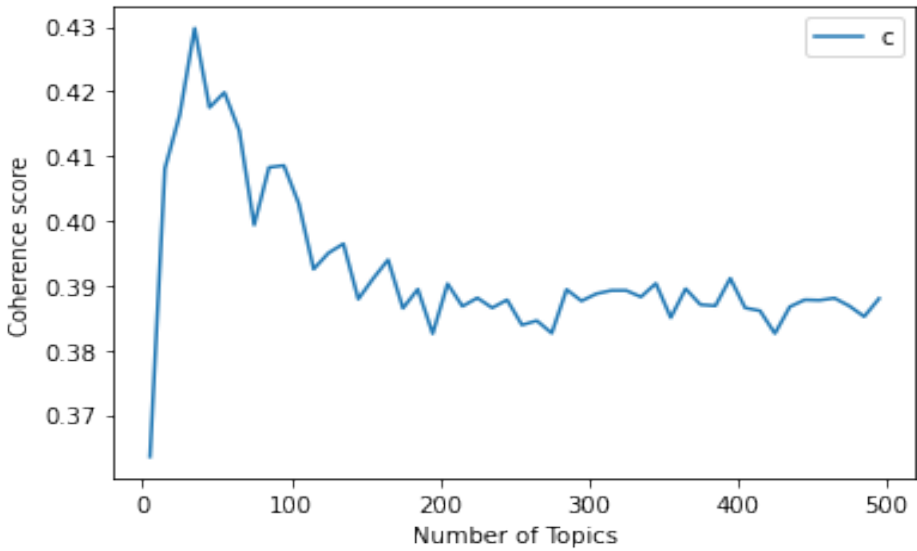


Fig-1.4(a)

All the 35 topics have been printed in the topic. Let us select any 10 out of them. It is very difficult to get something significant out of these 10 topics. All the topics generated seem to be a mix of 'government', 'united', 'people', 'state' and 'congress'. The performance of LDA is very poor as compared to LSI, here. The ten topics which we sampled are nearly impossible to differentiate from each other. The performance of LDA is very much influenced by modals- 'will', 'can' and 'may', pronouns and prepositions. The stoplist was same for LDA and LSI. Even then, LSI produced far better results.

```
[ ] lda_model.show_topic(15)
[('will', 0.016510025),
 ('must', 0.006928436),
 ('can', 0.0061527365),
 ('united', 0.0060042962),
 ('government', 0.0059873736),
 ('people', 0.005278812),
 ('congress', 0.0051245587),
 ('us', 0.0050686435),
 ('year', 0.005042023),
 ('one', 0.0043523307)]

[ ] lda_model.show_topic(19)
[('will', 0.018163654),
 ('government', 0.0077068997),
 ('must', 0.0071482946),
 ('states', 0.0051086713),
 ('us', 0.0046383967),
 ('world', 0.0046247444),
 ('can', 0.004391574),
 ('congress', 0.004288021),
 ('year', 0.0041529457),
 ('people', 0.004143453)]

[ ] lda_model.show_topic(17)
[('will', 0.01479972),
 ('must', 0.007606434),
 ('government', 0.0072158566),
 ('us', 0.005724601),
 ('people', 0.005284076),
 ('congress', 0.0052839634),
 ('united', 0.0051261256),
 ('can', 0.0050686053),
 ('world', 0.004887835),
 ('states', 0.0043795537)]

[ ] lda_model.show_topic(8)
[('will', 0.011218465),
 ('government', 0.008520857),
 ('united', 0.006148245),
 ('states', 0.0058669923),
 ('people', 0.0052705654),
 ('can', 0.005111385),
 ('may', 0.0046788733),
 ('congress', 0.0046211877),
 ('country', 0.004057839),
 ('must', 0.003949522)]

[ ] lda_model.show_topic(7)
[('will', 0.012592978),
 ('government', 0.0071704597),
 ('states', 0.0068256347),
 ('congress', 0.006195552),
 ('can', 0.005533145),
 ('people', 0.005294831),
 ('new', 0.0048970347),
 ('world', 0.004655728),
 ('must', 0.0045595467),
 ('america', 0.004422013)]

[ ] lda_model.show_topic(32)
[('will', 0.013705566),
 ('government', 0.006362667),
 ('congress', 0.0063563953),
 ('united', 0.006173442),
 ('can', 0.006073642),
 ('states', 0.005350847),
 ('must', 0.004613758),
 ('people', 0.0042024194),
 ('country', 0.0039204345),
 ('may', 0.003844109)]
```

Fig-1.4(b),(c) For LDA

## 1.5 A General Analysis: Task 4

For visualising data obtained from LDA, we have used the library 'pyLDavis'. It produces marvellous visualisations and graphics. They give a very good insight into the significance of all the topics. The html page which is generated using this library has been embedded in the folder submitted. On that page, one can vary relevance metric for any topic and can see top 30 words that are contained in that topic. Also, its image has been attached in this report.

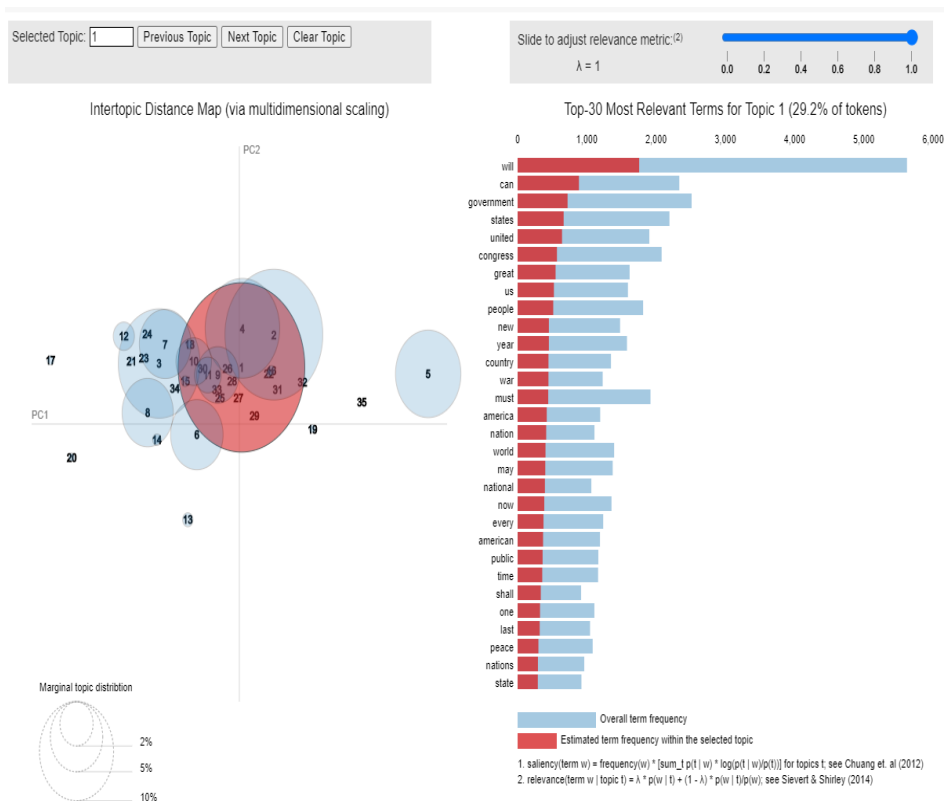


Fig-1.5(a) Number of topics=35

## 2 Dataset 2: Task 5

Dataset chosen for this task is the collection of AP wire stories. Data preprocessing is done in same way as dataset 1. After that, tf-idf weighted vectors are generated. Then, the optimal number of topics is chosen for LDA by following the same procedure as task 3. From the plot, it is clear that the peak is observed at number of topics equal to 15.

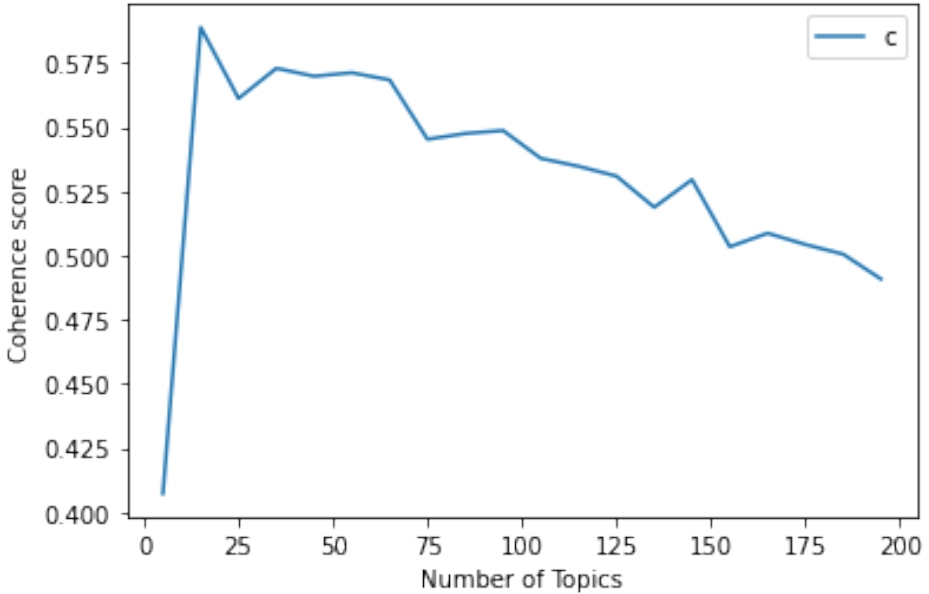


Fig-2(a)

```
[0,
'0.015**said" + 0.008**\'\' + 0.005**stock" + 0.004**corroon" + 0.004**cdy" + 0.004**percent" + 0.004**market" + 0.004**black" + 0.004**clr" + 0.004**2**'),
(1,
'0.015**said" + 0.011**\'\' + 0.005**poll" + 0.005**dukakis" + 0.005**new" + 0.005**percent" + 0.004**will" + 0.003**party" + 0.003**also" + 0.003**last**'),
(2,
'0.024**said" + 0.012**\'\' + 0.007**year" + 0.006**percent" + 0.006**s" + 0.006**u" + 0.004**will" + 0.004**last" + 0.004**owen" + 0.003**united**'),
(3,
'0.027**said" + 0.011**percent" + 0.010**\'\' + 0.005**million" + 0.004**1" + 0.004**u" + 0.003**will" + 0.003**soviet" + 0.003**two" + 0.003**president**'),
(4,
'0.021**said" + 0.010**\'\' + 0.005**president" + 0.005**bush" + 0.005**dresses" + 0.005**will" + 0.004**at&t" + 0.004**s" + 0.004**macmillan" + 0.003**u**'),
(5,
'0.019**said" + 0.011**\'\' + 0.006**will" + 0.006**0000" + 0.005**u" + 0.005**s" + 0.004**pacs" + 0.004**new" + 0.004**court" + 0.003**united**'),
(6,
'0.014**said" + 0.010**\'\' + 0.009**1" + 0.007**late" + 0.006**dollar" + 0.005**yen" + 0.005**000" + 0.005**million" + 0.004**hyundai" + 0.004**police**'),
(7,
'0.023**\'\' + 0.021**said" + 0.006**south" + 0.004**people" + 0.004**new" + 0.004**primary" + 0.004**state" + 0.003**i" + 0.003**president" + 0.003**bush**'),
(8,
'0.023**said" + 0.020**\'\' + 0.006**u" + 0.006**government" + 0.005**trade" + 0.005**will" + 0.005**s" + 0.005**states" + 0.005**united" + 0.003**new**'),
(9,
'0.010**said" + 0.009**\'\' + 0.005**will" + 0.004**000" + 0.004**president" + 0.003**people" + 0.003**new" + 0.003**million" + 0.003**also" + 0.003**two**'),
(10,
'0.015**said" + 0.012**\'\' + 0.005**new" + 0.005**million" + 0.004**government" + 0.004**1" + 0.003**u" + 0.003**today" + 0.003**market" + 0.003**stock**'),
(11,
'0.020**said" + 0.015**\'\' + 0.012**percent" + 0.007**billion" + 0.005**year" + 0.004**million" + 0.004**000" + 0.004**last" + 0.003**1" + 0.003**will**'),
(12,
'0.020**said" + 0.020**\'\' + 0.007**will" + 0.005**one" + 0.004**new" + 0.003** " + 0.003**s" + 0.003**000" + 0.003**years" + 0.003**can**),
(13,
'0.031**said" + 0.013**\'\' + 0.004**one" + 0.003**new" + 0.003**will" + 0.003**also" + 0.003**1" + 0.003**two" + 0.003**tuesday" + 0.002**percent**'),
(14,
'0.027**said" + 0.015**\'\' + 0.004**police" + 0.004**-" + 0.004**people" + 0.004**two" + 0.003**000" + 0.003**united" + 0.003**government" + 0.003**will**')]
```

Fig-2(b) All 15 topics

The results of LDA are a lot better here. A wide variety of words is captured here. The effect of modals is very diminished as compared to the previous dataset. The visualisations obtained can be viewed much better in the web page embedded in the folder submitted. Nevertheless, an image of the page is attached here. The term 'better' means that a variety of words are captured with proper significance on the necessary words.

Selected Topic:  Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric: <sup>(2)</sup>   $\lambda = 1$

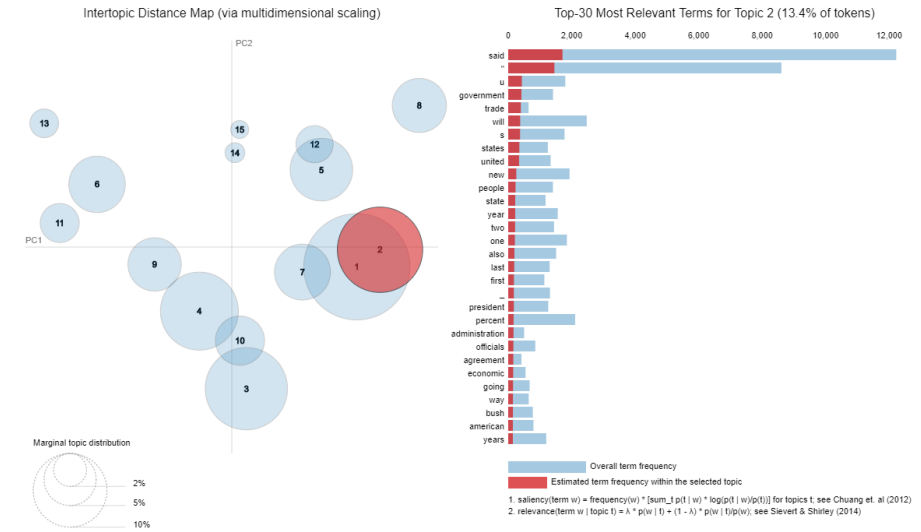


Fig-2(c)