

```
In [1]: 1 # !pip install gensim==3.8.3
        2 # !pip install pyLDAvis==3.3.1
```

```
In [2]: 1 import nltk
        2 from nltk.tokenize import word_tokenize, sent_tokenize
        3 from IPython.display import HTML, display
        4 import tabulate
        5 import pandas as pd
        6 import numpy as np
        7 from PIL import Image
        8 from wordcloud import WordCloud
        9 import seaborn as sns
       10 from gensim.models.coherencemodel import CoherenceModel
       11 import pyLDAvis
       12
       13 # from nltk.tokenize import word_tokenize, sent_tokenize
       14 from nltk.corpus import stopwords
       15 from nltk.stem import WordNetLemmatizer, PorterStemmer
       16 # from nltk.stem.porter import *
       17
       18 import gensim
       19 from gensim.models import Phrases
       20 #Prepare objects for LDA gensim implementation
       21 from gensim import corpora
       22 #Running LDA
       23 from gensim import models
       24
       25 import warnings
       26 warnings.filterwarnings('ignore')
       27
       28 import matplotlib.pyplot as plt
       29 import re
       30 %matplotlib inline
```

C:\Users\devas\anaconda3\lib\site-packages\gensim\similarities__init__.py:15: UserWarning: The gensim.similarities.levenshtein submodule is disabled, because the optional Levenshtein package <<https://pypi.org/project/python-Levenshtein/>> is unavailable. Install Levenshtein (e.g. `pip install python-Levenshtein`) to suppress this warning.
warnings.warn(msg)

```
In [3]: 1 data=pd.read_csv('C:/Users/devas/Documents/Pattern/omicron_tweets.csv', engine='python')
        2 # dropping empty rows
        3 data = data.dropna(subset=['text'])
        4 # dropping duplicates
        5 data = data.drop_duplicates()
        6 print('Rows: {}, columns: {}'.format(data.shape[0], data.shape[1]))
        7 data.head(2)
```

Rows: 8066, columns: 13

Out[3]:

	tweet_id	date	text	user_name	user_location	user_description	user_created	user_followers	user_friends	user_favourites	hashtags	source	is_retweet
0	1465162850457071621	2021-11-29 03:36:59+00:00	Will Boris Johnson ever learn his lesson after...	James robertson	NaN	NaN	2013-04-22 19:41:30+00:00	303.0	188	84059	NaN	Twitter for Android	False
1	1465162841665769474	2021-11-29 03:36:57+00:00	#Omicron so hot; it's being rush shipped in De...	Gene Bulmer	East Coast	#ChristFollower, #Leader, Husband, Father, #Ru...	2014-11-10 22:43:27+00:00	99.0	162	745	['Omicron']	Twitter for iPhone	False

Data Manipulation

In [4]:

```
1 data['processed'] = ''
2 nltk.download('wordnet')
3 nltk.download('stopwords')
4
5 stop_words = stopwords.words('english')
6
7 stop_words.extend(['https'])
8 # initializing the wordnet lemmatizer
9 lm = WordNetLemmatizer()
10
11 def processing(content):
12
13     content = content.replace('\n', ' ').split(' ')
14     # removing these punctuations from tokens like it will convert the word mode? into mode
15     rx = re.compile('([&#.?!\-()])*)')
16     content = [rx.sub('', word) for word in content]
17
18     # removing stopwords
19     content = [word.strip().lower() for word in content if word.strip().lower() not in stop_words]
20     # remove words whose length is greater than 1 and or alphabets only
21     content = [word for word in content if len(word)>1 and word.isalpha()]
22     # lemmatizing the words to their basic form
23     content = [lm.lemmatize(word) for word in content]
24
25     return ' '.join(content)
26
27 for k in range(len(data)):
28     data.iloc[k,-1] = processing(data.iloc[k,2])
29 # processed data
30 data.head()
```

```
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\devas\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\wordnet.zip.
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\devas\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
```

Out[4]:

	tweet_id	date	text	user_name	user_location	user_description	user_created	user_followers	user_friends	user_favourites	hashtags	source	is_retweet
0	1465162850457071621	2021-11-29 03:36:59+00:00	Will Boris Johnson ever learn his lesson after...	James robertson	NaN	NaN	2013-04-22 19:41:30+00:00	303.0	188	84059	NaN	Twitter for Android	False
1	1465162841665769474	2021-11-29 03:36:57+00:00	#Omicron so hot; it's being rush shipped in De...	Gene Bulmer	East Coast	#ChristFollower, #Leader, Husband, Father, #Ru...	2014-11-10 22:43:27+00:00	99.0	162	745	['Omicron']	Twitter for iPhone	False
2	1465162837047726094	2021-11-29 03:36:56+00:00	@JoeBiden How To Stop The Spread of Coronaviru...	Andrew Arcie Galendez (F+WRTIP-)	NaN	I AM PURE FILIPINO; BOTH MY PARENTS ARE FILIPI...	2021-01-28 11:31:59+00:00	114.0	3223	4830	NaN	Twitter for Android	False
3	1465162836498432005	2021-11-29 03:36:56+00:00	Gold Coast leaders urge against overcautious r...	myGC.com.au	Gold Coast, Australia	The Gold Coast's best news, local events, weat...	2009-03-30 03:54:34+00:00	12232.0	3191	271	['GoldCoastNews']	Zapier.com	False
4	1465162831444119562	2021-11-29 03:36:55+00:00	What Is Omicron Radiation...nhhttps://t.co/QfK...	Ghost Of Gus Hall	Poughkeepsie Peoples Republic	Freelancer who knows a name and is lying to po...	2011-10-17 03:25:39+00:00	603.0	960	819	NaN	Twitter for Android	False

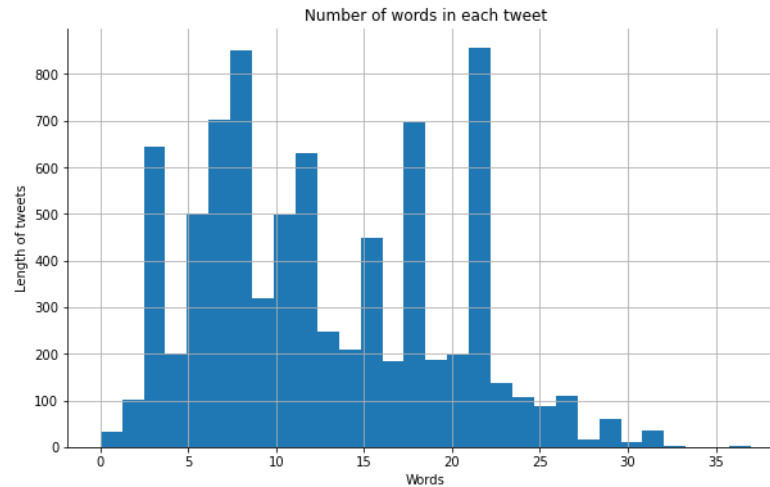
```

1 b = data['processed'].tolist()
2 b = ' '.join(map(str, b))
3 # b = b.replace(',', ' ').lower().replace('-', '')
4
5 wordcloud = WordCloud(max_font_size=40, max_words=1000, background_color="white", random_state=100,
6                       prefer_horizontal=0.60).generate(b.lower())
7 # max_font_size=40, max_words=1000, background_color="white",
8   # random_state=100, prefer_horizontal=0.50
9 plt.figure(figsize=(12,8))
10 plt.imshow(wordcloud, interpolation="bilinear")
11 plt.axis("off")
12 # plt.show()
13 plt.savefig('word_cloud.jpg')

```



```
In [6]: 1 lengths = [len(sentence.split(' ')) if len(sentence)>0 else 0 for sentence in data['processed']]
2 data['Lengths'] = lengths
3 plt.figure(figsize=(10,6))
4 data['Lengths'].hist(bins=30)
5 sns.despine(top=True, right=True, left=False, bottom=False)
6 plt.title('Number of words in each tweet')
7 plt.xlabel('Words')
8 plt.ylabel('Length of tweets')
9 plt.savefig('length_tweets.jpg')
10 plt.show()
```



```
In [7]: 1 import gensim.corpora as corpora
2
3 #decomposing sentences into tokens
4 tokens = [sentence.split(' ') for sentence in data['processed']]
5 # training a bi gram model in order to include those bigrams as tokens who occurred at least 6 times
6 # in the whole dataset
7 bigram = gensim.models.Phrases(tokens, min_count=2, threshold=100)
8 bigram_mod = gensim.models.phrases.Phraser(bigram)
9
10 # including bigrams as tokens
11 sents = [bigram_mod[token] for token in tokens]
12
13 # Create Dictionary to keep track of vocab
14 dct = corpora.Dictionary(tokens)
15
16 print('Unique words before filtering/after pre-processing', len(dct))
17 # no_below= 30
18 # filter the words that occur in less than 3 documents and in more than 60% of documents
19 dct.filter_extremes(no_below=3, no_above=0.60)
20 print('Unique words after filtering', len(dct))
21
22 # Create Corpus
23 corpus = [dct.doc2bow(sent) for sent in sents]
24
25 tfidf = gensim.models.TfidfModel(corpus)
26 corpus_tfidf = tfidf[corpus]
```

Unique words before filtering/after pre-processing 9834
Unique words after filtering 3263

```

In [8]: 1 %%time
2 from gensim.models import CoherenceModel
3 import time
4 import os
5
6 scores = []
7 for k in range(3,15):
8     # LDA model
9     lda_model = gensim.models.LdaModel( corpus=corpus_tfidf, num_topics=k,
10                                         id2word=dct, random_state=12)
11     # to calculate score for coherence
12     coherence_model_lda = CoherenceModel(model=lda_model, texts=sents, dictionary=dct, coherence='c_v')
13     coherence_lda = coherence_model_lda.get_coherence()
14     print(k, coherence_lda)
15     scores.append(coherence_lda)

```

```

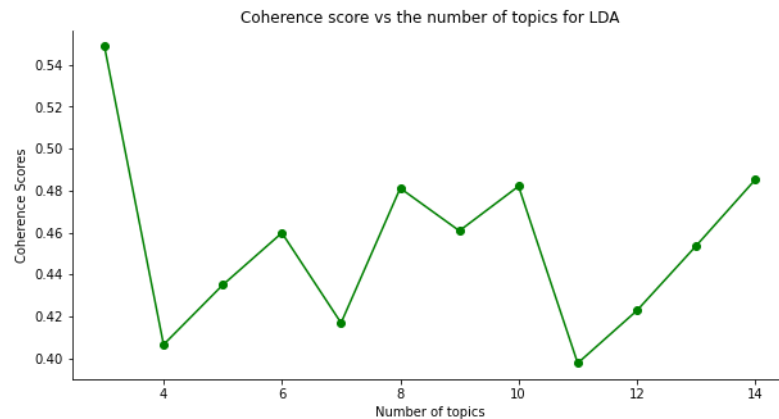
3 0.5488264668187827
4 0.4065263368076506
5 0.43498981193412023
6 0.4598313128230022
7 0.41688548620270843
8 0.48115228053009207
9 0.4608086515564191
10 0.4820849819803641
11 0.3976712431810696
12 0.4227588296332713
13 0.45360036270778586
14 0.48519719309739046
Wall time: 1min 17s

```

```

In [9]: 1 selected_topics = np.argmax(scores)+3
2
3 plt.figure(figsize=(10, 5))
4 plt.plot(list(range(3,15)), scores, marker='o', color='green')
5 sns.despine(top=True, right=True, left=False, bottom=False)
6
7 plt.locator_params(integer=True)
8 plt.title('Coherence score vs the number of topics for LDA')
9 plt.xlabel('Number of topics')
10 plt.ylabel('Coherence Scores')
11 plt.savefig('lda_scores.jpg')
12 plt.show()

```



```

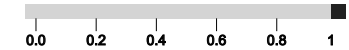
In [10]: 1 import pyLDavis.gensim_models
2
3 lda_model = gensim.models.LdaModel(corpus=corpus_tfidf, id2word=dct, num_topics=selected_topics,
4                                     random_state=12, chunksize=128, passes=10 )
5
6 pyLDavis.enable_notebook()
7 results = pyLDavis.gensim_models.prepare(lda_model, corpus_tfidf, dct, sort_topics=False)
8 pyLDavis.save_html(results, 'ldavis_english' + '.html')
9 results

```

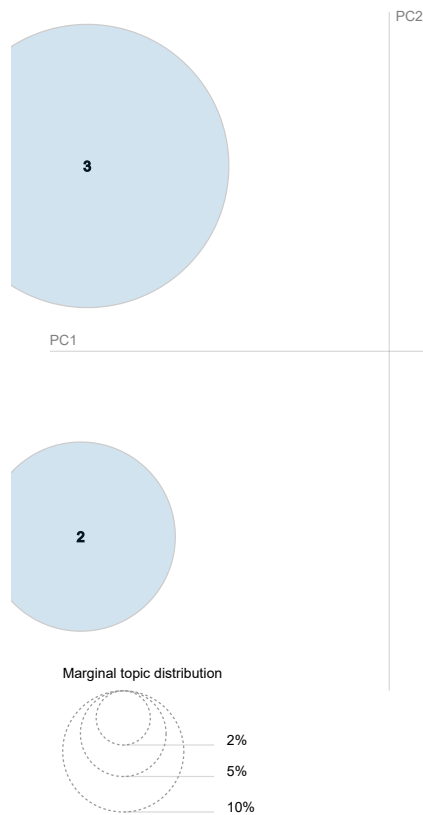
Out[10]:

Selected Topic:

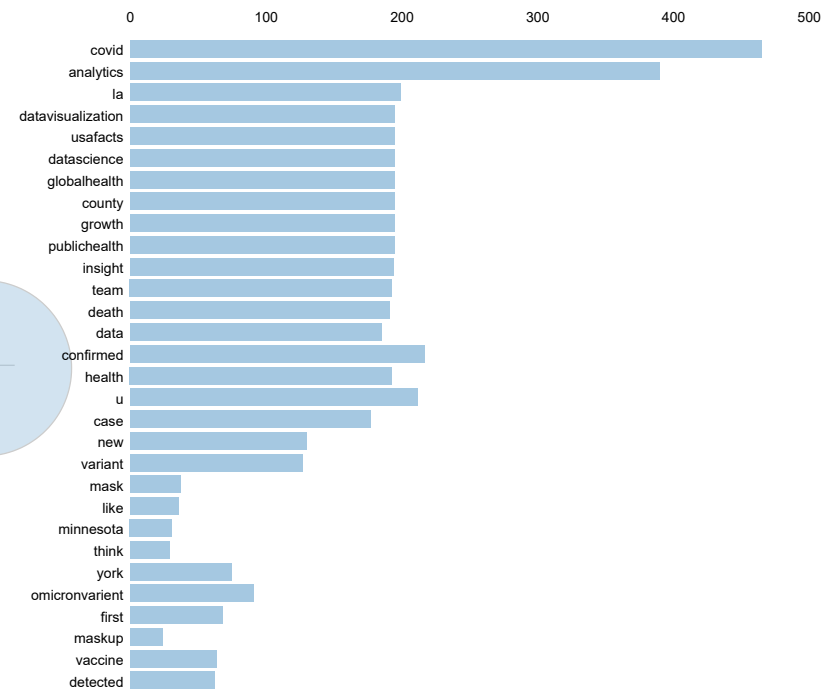
Slide to adjust relevance metric:(2)
 $\lambda = 1$



Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Salient Terms⁽¹⁾



Overall term frequency

Estimated term frequency within the selected topic

1. $saliency(\text{term } w) = \text{frequency}(w) * [\sum_t p(t|w) * \log(p(t|w)/p(t))]$ for topics t : see Chuang et. al (2012)

2. $relevance(\text{term } w | \text{topic } t) = \lambda * p(w|t) + (1 - \lambda) * p(w|t)/p(w)$: see Sievert & Shirley (2014)

```
In [11]: 1 # top words in each topic
         2 lda_model.print_topics()
```

```
Out[11]: [(0,
          '0.106*covid" + 0.095*analytics" + 0.049*la" + 0.048*usafacts" + 0.048*datascience" + 0.048*datavisualization" + 0.048*globalhealth" + 0.048*county" + 0.048*growth" + 0.048*publichealth'),
          (1,
          '0.008*mask" + 0.008*like" + 0.008*omicronvariant" + 0.006*minnesota" + 0.006*omicronvariant" + 0.006*think" + 0.005*maskup" + 0.005*anyone" + 0.004*auspol" + 0.004*biden'),
          (2,
          '0.017*case" + 0.012*new" + 0.012*variant" + 0.007*york" + 0.006*first" + 0.006*vaccine" + 0.006*state" + 0.006*detected" + 0.006*coronavirus" + 0.005*omicronvariant')]
```

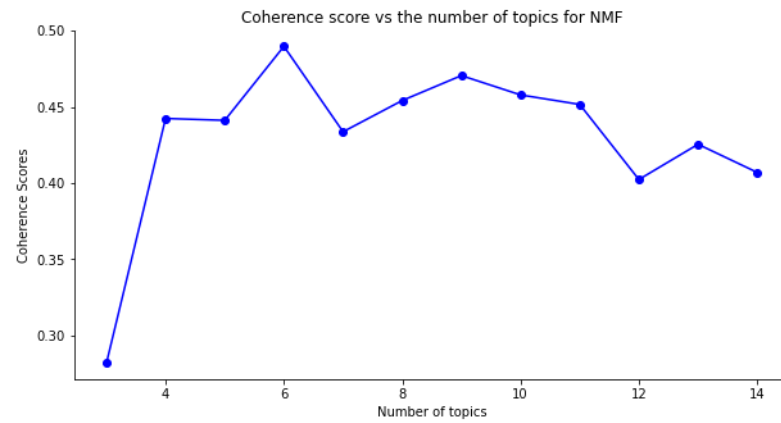
Non-Negative Matrix Factorization (NMF)

```
In [12]: 1 from gensim.models.nmf import Nmf
         2
         3 scores_nmf = []
         4 for k in range(3,15):
         5     # lda_mallet model
         6     nmf_model = Nmf(corpus_tfidf, num_topics=k, \
         7                  id2word=dct, \
         8                  passes=10)
         9     # to calculate score for coherence
        10     coherence_model_lda = CoherenceModel(model=nmf_model, texts=sents, dictionary=dct, coherence='c_v')
        11     coherence_lda = coherence_model_lda.get_coherence()
        12     print(k, coherence_lda)
        13     scores_nmf.append(coherence_lda)
```

```
3 0.28168517274532157
4 0.44236238790411175
5 0.4411030317965917
6 0.4897078960989399
7 0.43361447763211913
8 0.45412241563139766
9 0.4705144598195836
10 0.45781569456240473
11 0.45155744281011023
12 0.40227751202378337
13 0.4254008453564542
14 0.40692994154704376
```

In [13]:

```
1 plt.figure(figsize=(10, 5))
2 plt.plot(list(range(3,15)), scores_nmf, marker='o', color='blue')
3 sns.despine(top=True, right=True, left=False, bottom=False)
4
5 plt.locator_params(integer=True)
6 plt.title('Coherence score vs the number of topics for NMF')
7 plt.xlabel('Number of topics')
8 plt.ylabel('Coherence Scores')
9 plt.savefig('lda_scores.jpg')
10 plt.show()
```




```
In [14]: 1 selected_topics_nmf=np.argmax(scores_nmf)+3
2 nmf_model = Nmf(corpus=corpus_tfidf, id2word=dct, num_topics=selected_topics_nmf,
3               random_state=12, chunksize=128, passes=10 )
4
5 nmf_model.print_topics()
```

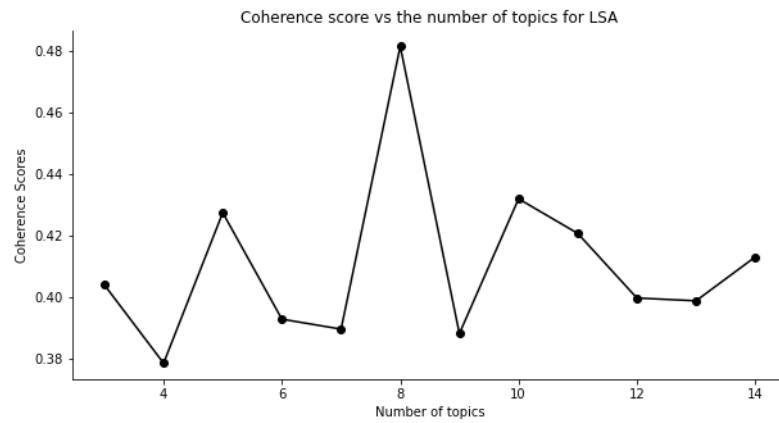
```
Out[14]: [(0,
'0.285*"poupon" + 0.285*"grey" + 0.247*"like" + 0.142*"gettin" + 0.002*"sound" + 0.001*"transformer" + 0.001*"name" + 0.001*"somewhat" + 0.001*"movie" + 0.001*"look"'),
(1,
'0.123*"case" + 0.044*"first" + 0.041*"two" + 0.030*"detected" + 0.025*"variant" + 0.019*"report" + 0.015*"canada" + 0.013*"york" + 0.013*"new" + 0.012*"confirmed"'),
(2,
'0.023*"variant" + 0.017*"new" + 0.010*"travel" + 0.010*"world" + 0.010*"omicronvariant" + 0.010*"coronavirus" + 0.009*"country" + 0.009*"vaccine" + 0.008*"know" + 0.007*"say"'),
(3,
'0.153*"cold" + 0.152*"common" + 0.144*"symptom" + 0.077*"routine" + 0.077*"majority" + 0.076*"similar" + 0.076*"present" + 0.073*"testing" + 0.073*"mild" + 0.066*"people"'),
(4,
'0.114*"covid" + 0.101*"analytics" + 0.051*"usafacts" + 0.051*"datavisualization" + 0.051*"datascience" + 0.050*"globalhealth" + 0.050*"growth" + 0.050*"publichealth" + 0.050*"insight" + 0.050*"county"'),
(5,
'0.297*"moderna" + 0.153*"candidate" + 0.153*"address" + 0.138*"booster" + 0.137*"via" + 0.077*"variant" + 0.003*"say" + 0.001*"pfizer" + 0.001*"vaccine" + 0.001*"patient"')]
```

Latent Samentic Analysis

```
In [15]: 1 from gensim.models.lsimodel import LsiModel
2
3 scores_lsi = []
4 for k in range(3,15):
5     # LSI model
6     lsi_model = LsiModel( corpus=corpus_tfidf, num_topics=k, power_iters=250,
7                       id2word=dct)
8     # to calculate score for coherence
9     coherence_model_lsi = CoherenceModel(model=lsi_model, texts=sents, dictionary=dct, coherence='c_v')
10    coherence_lsi = coherence_model_lsi.get_coherence()
11    print(k, coherence_lsi)
12    scores_lsi.append(coherence_lsi)
```

```
3 0.4039538049985449
4 0.3784976981923267
5 0.42744837035054306
6 0.3928130838248274
7 0.3895709231074262
8 0.48158829107852097
9 0.3881411519633174
10 0.43196943563017276
11 0.4207294653554794
12 0.3997005384465795
13 0.3987879310335635
14 0.41305537942116277
```

```
In [16]: 1 plt.figure(figsize=(10, 5))
2         plt.plot(list(range(3,15)), scores_lsi, marker='o', color='black')
3         sns.despine(top=True, right=True, left=False, bottom=False)
4
5         plt.locator_params(integer=True)
6         plt.title('Coherence score vs the number of topics for LSA')
7         plt.xlabel('Number of topics')
8         plt.ylabel('Coherence Scores')
9         plt.savefig('lda_scores.jpg')
10        plt.show()
```



```
In [17]: 1 selected_topics_lsi = np.nanargmax(scores_lsi)+3
2 lsi_model = LsiModel( corpus=corpus_tfidf, num_topics=selected_topics_lsi,
3 id2word=dct)
4 lsi_model.print_topics()
```

```
Out[17]: [(0,
'-0.577*"poupon" + -0.577*"grey" + -0.500*"like" + -0.289*"gettin" + -0.011*"cold" + -0.011*"common" + -0.010*"symptom" + -0.006*"covid" + -0.005*"routine" + -0.005*"majority"',
(1,
'-0.469*"cold" + -0.468*"common" + -0.444*"symptom" + -0.236*"routine" + -0.235*"majority" + -0.235*"similar" + -0.232*"present" + -0.225*"testing" + -0.224*"mild" + -0.204*"people"',
(2,
'-0.473*"covid" + -0.418*"analytics" + -0.209*"usafacts" + -0.209*"datavisualization" + -0.209*"datascience" + -0.209*"globalhealth" + -0.208*"growth" + -0.208*"publichealth" + -0.208*"insight" + -0.207*"county"',
(3,
'0.700*"moderna" + 0.360*"candidate" + 0.360*"address" + 0.327*"booster" + 0.325*"via" + 0.194*"variant" + 0.008*"case" + 0.007*"new" + 0.006*"say" + 0.006*"vaccine"',
(4,
'-0.492*"case" + -0.398*"variant" + -0.307*"new" + -0.198*"first" + -0.170*"coronavirus" + -0.167*"two" + -0.165*"detected" + -0.146*"omicronvariant" + -0.136*"say" + -0.123*"travel"',
(5,
'-0.679*"eth" + -0.362*"think" + -0.360*"mask" + -0.349*"luck" + -0.335*"dont" + 0.122*"case" + -0.081*"omicronvariant" + -0.047*"know" + 0.043*"first" + 0.040*"two"',
(6,
'-0.543*"omicronvariant" + 0.494*"case" + -0.396*"cake" + 0.156*"eth" + 0.147*"first" + 0.146*"two" + -0.133*"vaccine" + -0.109*"know" + -0.095*"travel" + -0.092*"new"',
(7,
'-0.537*"omicronvariant" + -0.445*"cake" + -0.379*"case" + 0.207*"new" + 0.191*"variant" + 0.154*"travel" + 0.143*"know" + -0.140*"first" + 0.129*"vaccine" + -0.119*"two"')]
```

Best Model Analysis

On the basis of coherence scores lda model have been chosen as the best model for the current data set

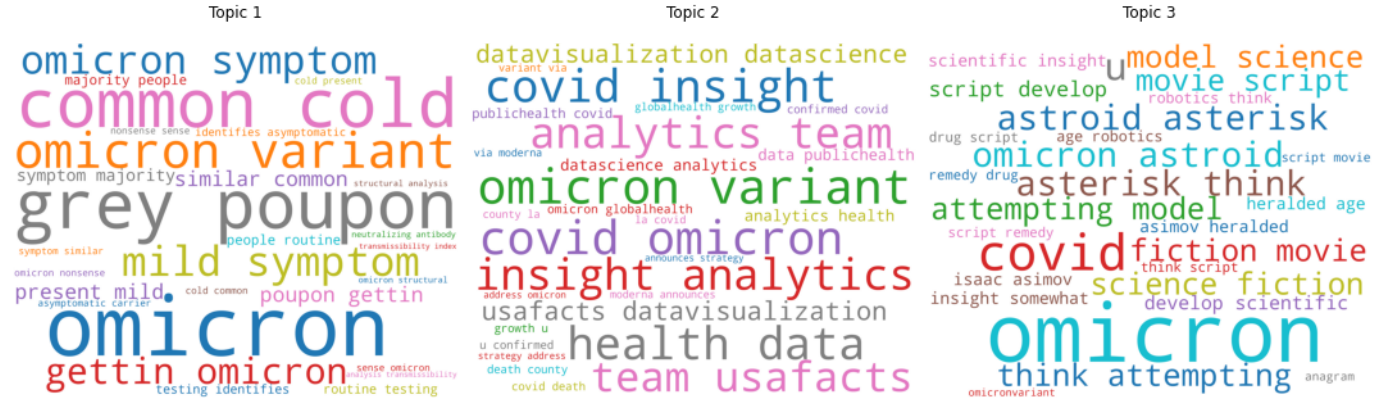
```
In [18]: 1 # judged manually from pyldavis
2 topics_name = ['Topic 1', 'Topic 2', 'Topic 3']
3
4 # getting the most dominant topics from a trained model
5 predicted_topics = lda_model[corpus_tfidf]
6
7 probs, topics = [], []
8 for k in predicted_topics:
9     # sorting the probabilities
10    k.sort(key=lambda x:x[1])
11    # selecting the topic with greatest probability
12    topics.append(topics_name[ k[0][0] ] )
13
14 data['Topics'] = topics
15 data.head(2)
```

```
Out[18]:
```

	tweet_id	date	text	user_name	user_location	user_description	user_created	user_followers	user_friends	user_favourites	hashtags	source	is_retweet	processed	Lengths	Topics
0	1465162850457071621	2021-11-29 03:36:59+00:00	Will Boris Johnson ever learn his lesson after...	James robertson	NaN	NaN	2013-04-22 19:41:30+00:00	303.0	188	84059	NaN	Twitter for Android	False	boris johnson ever learn lesson endless covid ...	10	Topic 1
1	1465162841665769474	2021-11-29 03:36:57+00:00	#Omicron so hot; it's being rush shipped in De...	Gene Bulmer	East Coast	#ChristFollower, #Leader, Husband, Father, #Ru...	2014-11-10 22:43:27+00:00	99.0	162	745	['Omicron']	Twitter for iPhone	False	omicron rush shipped package symptom barely kn...	9	Topic 1

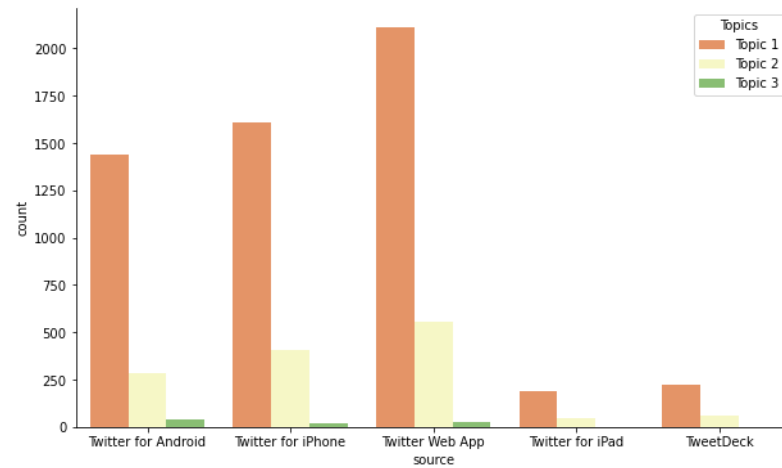
```
In [19]: 1 # 1. Wordcloud of Top N words in each topic
2
3 from matplotlib import pyplot as plt
4 from wordcloud import WordCloud
5
6 cloud = WordCloud(background_color='white', width=2500, height=2000,
7                   max_words=28, colormap='tab10', prefer_horizontal=1.0)
8
9 fig, axes = plt.subplots(1, 3, figsize=(15,8), sharex=True, sharey=True)
10
11 for i, ax in enumerate(axes.flatten()):
12
13     fig.add_subplot(ax)
14     if i>len(topics_name)-1:
15         continue
16     curr = data[data['Topics']==topics_name[i]]
17     print(curr.shape)
18     tokens = [tok for d in curr['processed'] for tok in d.split(' ')]
19     cloud.generate(' '.join( tokens ))
20
21     plt.gca().imshow(cloud)
22     plt.gca().set_title( topics_name[i]+'\\n')
23     plt.gca().axis('off')
24
25 plt.axis('off')
26 plt.tight_layout()
27 plt.show()
```

(6233, 16)
(1738, 16)
(95, 16)



In [20]:

```
1 def subcategory_plot(df, col):
2
3     plt.figure(figsize=(10,6))
4     bar_pub = list( df[col].value_counts().index[:5] )
5     temp2 = df[df[col].isin(bar_pub)]
6
7     sns.countplot(x=col, hue='Topics', data=temp2, palette="RdYlGn")
8     sns.despine()
9
10    plt.savefig('bars_{}.png'.format(col))
11    plt.show()
12
13    subcategory_plot(data, 'source')
```



In [21]:

```
1 # users have been categorized according to their number of followers
2 data['User type']=pd.cut(data.user_followers, [1,200, 1000, 2000, 10000, 100000000],
3                           labels=['Naive', 'Average', 'Popular', 'Micro-Influencer', 'Influencer'])
4
5 subcategory_plot(data, 'User type')
```

