National
College *of*
Ireland

# Configuration Manual

# Application of short text topic modelling techniques to Greta Thunberg discussion on Twitter

MSc Research Project
Masters in Data Analytics

## Sean Dingemans
Student ID: x18199089

School of Computing
National College of Ireland

Supervisor:     Dr Catherine Mulwa

**National College of Ireland**

**MSc Project Submission Sheet**

**School of Computing**

| | |
|---|---|
| **Student Name:** | Sean Dingemans…….…………………………………………………………………… |
| **Student ID:** | x18199089……………………………………………………………………… |
| **Programme:** | Masters in Data Analytics ……………………… **Year:** 2020…………………….. |
| **Module:** | MSc Research Project……………………………………………………………..……… |
| **Lecturer:** | Dr Catherine Mulwa……………………………………………………………………… |
| **Submission Due Date:** | 17 August 2020……………………………………………………………………..……… |
| **Project Title:** | Application of short text topic modelling techniques to Greta Thunberg discussion on Twitter………**.**……………………….………… |
| **Word Count:** | …2297…………………… **Page Count:** ………..…57… |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project.  All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section.  Students are required to use the Referencing Standard specified in the report template.  To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** ……………………

**Date:** 17 August 2020……………………………………………………………………………

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | ☐ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid.  It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| Office Use Only | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Sean Dingemans
Student ID: x18199089

## 1 Environment Setup

An AWS account was setup with billing information then filled in. Afterwards a virtual machine instance was setup through Amazon's Elastic Cloud Compute (EC2) device. The instance creation process begins with the selection of 'Launch Instance' as shown in Figure 1.
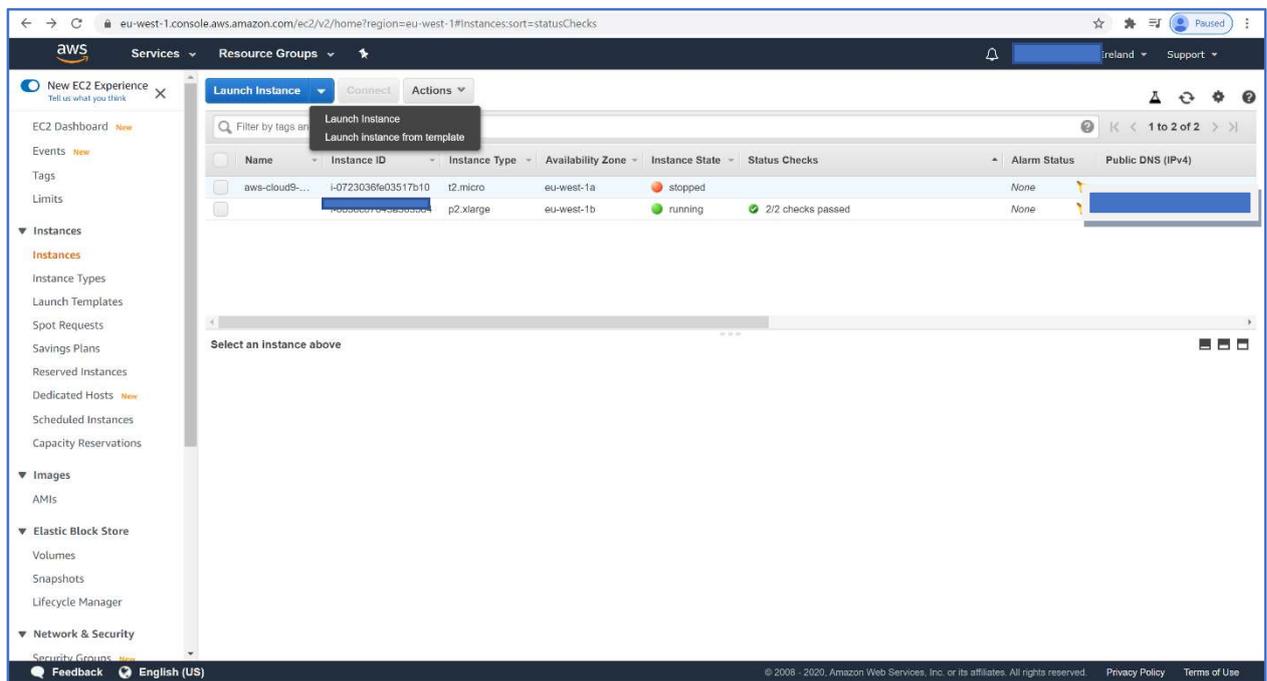


**Figure 1:** AWS creation of instance by launce instance

The first step of instance creation was to select an Amazon Machine Image (AMI). An AMI was chosen for the ec2 instance that had much of the required preinstalled software, the Deep-Learning Ubuntu AMI: pandas, numpy, python, jupyter notebook, java 1.8, Anaconda

**Figure 2:** AWS section of AWS AMI

An EC2 instance of the gpu family was chosen as these instances are optimised for computationally heavy tasks in addition to their hardware specs. The p2x instance was selected due to being affordable relative to other instances in the GPU family. Specs for the instance can be seen in Figure 3.



**Figure 3:** AWS GPU Instance Hardware Type selection

**Figure 4:** AWS Brief overview of some of the instance hardware specs.

A standard VPC was created. This was this first step in allowing the instance to have an IP address so that it could be accessed. The IPV4 CIDR block size provides 64000 different ip addresses. (see figure 4).



**Figure 5:** AWS VPC creation part 1

**Figure 6:** AWS CIDR Block specification for VPC part 2

For the VPC to work, a subnet must be assigned. A subnet utilising 16 of the 64000 possible Ip addresses was created. Only one IP address will be needed, however no additional costs are accrued. The region of the instance is also set here. It is set to eu-west-1b. The region specified is not critical, as there are not thousands of clients accessing the instance, and thus any lag is negligible.



**Figure 7:** AWS Subnet creation for VPC

**Figure 8:** AWS Subnet Creation

This step was just done to initialise the instance with an ip address.



**Figure 9:** AWS IP address allocation for instance creation

**Figure 10:** AWS Auto-assign IP address

The Created VPC along with its subnet could now be allocated to the instance to be created.



**Figure 11:** AWS Allocation of sub-net and vpc to instance

In step 3 of launching the instance, a list of installation commands were added for the environment. Installed software included nginx for the ability to run programs in the browser through the instance's associated Ip address, once the instance was running supervisor for starting up nginx automatically when starting the instance after it is stopped naano text editor.



**Figure 12:** AWS Configuration Details



**Figure 13:** AWS Configure Instance details 2

In stage 4, storage was allocated to the instance. An additional 90GB was allocated to the instance later on. All other settings were left default.



**Figure 14:** AWS Add storage

A security group had to be created to specify inbound rules to access Instance. Since Jupyter notebooks were to be accessed from the browser, https and http protocols were specified (the ubuntu shell by itself can't be used for interaction with jupyter notebooks, a browser is needed for this interactivity to be allowed). SSH was specified as well so that the ubuntu shell of the instance could be accessed – many programmes were to be run from the shell.



**Figure 15:** AWS security group creation for protocol specifications

Rules I used for my own project, a port is allocated for ssh purposes, port used to ssh into the instance



**Figure 16:** AWS protocols for inbound rules regarding access of Instance

Upon reviewing and launching the instance, a private key was allocated to the instance and automatically downloaded. The instance was not accessed yet



**Figure 17:** Creation of private hash key used to ssh into instance

Once stopped, an IP address from the subnet was then set to be an elastic IP address and was then allocated to the instance. Setting the ip address to static allowed for the instance to be accessed with the same ip address after the instance is stopped and started again. Nginx and supervisor programs (used for displaying the Jupyter notebook in the browser) would need to be reconfigured if IP address was not static.



**Figure 18:** Elastic IP address association

In order to conveniently access the EC2 instance via SSH, it was decided to use the client SSH software tool Bitvise as the tool also allowed for file transfer to the local computer. The host address and port was specified and the client key downloaded from Amazon was added for authentication purposes.



**Figure 19:** Using Bitvise to SSH into instance

**Figure 20:** Import PEM key into Bitvise



**Figure 21:** Key pair that was imported int Bitvise

**Figure 22:** Jupyter notebook configuration initialisation

A hashed password that was generated is put into the config file the notebook.  The network of the notebook is set to the ip address of the instance to prevent the unlikely sharing of resources with other external sites



**Figure 23:** Generation of hashed password for notebook

**Figure 24:** Hashed key notebook password and localhost Port to access Notebook

Nginx is now configured to display Jupyter notebook in the browser via the public ip address of the instance. The EC2 instance does not have visual rendering capabilities and thus    First the user must change directories to the sites-available under nginx directory.



**Figure 25:** Nginx directory

**Figure 26:** Configuration file created for nginx

Specify the listening ports of the ec2 instance to allow nginx to display on the localhost port 8888 what is being received on ec2 port 80. The Jupyter notebook (which was configured to render visually on port 8888) can be rendered visually through the proxy server nginx on localhost port 8888



**Figure 27:** Port specification for nginx Jupyter notebook rendering

Once the sites-available file has been saved, the file can be copies to the sites-enabled directory and a symbolic link can be used to create a link between jupyter_app.conf in both sites-available and sites-enabled directories.



**Figure 28:** Symbolic link creation between sites available and sites-enabled

Command used to reload and update nginx configuration settings to recognise the new site configuration.

**Figure 29:** Updating of nginx sites available / sites enabled

Jupyter notebook can now be run on local host. However, we still need to ssh into the instance's ubuntu terminal and start jupyter notebook before accessing the notebook via localhost with the browser. It is more convenient to configure the notebook to run in the background of the instance. For this, supervisor will be needed.



**Figure 30:** Configuration moving to supervisor directory

Use sudo nano to create my_jupyter.conf



**Figure 31:** Use sudo nano to create my_jupyter.conf

Specify supervisor to run the jupyter notebook command upon starting up



**Figure 32:** Configure Jupiter notebook to run when instance is started

**Figure 33:** Create log file directory

Create log folder, for output as specified in my_jupyter.conf



**Figure 34:** Refresh systemctl daemon that uses supervisor daemon
process to re-initialise supervisor

**Figure 35:** Re-read configuration file for supervisor daemon



**Figure 36:** Refresh supervisor daemon

**Figure 37:** Running supervisor daemon

The notebook will now run upon starting the instance



**Figure 38:** Selection of Elastic IP addresses

**Figure 39:** List of generated Elastic IP addresses



**Figure 40:** IP address association to current instance

**Figure 41:** AWS Launch Instance



**Figure 42:** Private IP address

# 2 Data Collection and Processing

twint -s 'greta' --since 2019-09-20 --until 2019-09-30 -l en -d true -o Sep2030Greta3.json



**Figure 43:** Tmux for creating processes that can continue running uninterrupted and be returned to later on

Twint was collected with the following command



**Figure 44:** TMUC process for collection of Twitter data with TWINT

Time plot of collected data

```
In [1]:   ▶    1  import twint
               2  import nest_asyncio
               3  import pandas as pd
               4  nest_asyncio.apply()
               5  import seaborn as sns
               6  # package to clean text
               7  import re
               8  from nltk import FreqDist
               9  pd.set_option('display.max_colwidth', -1)
```

```
In [2]:   ▶    1  c = twint.Config()
               2  c.Search = 'greta'
               3  c.Lang = 'en'
               4  c.Since = '2019-09-20'
               5  c.Until = '2019-10-01'
               6
               7  c.Hide_output = True
               8  c.Store_json = True
               9  c.Output = 'Greta27MaySepJup.json'
              10  c.Resume = 'GT_last2.csv'
              11  c.Debug = True
```

```
In [ ]:   ▶    1  twint.run.Search(c)
```

```
In [3]:   ▶    1  df = pd.read_json('Greta27MaySepJup.json' , lines = True)
               2  df.info()
```

```
. . .
```

```
In [29]:  ▶    1  %matplotlib inline
               2  df['date'].value_counts().plot()
```

Out[29]: <matplotlib.axes._subplots.AxesSubplot at 0x7fbc6b80b4a8>



**Figure 45:** Initial exploration for the analysis of traffic

25

```
In [63]:    1  #SelectedCols = ['cashtags','hashtags','Link','mentions','quote_url','source','translate']
            2  #SelectedCols = ['hashtags','tweet']

In [4]:     1  pd.set_option('display.max_colwidth', -1)

In [64]:    1  df[SelectedCols].sample(100)

                                                    ...

In [6]:     1  # create newdataframe to count total number of hashtags
            2  flattened_hashtags_df = pd.DataFrame(
            3      [hashtag for hashtags_list in df.hashtags
            4       for hashtag in hashtags_list],
            5      columns=['hashtag'])

In [8]:     1  #totla hashtag count
            2  flattened_hashtags_df['hashtag']

   Out[8]:  227455

In [10]:    1  allWords = [word for item in list(flattened_hashtags_df['hashtag']) for word in item]

In [16]:    1  fdist = FreqDist(flattened_hashtags_df['hashtag'])

In [17]:    1  len(fdist)

   Out[17]: 36507

In [18]:    1  print(fdist)

            <FreqDist with 36507 samples and 227455 outcomes>

In [19]:    1  #Looking at Hashtags Still
            2  k = 6000
            3  top_k_words = fdist.most_common(k)
            4  char_frequency = top_k_words

In [20]:    1  print(char_frequency)
```

**Figure 46:** Examination of hashtag counts

```
In [76]:  ▶  1  popular_hashtags
```

Out[76]:

| | hashtag | counts |
|---|---|---|
| 0 | #climatestrike | 15964 |
| 1 | #climatechange | 11604 |
| 2 | #greta | 7814 |
| 3 | #gretathunberg | 7282 |
| 4 | #climateaction | 6263 |
| 5 | #fridaysforfuture | 6227 |
| 6 | #howdareyou | 6023 |
| 7 | #climatecrisis | 3764 |
| 8 | #climateactionsummit | 3200 |
| 9 | #climateemergency | 2753 |
| 10 | #gretathurnberg | 2671 |
| 11 | #unga | 2105 |
| 12 | #globalclimatestrike | 1929 |
| 13 | #climate | 1891 |
| 14 | #gretathunbergoutdidtrump | 1408 |
| 15 | #auspol | 1390 |
| 16 | #naturenow | 1231 |
| 17 | #climateactionnow | 1173 |
| 18 | #climatehoax | 1040 |
| 19 | #climatestrikecanada | 1037 |
| 20 | #trump | 987 |
| 21 | #climatechangeisreal | 979 |
| 22 | #schoolstrike4climate | 939 |
| 23 | #globalwarming | 924 |
| 24 | #gretathumberg | 879 |
| 25 | #cdnpoli | 834 |
| 26 | #greta4nobelprize | 775 |
| 27 | #climatestrikes | 774 |
| 28 | #maga | 748 |
| 29 | #climatechangehoax | 712 |
| ... | ... | ... |

**Figure 47:** Most popular hashtags discovered

**Figure 48:** No retweets were found for the 740000 collected tweets

```
3   # package to clean text
4   import re
5   from nltk import FreqDist
6   import demoji
7   demoji.download_codes()
8   import nltk
9   from nltk.corpus import wordnet
10  import numpy as np
11  import matplotlib.pyplot as plt
12  import numpy as np
13  import pandas as pd
14  pd.set_option('display.max_colwidth', -1)

Downloading emoji data ...
... OK (Got response in 0.28 seconds)
Writing emoji data to /home/ubuntu/.demoji/codes.json ...
... OK
```

```
1   dfCheck6Jun = pd.read_json('Greta27MaySepJup.json' , lines = True)
```

```
1   # Processing takes too long sometimes - I need a test dataset
2   dummyCheck = dfCheck6Jun.sample(100)
```

```
1   def remove_links(tweet):
2       '''Takes a string and removes web links from it'''
3       tweet = re.sub(r'http\S+', '', tweet) # remove http links
4       tweet = re.sub(r'bit.ly/\S+', '', tweet) # rempve bitly links
5       tweet = re.sub(r'pic.twitter.com\/[^\s]{2,}','', tweet) # remove pictures
6       tweet = re.sub(r'#[^\s]{2,}','', tweet) # remove hashtags
7       tweet = tweet.strip('[link]') # remove [links]
8       return tweet
9
10  def remove_users(tweet):
11      '''Takes a string and removes retweet and @user information'''
12      tweet = re.sub('(RT\s@[A-Za-z]+[A-Za-z0-9-_]+)', '', tweet) # remove retweet
13      tweet = re.sub('(@[A-Za-z]+[A-Za-z0-9-_]+)', '', tweet) # remove tweeted at
14      return tweet
15
16
```

```
1   nltk.download('stopwords')
2   nltk.download('averaged_perceptron_tagger')
3   nltk.download('wordnet')

[nltk_data] Downloading package stopwords to /home/ubuntu/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]     /home/ubuntu/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]       date!
```

**Figure 49:** Functions used to clean tweets of metadata

28

```python
my_stopwords = nltk.corpus.stopwords.words('english')
word_rooter = nltk.stem.snowball.PorterStemmer(ignore_stopwords=False).stem
my_punctuation = '!"$%&\'()*+,-./:;<=>?[\\]^_`{|}~•@'

# cleaning master function
def tokenCheck(tweet, bigrams=False):
    tweet = remove_users(tweet)
    tweet = remove_links(tweet)
    tweet = demoji.replace(tweet)
    tweet = tweet.lower() # Lower case
    tweet = re.sub('\s+', ' ', tweet) #remove double spacing

    tweet = re.sub('['+my_punctuation + ']+', ' ', tweet) # strip punctuation
    tweet = re.sub('\s+', ' ', tweet) #remove double spacing
    tweet = re.sub('([0-9]+)', '', tweet) # remove numbers

    tweet = nltk.word_tokenize(tweet)
    return [word for word in tweet if word not in my_stopwords and re.match(r'(?<!\S)[a-z-]+(?!\S)', word)]
```

```python
my_stopwords = nltk.corpus.stopwords.words('english')
```

```python
my_stopwords.extend(['greta','thunberg','via', 'cphqnukxenzvwvqmwabitddcoet', 'xfdfadcbaeafceedac', 'etc', 'th','ur','in
'sn','cf','hd','hp','nk','oy','sg','ic','lw','ut','wt','ak','gl','od','po','sb','ys','zu','cv','dx','fe','fw','mj','ml',
'lv','oo','pi','xs','ze','aj','ec','eq','lb','mb','mx','rc','tt','vw','bj','cs','ly','oc','oi','qu','wp','ws','fc','hc',
'nm','ri','tc','hy','jb','kc','nr','ox','pj','rr','rx','sk','tj','ao','cp','db','dp','fy','gs','jw','nv','vc','bk','cw',
'ci','cz','dg','dv','gi','gu','iz','rm','rp','wm','wr','ax','bm','hf','jd','jj','lm','pf','rf','rl','tg','vo','wc','za',
'vt','xl','xy','ae','bg','bl','cb','dh','ee','ei','ez','fd','fg','fx','gn','hb','io','kn','lp','ls','mz','pb','pn','rb',
])
```

**Figure 50:** Central tweet processing method.

All numbers, punctuation and non-Latin characters are removed. Tweets are all converted to lower case



```python
# function to remove single characters
def Filter_Single_Char_words(text):
    return [word for word in text if len(word) > 1]
```

```python
# get rid of single characters
dfJSONN9Jun['TokenCheck2'] = dfJSONN9Jun['TokenCheck2'].apply(Filter_Single_Char_words)
```

```python
#extended stopwords
def Extended_Stop_Words(text):
    return [word for word in text if word not in my_stopwords]
```

```python
print(my_stopwords)
```

```python
# get rid of single characters #This method took 2hrs!!!
dfJSONN9Jun['TokenCheck2'] = dfJSONN9Jun['TokenCheck2'].apply(Extended_Stop_Words)
```

**Figure 51:** After initial preprocessing and examination, smaller more specific functions were run.

```
In [14]:  ▶  1  # God, the run above went on from at least 1430 - 1600 - better write it to a file (check if it overwrites)
             2  #dfJSONN9Jun.to_json('8JuneNoHashtagsAllLatin.json', orient='records', lines = True)
             3  dfJSONN9Jun.to_json('12JuneTop29870Words.json', orient='records', lines = True)

In [ ]:   ▶  1  dfCheck6Jun['TokenCheck2'] = dfCheck6Jun.tweet.apply(tokenCheck)

In [ ]:   ▶  1  dfCheck6Jun[['tweet','TokenCheck2']].sample(100)

In [ ]:   ▶  1  dfCheck6Jun.to_json('8JuneNoHashtagsAllLatin.json', orient='records', lines = True)

In [10]:  ▶  1  dfJSONN9Jun = pd.read_json('12JunNoSinglesLessStopWds.json', lines = True)

In [15]:  ▶  1  allWords = [word for item in list(dfJSONN9Jun['TokenCheck2']) for word in item]

In [16]:  ▶  1  fdist = FreqDist(allWords)

In [17]:  ▶  1  len(fdist)

Out[17]:  110517

In [ ]:   ▶  1  print(fdist)

In [19]:  ▶  1  k = 29786
             2  top_k_words = fdist.most_common(k)
             3  MaybeNoLatin = top_k_words
```

**Figure 52:** All unique tokens and their frequencies were examined.

29786 tokens had frequencies greater than 5

```
In [34]:  ▶  1  print(MaybeNoLatin)

[('greta', 365221), ('thunberg', 153890), ('climate', 114289), ('people', 66511), ('world', 63288), ('change', 62544), ('lik
e', 55362), ('one', 37251), ('us', 35489), ('trump', 30876), ('young', 29981), ('go', 29848), ('child', 29600), ('know', 293
37), ('get', 28883), ('right', 27901), ('think', 27435), ('would', 27137), ('old', 25736), ('speech', 24414), ('see', 2424
4), ('girl', 24184), ('un', 23216), ('make', 23035), ('need', 22427), ('children', 22001), ('year', 21508), ('leaders', 2135
9), ('thank', 21357), ('good', 20723), ('planet', 20284), ('time', 20280), ('much', 18420), ('activist', 18208), ('well', 17
978), ('going', 17902), ('many', 17873), ('say', 17827), ('want', 17643), ('way', 17473), ('parents', 17162), ('years', 1695
3), ('take', 16908), ('even', 16837), ('future', 16701), ('really', 16592), ('kids', 16550), ('school', 16257), ('keep', 162
09), ('love', 16061), ('back', 15986), ('great', 15716), ('adults', 15539), ('dare', 15241), ('also', 15105), ('little', 150
18), ('global', 14883), ('stop', 14840), ('let', 14482), ('never', 14144), ('new', 14008), ('news', 13686), ('look', 13593),
('said', 13426), ('today', 13350), ('science', 13226), ('hope', 13221), ('something', 12747), ('better', 12680), ('please',
12671), ('real', 12662), ('nothing', 12443), ('action', 12424), ('could', 11904), ('life', 11796), ('earth', 11540), ('wor
k', 11434), ('believe', 11178), ('used', 11155), ('listen', 11149), ('thing', 11046), ('amazing', 10999), ('done', 10905),
('left', 10701), ('got', 10585), ('still', 10581), ('yes', 10564), ('every', 10517), ('day', 10416), ('help', 10333), ('powe
r', 10280), ('give', 10217), ('use', 10202), ('says', 10126), ('wrong', 10124), ('read', 10074), ('money', 9984), ('made', 9
923), ('kid', 9869), ('support', 9842), ('china', 9840), ('save', 9715), ('someone', 9675), ('nobel', 9657), ('must', 9602),
('using', 9577), ('feel', 9543), ('needs', 9514), ('crisis', 9483), ('summit', 9467), ('person', 9385), ('actually', 9343),
('tell', 9298), ('saying', 9198), ('words', 9177), ('around', 9113), ('come', 9111), ('things', 9091), ('message', 8983),
('childhood', 8903), ('trying', 8878), ('first', 8791), ('ever', 8785), ('prize', 8745), ('strike', 8744), ('making', 8663),
('enough', 8658), ('swedish', 8644), ('generation', 8603), ('everyone', 8543), ('anything', 8509), ('care', 8489), ('sure',
8470), ('age', 8379), ('point', 8348), ('truth', 8307), ('anyone', 8217), ('twitter', 8197), ('problem', 8102), ('men', 809
3), ('watch', 8070), ('talk', 8064), ('man', 8060), ('media', 7953), ('president', 7794), ('scientists', 7791), ('live', 769
2), ('maybe', 7676), ('political', 7661), ('talking', 7654), ('sorry', 7580), ('best', 7543), ('big', 7480), ('understand',
```

**Figure 53:** Tokens and their Frequencies

```
In [20]:  ▶  1  # define a function only to keep words in the top k words
             2  top_k_words,_ = zip(*fdist.most_common(k))
             3  top_k_words = set(top_k_words)
             4  def keep_top_k_words(text):
             5      return [word for word in text if word in top_k_words]

In [21]:  ▶  1  dfJSONN9Jun['TokenCheck2'] = dfJSONN9Jun['TokenCheck2'].apply(keep_top_k_words)
```

**Figure 54:** Method to only keep words in the corpus with frequencies greater than 5

```
In [22]: ▶  1  # document length for the histogram visualisation
            2  dfJSONN9Jun['docLen'] = dfJSONN9Jun['TokenCheck2'].apply(lambda x: len(x))
            3  doc_lengths = list(dfJSONN9Jun['docLen'])
            4  #dfJSONN9Jun.drop(labels='docLen', axis=1, inplace=True)
            5
            6  print("length of list:",len(doc_lengths),
            7        "\naverage document length", np.average(doc_lengths),
            8        "\nminimum document length", min(doc_lengths),
            9        "\nmaximum document length", max(doc_lengths))

         length of list: 743797
         average document length 10.53663432361249
         minimum document length 0
         maximum document length 52
```

**Figure 55:** Examination of average document length

```
In [24]: ▶  1  # plot a histogram of document length
            2  num_bins = 60
            3  fig, ax = plt.subplots(figsize=(12,6));
            4  # the histogram of the data
            5  n, bins, patches = ax.hist(doc_lengths, num_bins, normed=1)
            6  ax.set_xlabel('Document Length (tokens)', fontsize=15)
            7  ax.set_ylabel('Normed Frequency', fontsize=15)
            8  ax.grid()
            9  ax.set_xticks(np.logspace(start=np.log10(2),stop=np.log10(50),num=7, base=10.0))
           10  plt.xlim(0,60)
           11  ax.plot([np.average(doc_lengths) for i in np.linspace(0.0,0,100)], np.linspace(0.0,0.0035,100), '-',
           12        label='average doc length')
           13  ax.legend()
           14  ax.grid()
           15  fig.tight_layout()
           16  plt.show()

         /usr/local/lib/python3.6/dist-packages/matplotlib/axes/_axes.py:6521: MatplotlibDeprecationWarning:
         The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.
           alternative="'density'", removal="3.1")
```



**Figure 56:** Document length distribution

**Figure 57:** Create Token String Variable for Corpus format

Only keep documents with 3 or more tokens.



**Figure 58:** Corpus was examined for Duplicate Documents.

Despite no retweets, there were many long duplicated documents



**Figure 59:** Tweets with no duplicates tweets were saved to a new dataframe



**Figure 60:** Final Corpus Saved

```
length of list: 570512
average document length 11.6181605294893
minimum document length 3
maximum document length 52
```

```
In [9]:    1  # plot a histogram of document length
           2  num_bins = 60
           3  fig, ax = plt.subplots(figsize=(12,6));
           4  # the histogram of the data
           5  n, bins, patches = ax.hist(doc_lengths, num_bins, normed=1)
           6  ax.set_xlabel('Document Length (tokens)', fontsize=15)
           7  ax.set_ylabel('Normed Frequency', fontsize=15)
           8  ax.grid()
           9  #ax.set_xticks(np.logspace(start=np.log10(2),stop=30,num=7, base=10.0))
          10  #ax.set_xticks(np.logspace(start=np.log10(2),stop=30,num=7, base=10.0))
          11  plt.xlim(0,60)
          12  ax.plot([np.average(doc_lengths) for i in np.linspace(0.0,0,35)], np.linspace(0.0,0.0035,35), '-',
          13          label='average doc length')
          14  ax.legend()
          15  ax.grid()
          16  fig.tight_layout()
          17  plt.show()
```

```
/home/ubuntu/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:5: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.
  """
```



**Figure 61:** Average document length after removal of duplicate tweets

# 3 Implementation of Topic Models

**WNTM[1]:**

2. Installation ---------------- Straightforward Java compilation can be done with the following commands: > tar -xzcf wntm.tar.gz > cd wntm > javac *.java

**Figure 62:** Installation of WNTM. Code is unzipped and then Compiled with Java



```
ubuntu@ec2-52-48-23-220.eu-west-1.compute.amazonaws.com:22 - Bitvise xterm - ubuntu@ip-172-31-3-219: ~/wntmTest       —    □    ×
Last login: Thu Aug  6 11:04:54 2020 from 86.45.79.103
(base) ubuntu@ip-172-31-3-219:~$ cd wntmTest/
(base) ubuntu@ip-172-31-3-219:~/wntmTest$ ls
InferenceTopicsForOrgDocs.class  README            sample.words       try2this.words
InferenceTopicsForOrgDocs.java   Tweet.txt         take7Write.txt     try8.adjacent
JGibbLDA-v.1.0                   WTF.theta         try                try8.words
LDARunForWNTM1.theta             model-final.theta try10.adjacent     trythis.adjacent
PrepareInput.class               model-final.twords try10.words        trythis.words
PrepareInput.java                sample.adjacent   try2this.adjacent  wntm.tar.gz
(base) ubuntu@ip-172-31-3-219:~/wntmTest$
```

**Figure 63:** WNTM directory

---

34

```
tmux a -t WNTMTry
ls
nano java -Xmx16g PrepareInput take7Write.txt ./ sample 10
sudo java -Xmx16g PrepareInput take7Write.txt ./ sample 10
ls
nano sample.adjacent
cp sample.adjacent JGibbLDA-v.1.0/testDir/
cd JGibbLDA-v.1.0/
java -Xmx16G  -cp bin:lib/args4j-2.0.6.jar jgibblda.LDA -est -alpha 0.1 -beta 0.01 -ntopics 60 -niters 2000 -savestep 200 -twords 20 -dir try -dfile sample.adjacent
java -Xmx16G  -cp bin:lib/args4j-2.0.6.jar jgibblda.LDA -est -alpha 0.1 -beta 0.01 -ntopics 60 -niters 2000 -savestep 200 -twords 20 -dir testDir -dfile sample.adjacent
ls
java -Xmx16G  -cp bin:lib/args4j-2.0.6.jar jgibblda.LDA -est -alpha 0.1 -beta 0.01 -ntopics 60 -niters 2000 -savestep 200 -twords 20 -dir testDir -dfile sample.adjacent
ls
java -Xmx16G  -cp bin:lib/args4j-2.0.6.jar jgibblda.LDA -est -alpha 0.1 -beta 0.01 -ntopics 60 -niters 2000 -savestep 200 -twords 20  -dfile testDir/sample.adjacent
sudo java -Xmx16G  -cp bin:lib/args4j-2.0.6.jar jgibblda.LDA -est -alpha 0.1 -beta 0.01 -ntopics 60 -niters 2000 -savestep 200 -twords 20 -dir testDir -dfile testDir/sample.adjacent
cd testDir/
ls
nano sample.adjacent
cd ..
sudo java -Xmx16G  -cp bin:lib/args4j-2.0.6.jar jgibblda.LDA -est -alpha 0.1 -beta 0.01 -ntopics 60 -niters 2000 -savestep 200 -twords 20 -dir testDir -dfile sample.adjacent
logout
ls
df -h
cd wntmTest/
```

**Figure 64:** Creation of Word nodes for LDA topic generation

```
sudo java InferenceTopicsForOrgDocs model-final.twords model-final.theta take7Write.txt LDARunForWNTM1.theta
```

**Figure 65:** Inference of theta from LDA with Word Node Corpus onto the documents themselves

**BTM[2]:**

```
unzip BTM-master.zip
ls
cd BTM-master/
ls
cd script/
ls
nano runExample.sh
ls
ps -aux
tmux a -t DMMRunNow
ls
cd BTM-master/
ls
cd script
ls
nano runExample.sh
ls
nano topicDisplay.py
nano indexDocs.py
nano runExample.sh
ls
cd ..
ls
cd ..
ls
cp STTMClone/STTM/take7Write.txt .
ls
cp take7Write.txt BTM-master
cd BTM-master/
ls
cd sample-data/
ls
nano doc_info.txt
cd ..
ls
cd script/
ls
nano runExample.sh
```

**Figure 66:** Unzipping and running BTM

2 https://github.com/xiaohuiyan/BTM/tree/master/src

```bash
#!/bin/bash
# run an toy example for BTM

K=60    # number of topics

alpha=`echo "scale=3;6/$K"|bc`
beta=0.01
niter=2000
save_step=100

input_dir=../sample-data/
output_dir=../output/
model_dir=${output_dir}model/
mkdir -p $output_dir/model

# the input docs for training
doc_pt=${input_dir}take7Write.txt
#doc_pt=take7Write.txt

echo "=============== Index Docs ============="
# docs after indexing
dwid_pt=${output_dir}doc_wids.txt
# vocabulary file
voca_pt=${output_dir}voca.txt
python3 indexDocs.py $doc_pt $dwid_pt $voca_pt

## learning parameters p(z) and p(w|z)
echo "=============== Topic Learning ============="
W=`wc -l < $voca_pt` # vocabulary size
make -C ../src
echo "../src/btm est $K $W $alpha $beta $niter $save_step $dwid_pt $model_dir"
../src/btm est $K $W $alpha $beta $niter $save_step $dwid_pt $model_dir

## infer p(z|d) for each doc
echo "=============== Infer P(z|d)============="
echo "../src/btm inf sum_b $K $dwid_pt $model_dir"
../src/btm inf sum_b $K $dwid_pt $model_dir

## output top words of each topic
echo "=============== Topic Display ============="
python3 topicDisplay.py $model_dir $K $voca_pt
```

**Figure 67:** Parameter Specification for BTM. Corpus location and output location are also specified

## LDA[3]:

```
java -Xmx42G -jar jar/ jLDADMM.jar -model LDA -corpus take7Write.txt -ntopics 60 -
alpha 0.1 -beta 0.01 -initers 2000 -niters 50 -name LDA3July
```

**Figure 68:** Execution of LDA algorithm. Parameters and JVM memory allocation are specified

---

[3] https://github.com/datquocnguyen/jLDADMM

## DMM[4]:

```
java -Xmx42G -jar jar/jLDADMM.jar -model DMM -corpus take7Write.txt -ntopics 60 -
alpha 0.1 -beta 0.01 -initers 2000 -niters 50 -name DMMDat3July
```

**Figure 69:** Execution of DMM algorithm. Parameters and JVM memory allocation are specified

## LF-DMM[5]:

```
java -Xmx42G -jar jar/LFTM.jar -model LFDMM -corpus take7Write.txt -vectors
WordVec100skiptake1.txt -ntopics 60 -alpha 0.1 -beta 0.01 -lambda 0.5 -initers 2000 -
niters 200 -name LFDMM1July
```

**Figure 70:** Execution of LFDMM. Corpus, parameters and Trained word embedding are specified



**Figure 71:** Libraries for word embedding generation and corpus for word embedding

---

```
In [5]:  ▶|  1  sent = [row for row in WordEmbeddingDfAtt['TokenCheck2']]

In [ ]:  ▶|  1  print(sent)

In [6]:  ▶|  1  WordEmbedModel = Word2Vec(sent, min_count=1,size= 100,workers=3, window =5, sg = 1)

In [7]:  ▶|  1  WordEmbedModel.wv.save_word2vec_format('WordVec100skiptake2.txt', binary=False)
```

Figure 72: Word Embedding Training for LFDMM

## GPU-DMM[6]:

```
cp WHUCorp7.txt GPUDMM-master/src/
cd GPUDMM-master//src/
ls
cd ..
cd src/
ls
cd RatioGPUDMM/
ls
nano RatioGPUDMM.java
javac -Xlint *.java
cd ..
java RatioGPUDMM.RatioGPUDMM
```

**Figure 73:** RatioGPUDMM parameters changed and code is recompiled

```
public static void main(String[] args) {

    ArrayList<Document> doc_list = Document.LoadCorpus("WHUCorp7.txt");
    //here
    int num_iter = 2000, save_step = 200;
    double beta = 0.01;
    String similarityFileName = "testing2.txt";
    double weight = 0.1;
    double threshold = 0.75;
    int filterSize = 20;

    for (int round = 1; round <= 5; round += 1) {
        for (int num_topic = 60; num_topic <= 60; num_topic += 20) {
            double alpha = 1.0 * 6 / num_topic;
            RatioGPUDMM gsdmm = new RatioGPUDMM(doc_list, num_topic, num_iter, save_step, beta, alpha, threshold);
            gsdmm.word2idFileName = "FFS.txt";
            gsdmm.topWords = 100;

            //here
            gsdmm.filterSize = filterSize;
            gsdmm.roundIndex = round;
            gsdmm.similarityFileName = similarityFileName;
            gsdmm.weight = weight;
            gsdmm.initNewModel();
            gsdmm.init_GSDMM();
            gsdmm.run_iteration();
            String flag = round+"round_"+num_topic + "topic_weight05_snippet" + "_filter20_iter1000_gpudmm";
            flag = "snippetUpdateNow/" + flag;
            gsdmm.saveModel(flag);
        }
    }
}
```

**Figure 74:** Specification of parameters for GPUDMM

---

[6] https://github.com/WHUIR/GPUDMM

39

**Figure 75:** Corpus Format Preparation for GPUDMM



**Figure 76:** Word Embedding was converted to cosine similarities, as required by GPUDMM

Finding Top 20 weighted words per topic from phi matrix – the matrix was 29000 columns by 60 rows and had to be transformed to extract the words per topic:

```
fprints values > .0025 and puts in a list


 Modified code from Prof. S.S.Shylaja, Head, Dept. of ISE, PESIT.
|
awk '
{
    for (i=1; i<=NF; i++)  {
        a[NR,i] = $i
    }
}
NF>p { p = NF }
END {
    for(j=1; j<=p; j++) {str=a[i,j]
        for(i=1; i<=NR; i++){ if(a[i,j] > 0.0025){
            str= str "   "i" "j" "a[i,j] "\n";
 }
        }
        print str
    }
}' 5round_60topic_weight05_snippet_filter20_iter1000_gpudmm_phi.txt > bugger.txt
```

**Figure 77:** Code to Invert word topic matrix to optimise file

```
sort -t$'\t' -k1,1n -k3'\t'3n phiVals1.txt > phiVals1.Sorted3.txt

sort -t$'\t' -k1,1n phiVals1.txt > phiVals1.Sorted3.txt
```

**Figure 78**: Sort words and weights by topic for output of topic model

# 4 Evaluation of Topic Models

---

**TestPhase1 - find max value**

awk '{m=$1;for(i=1;i<=NF;i++)if($i>=m) { m=$i; idx=I}; print $1,"val"(idx-1),m}' LFLDA3July.theta > test.txt

**keep 3rd column with the max values – isolate it**

awk -F" " '{$1=$2=""; print $3}' testPhase1LFLDA.txt > testPhase2.txt

**concatenate 2 files together**

O
paste testPhase2.txt LFLDA3July.theta > 3PhaseTest.txt

**find the max values and their repeats – get their counts**

awk '{m=$1;count=0;for(i=2;i<=NF;i++)if($i==m) { m=$i; idx=I; count+=1 }; print $1,"val"(idx-1),m, count}'
3PhaseTest.txt > 4PhaseTest.txt

**Examine no of topics per doc – isolate and group their counts**

awk '{a[$4]++;} END{for(i in a) print a[i]" "i}'4PhaseTest.txt > topicFrequencies.txt

---

**Figure 79:** Awk Commands to find document counts per Topic

```
In [1]:     1  import numpy as np
            2  from sklearn.feature_extraction.text import CountVectorizer
            3  import pandas as pd
            4  pd.set_option('display.max_colwidth', -1)
            5  pd.set_option('display.max_rows', None)
            6  from collections import Counter
            7  from sklearn.model_selection import train_test_split
            8  from sklearn.feature_extraction.text import TfidfTransformer
            9  from sklearn import svm
           10

           /home/ubuntu/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:5: FutureWarning: Passing a negative integer is dep
           recated in version 1.0 and will not be supported in future version. Instead, use None to not limit the column width.
             """

In [9]:     1  data1 = pd.read_json('../16JuneTopNoGreta5hundred.json', lines = True)

In [10]:    1  data2 = pd.read_csv('../jLDADMM-master/max_indicesWNTM.txt', header=None)

In [11]:    1  data2.columns = ['label']

In [ ]:     1  data3['date'].sample(100)

In [ ]:     1  %matplotlib inline
            2
            3  data3['date'].value_counts().plot()

In [14]:    1  data3 = pd.concat([data1, data2], axis=1)
```

**Figure 80:** Joining assigned document labels to corpus

**Figure 81:** Addressing Class imbalance in WNTM dataset for SVC classification



**Figure 82:** Classification model training

```
In [42]:  ▶  1  from sklearn import metrics
             2  print(metrics.classification_report(test_y, y_score))

                      precision    recall  f1-score   support

                1         0.84      0.85      0.84       753
                2         0.77      0.84      0.80       847
                3         0.71      0.65      0.68       928
                4         0.74      0.83      0.78       866
                5         0.89      0.91      0.90       885
                6         0.77      0.81      0.79       609
                7         0.69      0.55      0.61        89
                8         0.75      0.66      0.70       598
                9         0.66      0.72      0.69       294
               10         0.68      0.64      0.66       877
               11         0.88      0.88      0.88       718
               12         0.84      0.81      0.83       903
               13         0.83      0.79      0.81       734
               14         0.43      0.38      0.40       916
               15         0.70      0.67      0.68       215
               16         0.73      0.87      0.79       316
               17         0.64      0.51      0.56        95
               18         0.54      0.47      0.51       929
               19         0.82      0.85      0.84       896
               20         0.78      0.85      0.81       351
               21         0.74      0.80      0.77       919
               22         0.95      0.76      0.84       464
               23         0.86      0.92      0.89       956
               24         0.78      0.91      0.84       881
               25         0.81      0.91      0.85       914
               26         0.87      0.81      0.84       527
               27         0.93      0.97      0.95       895
               28         0.93      0.88      0.90       742
               29         0.84      0.79      0.81       393
               30         0.93      0.92      0.92       609
               31         0.45      0.57      0.51       883
               32         0.90      0.81      0.85       329
               33         0.87      0.80      0.83       860
               34         0.75      0.74      0.74       922
               35         0.31      0.39      0.35       893
               36         0.93      0.79      0.85       806
               37         0.91      0.89      0.90       636
               38         0.86      0.87      0.87       893
               39         0.83      0.78      0.80       942
               40         0.77      0.80      0.78       880
               41         0.78      0.77      0.78       759
               42         0.73      0.71      0.72       937
               43         0.89      0.89      0.89       551
               44         0.91      0.95      0.93       896
               45         0.41      0.38      0.40       875
               46         0.93      0.98      0.95       883
               47         0.88      0.89      0.89       516
               48         0.53      0.43      0.48       885
               49         0.74      0.80      0.77       939
               50         0.87      0.89      0.88       794
               51         0.80      0.75      0.77       639
               52         0.91      0.95      0.93       864
               53         0.85      0.78      0.82       384
               54         0.88      0.84      0.86       905
               55         0.78      0.82      0.80       611
               56         0.86      0.88      0.87       830
               57         0.82      0.76      0.79       895
               58         0.84      0.78      0.81       947
               59         0.87      0.92      0.90       922
               60         0.49      0.43      0.46       901

         accuracy                             0.77     44096
        macro avg         0.78      0.77      0.77     44096
     weighted avg         0.77      0.77      0.77     44096
```

**Figure 83:** Classification recall results per topic

```
In [43]: ▶  1  np.savetxt("Two_y_score3000MaxOutWNTMTDIDFFea.csv", y_score, delimiter=",")
            2  np.savetxt("Two_test_y3000MaxOutWNTMTDIDFFea.csv", test_y, delimiter=",")

In [6]:  ▶  1  from numpy import genfromtxt
            2  One_y_score = genfromtxt('y_score3000MaxOutWNTMTDIDFFea.csv', delimiter=',')

In [7]:  ▶  1  from numpy import genfromtxt
            2  One_test_y = genfromtxt('test_y3000MaxOutWNTMTDIDFFea.csv', delimiter=',')

In [8]:  ▶  1  from sklearn import metrics
            2  print(metrics.classification_report(One_test_y, One_y_score))
```

**Figure 84:** Saving recall result to file

It was decided to set the coherence window to 15 upon examination of average document length per topic. For all models, a few topics had average document length counts slightly above 15. To accommodate these topics in the coherence windows, the size of the context window was set to 15.

```
In [3]:  ▶  1  data1 = pd.read_json('../16JuneTopNoGreta5hundred.json', lines = True)

In [6]:  ▶  1  data2 = pd.read_csv('../ClassificationRuns/LDAHardVals.txt', header=None)

In [ ]:  ▶  1  data2 = pd.read_csv('../ClassificationRuns/LDAHardVals.txt', header=None)

In [14]: ▶  1  data2 = pd.read_csv('../ClassificationRuns/max_indicesBTM.txt', header=None)

In [20]: ▶  1  data2 = pd.read_csv('../LFTM-master/FDMMTopicAssignmentsCleaned1.txt', header=None)

In [26]: ▶  1  data2 = pd.read_csv('../ClassificationRuns/max_indicesBeta01.LABEL', header=None) #gpudmm

In [34]: ▶  1  data2 = pd.read_csv('../jLDADMM-master/DMMDat3JulyTopicLabel.txt', header=None)

In [35]: ▶  1  data2.columns = ['label']

In [ ]:  ▶  1  data3['date'].sample(100)

In [ ]:  ▶  1  %matplotlib inline
            2
            3  data3['date'].value_counts().plot()

In [36]: ▶  1  data3 = pd.concat([data1, data2], axis=1)

In [18]: ▶  1  data3.drop('label', axis=1)
                                    . . .

In [ ]:  ▶  1  #data3.sample(5000).to_json('20July5000TestWNTMGraphPlot.json', orient='records', lines = True)

In [ ]:  ▶  1  #data3.to_csv('ALL_CV20JulyWNTMGraphPlot.csv')

In [27]: ▶  1  Filter27 = [3]
            2  data3FilterBool = data3.label.isin(Filter27)
            3  data3_SomeTopic_s = data3[data3FilterBool]

In [10]: ▶  1  FilterLDA = [1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,
            2                23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,
            3                42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60]
```

**Figure 85:** Loading of consecutive label allocation per topic model method

```
In [37]:  ▶  1  FilterLDA
              2  #N = 100
              3  for  xxx in FilterLDA:
              4      Filter27 = [xxx]
              5      data3FilterBool = data3.label.isin(Filter27)
              6      data3_SomeTopic_s = data3[data3FilterBool]
              7      data3_SomeTopic_s['docLen'] = data3_SomeTopic_s['TokenCheck2'].apply(lambda x: len(x))
              8      doc_lengths = list(data3_SomeTopic_s['docLen'])
              9
             10      print(xxx,"\nlength of list:",len(doc_lengths),
             11          "\naverage document length", np.average(doc_lengths),
             12          "\nminimum document length", min(doc_lengths),
             13          "\nmaximum document length", max(doc_lengths), "\n")
```

```
/home/ubuntu/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:7: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view
-versus-a-copy
  import sys
```

```
1
length of list: 1531
average document length 7.859568909209667
minimum document length 3
maximum document length 28

2
length of list: 9259
average document length 13.973539259099255
minimum document length 3
maximum document length 38

3
length of list: 658
average document length 7.729483282674772
minimum document length 3
maximum document length 28

4
length of list: 21755
average document length 10.617421282463802
minimum document length 3
maximum document length 33

5
length of list: 5239
average document length 9.729719412101547
minimum document length 3
maximum document length 32
```

**Figure 86:** Average document length per topic for each method

```
In [1]:  ▶  1  import numpy as np
            2  from sklearn.feature_extraction.text import CountVectorizer
            3  import pandas as pd
            4  pd.set_option('display.max_colwidth', -1)
            5  from collections import Counter
            6  from gensim.models.coherencemodel import CoherenceModel
            7  from gensim.corpora.dictionary import Dictionary
            8  from numpy import array
```

```
/home/ubuntu/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:5: FutureWarning: Passing a negative integer is dep
recated in version 1.0 and will not be supported in future version. Instead, use None to not limit the column width.
  """
```

```
In [2]:  ▶  1  data1 = pd.read_json('~/16JuneTopNoGreta5hundred.json', lines = True)
```

```
In [3]:  ▶  1  textData = data1['TokenCheck2']
```

```
In [4]:  ▶  1  CreatedDictionary = Dictionary(textData)
```

```
In [5]:  ▶  1  CreatedCorpus = [CreatedDictionary.doc2bow(doc) for doc in data1['TokenCheck2']]
```

```
In [6]:  ▶  1  from csv import reader
            2  # read csv file as a list of lists
            3  with open('WTNMM_DataDumpText.prn', 'r') as read_obj:
            4      # pass the file object to reader() to get the reader object
            5      csv_reader = reader(read_obj)
            6      # Pass reader object to list() to get a list of lists
            7      list_of_rows = list(csv_reader)
            8      #print(list_of_rows)
```

**Figure 87:** Data preparation for coherence testing

```
In [8]:    1  #from gensim.test.utils import common_corpus, common_dictionary
           2  from gensim.models.coherencemodel import CoherenceModel
           3   topics = ListJerr
           4
           5  cm = CoherenceModel(topics=ListJerr, texts=textData, dictionary=CreatedDictionary, window_size=15, coherence='c_v')
           6  coherence = cm.get_coherence()  # get coherence value

In [ ]:    1  print(ListJerr)

In [9]:    1  coherence

Out[9]:   0.49974832915795103

In [10]:   1  cm.get_coherence_per_topic()
                                      ...

In [9]:    1  import csv
           2  with open('WNTM_CoherenceScores.csv', 'w') as CoherenceFile:
           3      wr = csv.writer(CoherenceFile, delimiter="\n",quoting=csv.QUOTE_NONE)   #QUOTE_NONE
           4      wr.writerow(cm.get_coherence_per_topic())
```

**Figure 88:** Meta-weighted average coherence for WNTM topics

```
In [4]:  ▶  1  with open('LDA_CoherenceScores.csv') as f:
            2      w = [float(x) for x in next(f).split()]
            3      LDAScoreArrayarray = [[float(x) for x in line.split()] for line in f]
            4  LDAScoreArrayFlat = []
            5  for lowerList in LDAScoreArrayarray:
            6      for element in lowerList:
            7          LDAScoreArrayFlat.append(element)

In [5]:  ▶  1  LDAScoreArrayFlat
```

```
                                        . . .
```

```
In [7]:  ▶  1  with open('LF_DMM_CoherenceScores.csv') as f:
            2      w = [float(x) for x in next(f).split()]
            3      LF_DMMScoreArrayarray = [[float(x) for x in line.split()] for line in f]
            4  LF_DMMScoreArrayFlat = []
            5  for lowerList in LF_DMMScoreArrayarray:
            6      for element in lowerList:
            7          LF_DMMScoreArrayFlat.append(element)
```

```
In [8]:  ▶  1  with open('GPU_DMM_CoherenceScores.csv') as f:
            2      w = [float(x) for x in next(f).split()]
            3      GPU_DMMScoreArray = [[float(x) for x in line.split()] for line in f]
            4  GPU_DMMScoreArrayFlat = []
            5  for lowerList in GPU_DMMScoreArray:
            6      for element in lowerList:
            7          GPU_DMMScoreArrayFlat.append(element)
```

```
In [9]:  ▶  1  with open('DMM_CoherenceScores.csv') as f:
            2      w = [float(x) for x in next(f).split()]
            3      DMMScoreArray = [[float(x) for x in line.split()] for line in f]
            4  DMMScoreArrayFlat = []
            5  for lowerList in DMMScoreArray:
            6      for element in lowerList:
            7          DMMScoreArrayFlat.append(element)
```

```
In [10]:  ▶  1  with open('BTM_CoherenceScores.csv') as f:
             2      w = [float(x) for x in next(f).split()]
             3      BTMScoreArray = [[float(x) for x in line.split()] for line in f]
             4  BTMScoreArrayFlat = []
             5  for lowerList in BTMScoreArray:
             6      for element in lowerList:
             7          BTMScoreArrayFlat.append(element)
```

```
In [11]:  ▶  1  with open('WNTM_CoherenceScores.csv') as f:
             2      w = [float(x) for x in next(f).split()]
             3      WNTMScoreArray = [[float(x) for x in line.split()] for line in f]
             4  WNTMScoreArrayFlat = []
             5  for lowerList in WNTMScoreArray:
             6      for element in lowerList:
             7          WNTMScoreArrayFlat.append(element)
```

**Figure 89:** Recall based on Topic Proportions (extracted plot of documents < 7000)

Loading of saved coherence score distributions per topic model method into Lists for post-hoc analysis

```
In [12]:  ▶  1  stats.kruskal(BTMScoreArrayFlat, WNTMScoreArrayFlat,DMMScoreArrayFlat,GPU_DMMScoreArrayFlat,LF_DMMScoreArrayFlat,LDAScore

Out[12]:  KruskalResult(statistic=83.9192783621761, pvalue=5.526965504325747e-16)

In [18]:  ▶  1  CoherencePH6 = [BTMScoreArrayFlat, WNTMScoreArrayFlat,DMMScoreArrayFlat,GPU_DMMScoreArrayFlat,LF_DMMScoreArrayFlat,LDAS
            2  sp.posthoc_conover(CoherencePH6, p_adjust = 'fdr_tsbh').to_csv('DoesThisWrite.csv')
            3  sp.posthoc_conover(CoherencePH6, p_adjust = 'fdr_tsbh')

Out[18]:
```

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | -1.000000e+00 | 6.278876e-02 | 4.846532e-02 | 0.002346 | 1.051897e-10 | 5.136865e-02 |
| 2 | 6.278876e-02 | -1.000000e+00 | 1.531696e-03 | 0.000024 | 1.695165e-14 | 2.942244e-01 |
| 3 | 4.846532e-02 | 1.531696e-03 | -1.000000e+00 | 0.078695 | 6.106969e-07 | 1.062163e-03 |
| 4 | 2.346491e-03 | 2.354119e-05 | 7.869510e-02 | -1.000000 | 8.287473e-05 | 1.496790e-05 |
| 5 | 1.051897e-10 | 1.695165e-14 | 6.106969e-07 | 0.000083 | -1.000000e+00 | 1.206095e-14 |
| 6 | 5.136865e-02 | 2.942244e-01 | 1.062163e-03 | 0.000015 | 1.206095e-14 | -1.000000e+00 |

**Figure 90:** Post-hoc testing on coherence scores

```
In [1]:  ▶  1  %matplotlib inline
            2  import logging
            3  logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.INFO)
            4  from gensim.corpora import Dictionary
            5  import matplotlib.pyplot as plt
            6  import csv
            7  from csv import reader
            8

In [2]:  ▶  1  with open('DistanceFiles/WNTM_DumpText.csv') as f:
            2      mylist = [tuple(map(float, i.split(','))) for i in f]

In [3]:  ▶  1  from gensim.matutils import hellinger

In [16]: ▶  1  from gensim.matutils import jaccard

In [ ]:  ▶  1  len(mylist)

In [5]:  ▶  1
            2
            3  letsSee = [mylist[i:i+20] for i in range(0, len(mylist), 20)]

In [6]:  ▶  1  len(letsSee)

Out[6]:  60

In [6]:  ▶  1  letsSee[59]

Out[6]:  [(227.0, 0.0109),
          (22.0, 0.009),
          (63.0, 0.008),
          (650.0, 0.0075),
          (73.0, 0.0075),
          (64.0, 0.0064),
          (27.0, 0.0059),
          (288.0, 0.0059),
          (20.0, 0.0055),
          (154.0, 0.0052),
          (83.0, 0.0048),
          (480.0, 0.0047),
          (10.0, 0.0045),
          (615.0, 0.0041),
          (612.0, 0.004),
          (142.0, 0.0039),
          (416.0, 0.0039),
          (410.0, 0.0037),
          (477.0, 0.0037),
          (656.0, 0.0035)]
```

**Figure 91:** Generation of distances amongst topics with their word weightings

```
In [8]:  ▶  1  dist = lambda p1, p2: hellinger(p1, p2)
            2  dm = np.asarray([[dist(p1, p2) for p2 in letsSee] for p1 in letsSee])

In [12]: ▶  1  dm

Out[12]: array([[0.       , 0.61185315, 0.66836125, ..., 0.64202903, 0.74194741,
                 0.66267262],
                [0.61185315, 0.       , 0.37513808, ..., 0.33885445, 0.48139125,
                 0.37894591],
                [0.66836125, 0.37513808, 0.       , ..., 0.40242918, 0.55394621,
                 0.43771583],
                ...,
                [0.64202903, 0.33885445, 0.40242918, ..., 0.       , 0.53271593,
                 0.41447557],
                [0.74194741, 0.48139125, 0.55394621, ..., 0.53271593, 0.       ,
                 0.54042576],
                [0.66267262, 0.37894591, 0.43771583, ..., 0.41447557, 0.54042576,
                 0.       ]])

In [9]:  ▶  1  #flatten array to create a list
            2  flatDm = dm.flatten()
```

**Figure 92:** Array flattening in preparation for Multi-dimensional Scaling

```
In [10]: ▶  1  %matplotlib inline
            2
            3  #data3['date'].value_counts().plot()
            4
            5
            6  flatDm
            7  #plt.hist(flatDm, bins = [0.01, 0.5, 0.8])
            8  #n, bins, patches = plt.hist(flatDm)
            9  plt.hist(flatDm,  bins=100)
           10  plt.show()
```

```
In [12]: ▶  1  %matplotlib inline
            2
            3  #data3['date'].value_counts().plot()
            4
            5
            6  flatDm
            7  #plt.hist(flatDm, bins = [0.01, 0.5, 0.8])
            8  #n, bins, patches = plt.hist(flatDm)
            9  plt.hist(flatDm, normed=True, bins=100)
           10  plt.show()
```

/home/ubuntu/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:9: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.
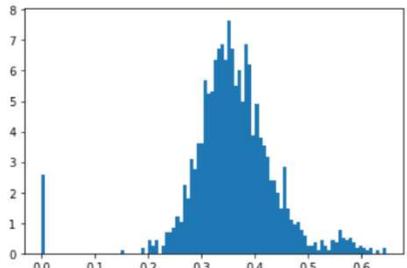  if __name__ == '__main__':



**Figure 93:** Distribution of distance similarities

```
In [32]:  ▶  1  mds_model = manifold.MDS(n_components = 2, random_state = 123,
              2      dissimilarity = 'precomputed')
              3  mds_fit = mds_model.fit(dm)
              4  mds_coords = mds_model.fit_transform(dm)

In [28]:  ▶  1  mds_model = manifold.MDS(n_components =2 ,
              2      dissimilarity = 'precomputed')
              3  mds_fit = mds_model.fit(dm)
              4  mds_coords = mds_model.fit_transform(dm)

In [127]: ▶  1  len(mds_coords)

Out[127]: 60

In [13]:  ▶  1  TopicNumbers = [str(x) for x in range(60)]

In [24]:  ▶  1  %matplotlib notebook
              2  fig = plt.figure(figsize=(30, 16))
              3  #fig.suptitle("Manifold Learning with %i points, %i neighbors"
              4  #              % (1000, n_neighbors), fontsize=14)
              5
              6  # Add 3d scatter plot
              7
              8  ax = fig.add_subplot(251, projection='3d')
              9  #n = 100
             10  #for c, m, zl, zh in [('r', 'o', -60, -25), ('b', '^', -30, -5)]:
             11  #      xs = randrange(n, 23, 50)
             12  #      ys = randrange(n, 0, 100)
             13  #      zs = randrange(n, zl, zh)
             14  #      ax.scatter(xs, ys, zs, c=c, marker=m)
             15  ax.scatter(mds_coords[:,0],mds_coords[:,1],mds_coords[:,2], cmap=plt.cm.Spectral)
             16  ax.view_init(20, -72)
```

**Figure 94**: Multidimensional Scaling specification

```
In [33]:  ▶  1  plt.figure()
              2  plt.scatter(mds_coords[:,0],mds_coords[:,1],
              3      facecolors = 'none', edgecolors = 'none')  # points in white (invisible)
              4  labels = TopicNumbers
              5  for label, x, y in zip(TopicNumbers, mds_coords[:,0], mds_coords[:,1]):
              6      plt.annotate(label, (x,y), xycoords = 'data')
              7  plt.xlabel('First Dimension')
              8  plt.ylabel('Second Dimension')
              9  plt.title('Dissimilarity amongst topics')
             10  plt.show()
```



**Figure 95:** Multidimensional scaling plot projected onto two dimensions

**Figure 96:** Creation of Word Features to be used in Chi test



**Figure 97:** Word vector preparation for examination of significant words per topic according to $Chi^2$ test of independence

```
In [17]:  ▶  1  FilterStrongWNTM
            2  N = 20
            3  for  xxx in FilterStrongWNTM:
            4      chiVals = chi2(tfidfTestChi, (labels == xxx))
            5      indices = np.argsort(chiVals[0])# gives some number, a  ranking to scores
            6      Pvals = chiVals[1]
            7      sortedPvals = Pvals[indices]
            8      feature_names = np.array(tfidf.get_feature_names())[indices] # find the names from word vec, order by indices ra
            9      words = [w for w in feature_names]
           10      a = words[-N:]
           11      b = sortedPvals[-N:]
           12      res = [str(i) + " " + str(j) for i, j in zip(a, b)]
           13      print("#############################################  "+ str(xxx)+ "  ###################################
           14      print("Terms Most Correlated:\n {}".format('\n. '.join(res)))
```

```
#################################################  1  #########################################
Terms Most Correlated:
 dears 0.0
. incandescent 0.0
. males 0.0
. freaking 0.0
. irrefutable 0.0
. emergence 0.0
. male 0.0
. represents 0.0
. misogyny 0.0
. white 0.0
. demon 0.0
. hurtling 0.0
. obsolescence 0.0
. convergence 0.0
. hints 0.0
. certain 0.0
. middle 0.0
. aged 0.0
. men 0.0
. triggering 0.0
#################################################  5  #########################################
Terms Most Correlated:
 trip 0.0
. jets 0.0
. hummer 0.0
. plane 0.0
. flying 0.0
. private 0.0
. sailboat 0.0
. sailing 0.0
```

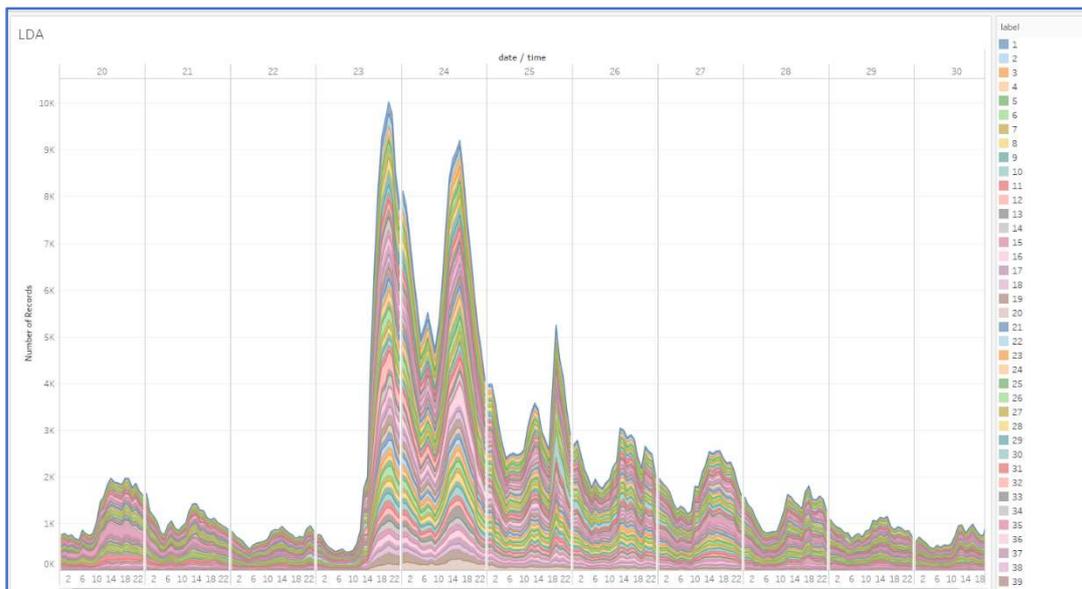**Figure 98:** Chi-squared test and relevant significant words per topic



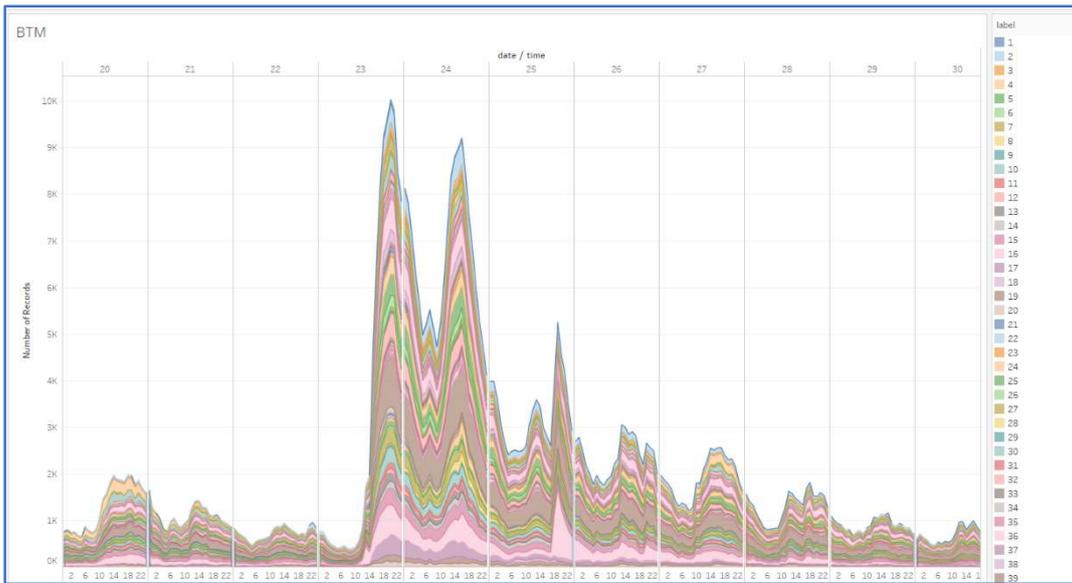**Figure 99:** Time series plot of the topics modelled by LDA

53

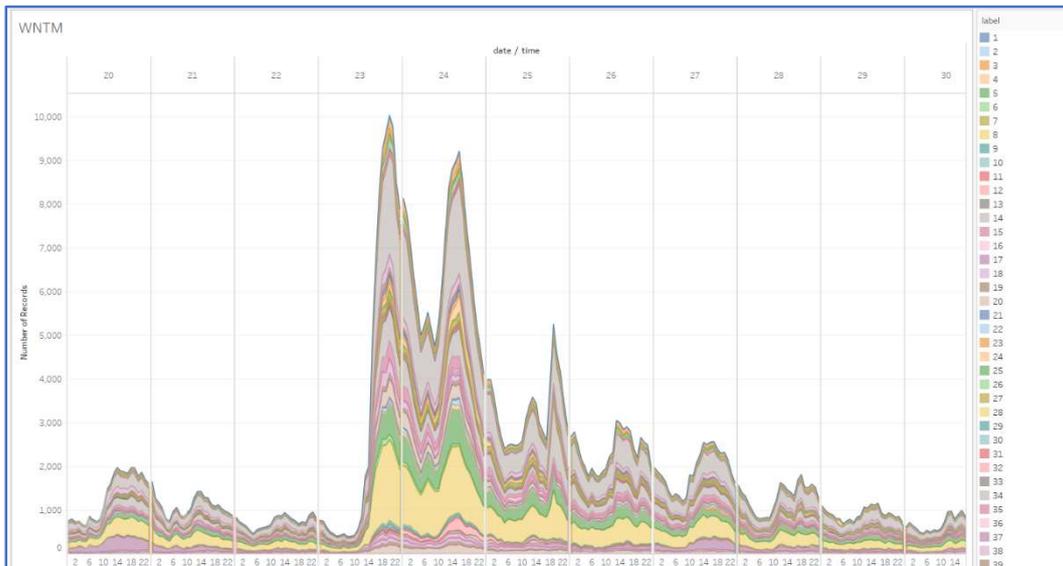**Figure 100:** Time series plot of the topics modelled by BTM



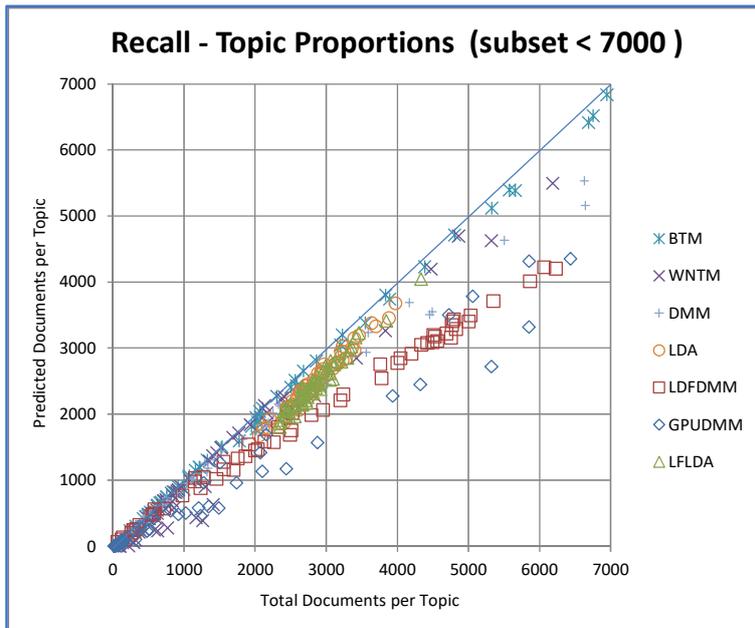**Figure 101:** Time series plot of the topics modelled by WNTM

**Figure 102:** Recall based on Topic Proportions (extracted plot of documents < 7000)
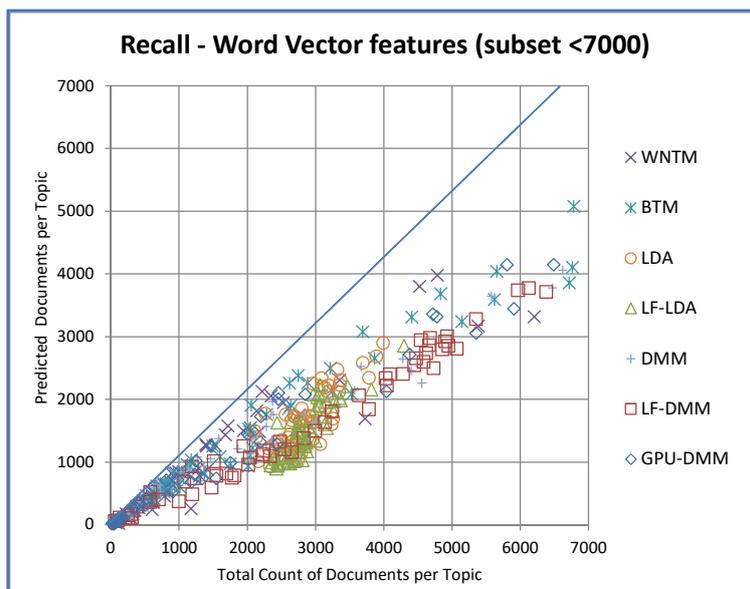


**Figure 103:** Recall based on Word Vectors (extracted plot of documents < 7000)