

# Topic Modelling of Tweets on Omicron

Venkata Shiva Sai Mallikarjun Devasani  
*Electrical and Computer Engineering*  
 University of Florida  
 Gainesville, FL, USA  
 vdevasani@ufl.edu

**Abstract** – In this project, the topic modeling analysis of Omicron's Twitter dataset taken from Kaggle is explained. One of the Covid-19 variants is Omicron. The dataset contains user tweets on this topic acquired from Kaggle. Three topic modeling techniques will be used to investigate the underlying distribution of topics across the tweets. Latent Dirichlet Allocation (LDA), Latent Semantic Analysis (LSA), and Non-negative Matrix Factorization (NMF) are those algorithms. To begin, preprocessing of the data and exploratory analysis will be performed to better comprehend the tweets. After preprocessing of the tweets data, they are inputted to the algorithms, which generates the C-V coherence plots. The score from these plots will be compared amongst all the three techniques to evaluate their performance and also to arrive at the optimum number of topics for the respective technique. Also, various visualizations were used to analyze the results.

**Index Terms**— LDA, NMF, LSA, Omicron variant, Tweets, Kaggle, Topic modelling.

## I. INTRODUCTION

THE B-1-1-529 Covid-19 variant is popularly known by the name, "the Omicron variant." Initially, on the 24<sup>th</sup> of November '21, it was first reported in South Africa. World Health Organization (WHO) has declared the Omicron to be a variant of concern on the 26<sup>th</sup> of November '21. The organization also declared that they are not completely certain about the severity or the transmissibility of this variant. Also, they were not sure whether the existing vaccinations and safety precautionary countermeasures would be sufficient to fight with the Omicron. Preliminary research showed that it is a variant posing much greater reinfection risk than other previously reported ones. This indicates that persons who have already been infected with COVID-19 may get readily reinfected, which has caused anxiety among people all around the world. The utilization of virtual entertainment stage media to examine individuals' impression of infections has shown to be essentially valuable previously. Illnesses like H1N1 flu, Zika infection and Ebola infection attracted public discussions to figure out worldwide opinions on sicknesses. This is especially entirely important on account of the Omicron infection which is generally new as such data can be valuable. I have chosen to apply unaided subject displaying methods on tweets information on clients about this variation and comprehend the main themes examined.

Topic Modelling approaches have been utilized in many distributed articles on various points, for example, in sciences, programming, interpersonal organizations, etc. (Daud et al., 2010, p. 280-301) utilized point demonstrating with delicate grouping to

characterize existing models in numerous classifications. Boundary assessment similar to Gibbs inspection and measures of execution assessment were utilized. Chen al., 2016 played out an overview on subject displaying with an emphasis on computer programming field. Top modelling was the subject of discussion across various articles from the December of 1999 to December of 2014.

From the inception of Covid-19 virus outbreak, numerous researchers have used social media data to acquire a deeper understanding of the disease. Here, Sentimental Analysis is the examination of the behavior of people towards the virus and handling the virus, their ideology regards the wearing of the facemasks. Heet al, 2021 is a researcher whose research is based on this sentimental analysis and filtering of the specific keywords. Ordun al, 2020 is another researcher whose work is based on one of the most popular topic modelling algorithms, LDA. He associated the tweets during the outbreak with certain events occurring in the world using this algorithm. Re-current neural networks in conjunction with the LDA algorithm is used by Jelodar al., 2020 to categorize the tweets into distinct topics and then perform sentimental analysis for each topic. Mackey al, 2020 employed Bi-term Topic Model for Covid19 tweets data by splitting them into topics consisting of certain words like topics clustering around symptoms, testing, and recovery.

## DATA MODIFICATION

### *Data Description*

The data collection for this project has been through the already available tweets dataset on Kaggle, as it is an online repository that has public datasets for free. The data collection comprises tweets regarding the new Covid19 version – Omicron – from people all around the world. As features or characteristics, the dataset includes all its critical information. The dataset is csv file with each tweet as rows and various information attributes regarding the tweet like the id of the tweet, location of the user, context of the tweet, username, date, hashtags, etc. All of this information is contributing to a total of 8073 tweet rows and 13 attribute columns for each row.

### *Preprocessing the data*

As we are dealing with the textual data here, preprocessing it before inputting it to the algorithms is very important. The textual data in its original form might be untidy. The preprocessing operations are carried out to assist in the cleaning and transformation of data into formats that may be better understood, hence assisting in the achievement of better results. Data preparation not only removes noise but also lowers data, which speeds up calculations.

tweet_id	date	text	user_name	user_location	user_description	user_created	followers	retweets	hashtags	source	is_reply
1.476-17	2021-11-21	Will Boris J James robertson				2013-04-22 1	303	188	84059	Twitter fo	FALSE
1.476-18	2021-11-21	Gene Bulmer East Coast			#AmericaForever, #2014-11-18		99	162	745	['Omiconr'] Twitter fo	FALSE
1.476-19	2021-11-21	@Joelbider Andrew Arie Galendes @3RfAI			AMP PURE	2014-01-10 2	114	3223	4830	Twitter fo	FALSE
1.476-20	2021-11-21	Gold Coast myGc.com.au Gold Coast, At The Gold Coast's			2009-03-30	12232	3191	271	['GoldCoastZapier.com']	Twitter fo	FALSE
1.476-21	2021-11-21	What is Ghost Of G I Poughkeepsie Freelancer who			2011-11-17 0	603	960	819	Twitter fo	FALSE	FALSE
1.476-22	2021-11-21	@Jonathan beccary			Love my God, Los 2020-11-13 2	114	461	174	Twitter fo	FALSE	FALSE
1.476-23	2021-11-21	Of recoup Calagra Herald Calagra, India			Latest local, natic 2008-10-16 2	201011	802	659	Echobox	FALSE	FALSE
1.476-24	2021-11-21	Indy today			Brings you news	2014-01-01 2	8664435	2448	['Canada', 'indytoday']	Twitter fo	FALSE
1.476-25	2021-11-21	torianakm Toronto, On			Imam Lutfullah PMSA 2015-06-23 1	51	612	4071	Twitter fo	FALSE	FALSE
1.476-26	2021-11-21	Johnson, M Bromio			city Manjy dave from br livi 2015-01-31	32	41	1162	['Omiconr'] Twitter fo	FALSE	FALSE
1.476-27	2021-11-21	@Attilio Corti Italia			Cardiovascular su 2012-01-29 1	363	1429	7368	['Omiconr'] Twitter fo	FALSE	FALSE
1.476-28	2021-11-21	COVID Sprakly Bright The blue birds Paddling			2021-10-03 0	58	403	1033	Twitter W	FALSE	FALSE
1.476-29	2021-11-21	@woopden11111111111111111111 Australia			2012-01-02 1	128	505	35151	Twitter fo	FALSE	FALSE
1.476-30	2021-11-21	While Apple or e=ed Japan			Japanesew=0Y 2019-09-21 2	2	36	68	Twitter W	FALSE	FALSE
1.476-31	2021-11-21	Omiconr Michael Hoff Pannsylvania			Ad.a d.huby, A 2008-10-32 1	198	301	7106	Twitter fo	FALSE	FALSE
1.476-32	2021-11-21	With the sTMX Capital			Blockchain Techn 2021-01-22 1	1037	7	2	Twitter W	FALSE	FALSE
1.476-33	2021-11-21	@Caraf6 The Plan			UK Somewhere James 1:12 Bless 2021-01-23 1	447	451	30573	Twitter fo	FALSE	FALSE
1.476-34	2021-11-21	Could the Matthew Kim Brisbane, Que Professor of IP & 2009-06-26 1			17072	18751	60557	Twitter fo	FALSE	FALSE	
1.476-35	2021-11-21	WFO SF Web Writ			Los Angeles, C 2015-06-23 1	17	17	478	Twitter W	FALSE	FALSE
1.476-36	2021-11-21	WFO The WTs Jo Brisbane, Que G O O 2015-10-14 1			56	281	976	Twitter W	FALSE	FALSE	
1.476-37	2021-11-21	In other ne Paul Reimer			Scarsfield MB 2012-11-01 1	99	65	1096	Twitter W	FALSE	FALSE
1.476-38	2021-11-21	On Canada's Sri Rathod			Bengaluru, Ind Keep Smiling, I 2021-08-22 20	20	1	7	Microsoft	FALSE	FALSE
1.476-39	2021-11-21	Canada's SiRiShta			Born Maritime host of My life 2011-02-07 2	20	185	6683	Twitter W	FALSE	FALSE
1.476-40	2021-11-21	@drjanyan Koemr Cornwall England			2017-03-01 1	78	272	1024	Twitter W	FALSE	FALSE
1.476-41	2021-11-21	Why is Post@G J Amer Woman Co Co Post-America w			2019-09-06 1	7597	287	174556	Twitter fo	FALSE	FALSE

Figure 1: Tweets data from Kaggle

Filtering/removing of punctuation, numbers, Stopwords pulled in from the NLTK library which does the natural language processing, unnecessary spaces, tokens that consist of only single characters, mathematical signs, links, hashtags, hyphens, underscores, slashes, underscores, periods, and index termed data is being lowercased. Lastly, using the WordNet lemmatize function, words were reduced to their most basic form.

### Analysis of tweet data

From the preprocessed data corpus from the previous stage, we are now going to generate a word cloud consisting most of the frequently repeated words so that we can have a deeper insight over the tweets about the Omicron variant made by Twitter users. The analysis's outcome is displayed below. The words that appeared the most frequently have a higher font size here. The term omicron, as predicted, dominates the word cloud, which appears in multiple bigrams.

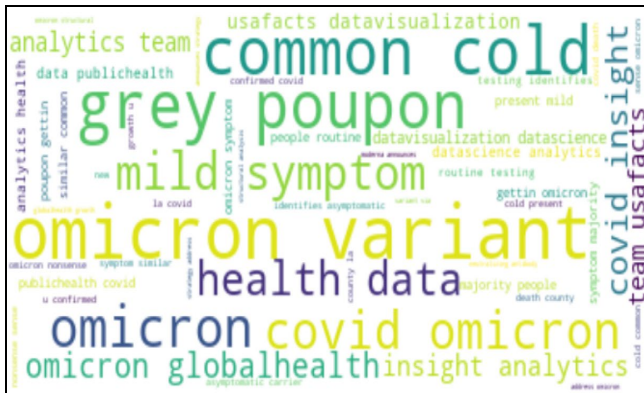


Figure 2 : Frequent words in the corpus

Another significant aspect to consider is the distribution of tweet durations in the dataset, which aids in understanding the length of user discussions. Figure 3 depicts the distributions, with the bulk of tweets ranging in length from 5 to 25 characters. Given Twitter's 280-character restriction, it was expected to fall somewhere in this region.

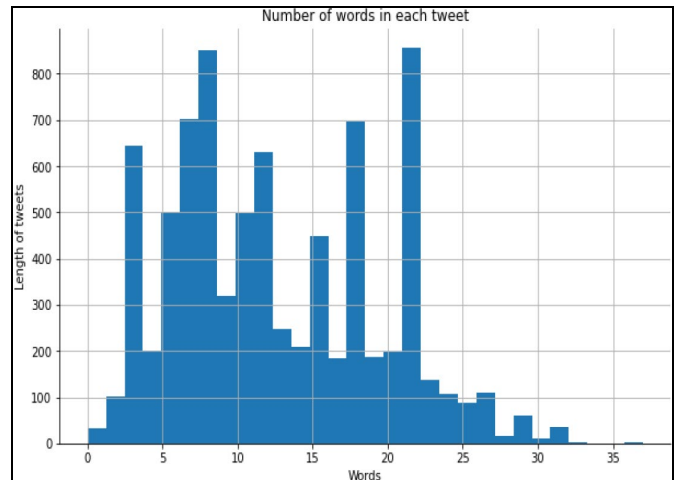


Figure 3: Tweet length vs words

### Selection of features and vectorization

Right after preprocessing the data, we face another challenge in the form of unwanted words throughout the data. Each tweet has lot of the unwanted context due to conjunctions or grammar for forming or bonding the sentence. All these unwanted contexts need to be eliminated using the feature selection. This helps us with having only significant terms in our data that is going to be inputted to the algorithm. Also, it massively speeds up the training process. If a word is appearing at more than 3 instances and less than 60% of all the tweets as they might be domain-specific, then that particular term is deemed significant. Another filtration that needed to be done is for vocabulary repetitions. For this selection, Bigrams are being used. Initially we had 9834 unique token items which were reduced to 3263 after applying the feature selection. This data is then taken as an input for LDA, I constructed a dictionary and corpus.

Because computers cannot operate directly with words, they must first be converted into some numerical form. The corpus's TF-IDF scores were computed, and the data was now ready to be input into the topic modeling algorithms.

## II. DESCRIPTION

To examine the dataset, I used three unsupervised topic modelling strategies. Topic modelling, also known as a probabilistic clustering technique, is a machine learning approach for categorizing massive collections of text corpora into abstract topics and themes.

### Latent Dirichlet Allocation

LDA is one of the most popular Topic Modelling Algorithm which is widely implemented across various fields as it has extensively well-suited libraries which can make its implementation very simple. LDA leverages the Bayesian so that it could spot hidden or latent topics in the corpus of data. Using combination of different topics and vector space representation of words to count the terms a Probabilistic Model is generated. The TD-IDF matrix is used to represent the text corpus which requires the number of topics to be specified in prior as this is an unsupervised modeling.

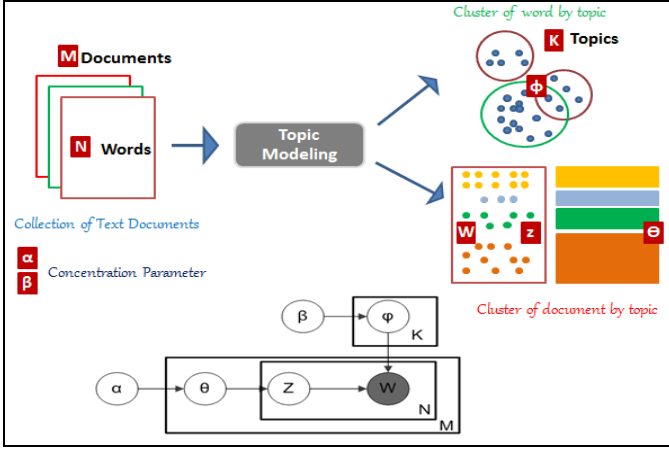


Figure 4: LDA Algorithm working model

### Latent Semantic Analysis (LSA)

LSA uses the LSI function and is formerly known to be called Latent Semantic Indexing is a Natural Language Processing (NLP) approach for creating semantic content by vectorizing the text corpus. LSA identifies a group of words and concepts that are hidden of the property to frequently appear together, and which are expected to be relatively linked using the term co-occurrence. For this process, Singular-Value Decomposition is used where the whole term-document matrix is broken down to learn the hidden groups of topics by matrix decomposition. Thereby, giving us the scope to identify the patterns, concepts, and the relationships amongst the words in the corpus. Initially, the documents existing in the term-document matrix are converted as frequency matrices. Followed by the reduction of high-dimensional text documents vector-space representations into lesser dimensional representations even may be to a fixed figure of dimensions.

#### 2.3 Non-Negative Matrix Factorization (NMF)

The NMF method is a linear algebra method for extracting important data information. Even though the NMF does not have any preceding knowledge or knows the data context, it is capable of extracting vital information from the input data matrix. Unlike the LSA or LDA, where the TD-IDF matrix is directly inputted to the algorithm, here, this input TD-IDF term-document matrix is further decomposed into topic matrix and topic-term matrix..

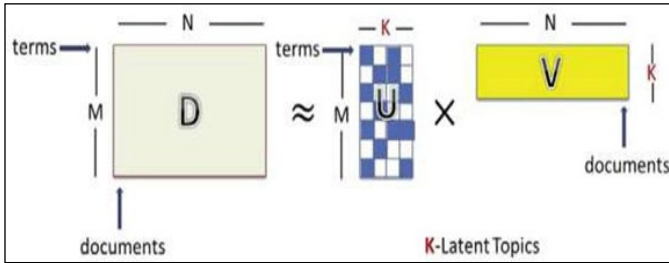


Figure 5: NMF Algorithm working model

### Libraries used

NumPy, Pillow, Tabulate, Pandas, wordcloud, nltk, seaborn, gensim, pyLDAvis, matplotlib.

## III. EVALUATION

### Latent Dirichlet Allocation (LDA)

The Gensim package was used to create the LDA model, and the corpus was provided as an input in the form of its TF-IDF. It's critical to choose the right number of topics in unsupervised modeling while keeping the underlying topic distribution in mind. We are going to use the Coherence scores to determine a topic's quality. C-V, a metric that uses sliding cosine similarity, windows, and Normalized Pointwise Mutual Information (NPMI) to measure the quality of subjects, was chosen for this purpose.

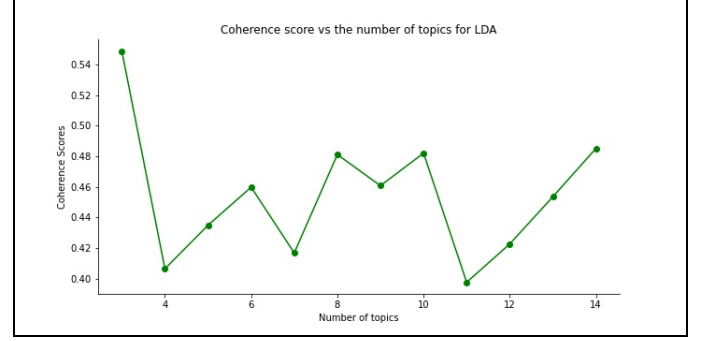


Figure 6: LDA coherence score

The coherence score was calculated on a scale of three to fifteen topics. Figure 6 shows the findings, as the topics are three, we find the coherence score to reach its maximum for this algorithm for the whole dataset.

### Latent Semantic Analysis (LSA)

The LSA's main premise is that words will be semantically related as soon as they appear in alike texts. Here, we again use the exact same TF-IDF corpus as that of the one used for the LDA. The LSI function is then supplied with this, and the genism is employed with finding the data clusters. The best number of topics for this method was determined by measuring coherence scores from 3 to 15. Figure 7 depicts the coherence scores in relation to the number of subjects.

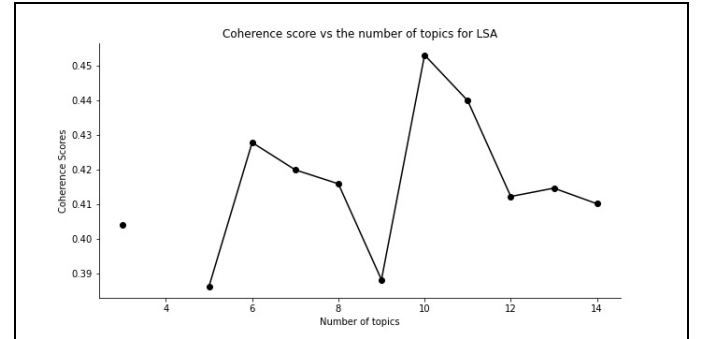


Figure 7: LSA coherence score

When the number of subjects was 10 and the coherence score was 0.45, the highest score was reached. In contrast to LDA, the coherence score here increased as the number of topics climbed until it reached ten, and then it fell.

### Non-Negative Matrix Factorization (NMF)

The same corpus that was utilized in the preceding section is used for Gensim. The best number of topics was determined using c v coherence ratings, which were recorded from topic 3 to 15. The resulting outcomes are depicted in Figure 8.



## VI. REFERENCES

- [1] Chew, C., & Eysenbach, G. (2010). Pandemics in the age of Twitter: content analysis of tweets during the 2009 H1N1 outbreak. PloS one. (URL: <https://journals.plos.org/plosone/article?id=10.1371%2Fjournal.pone.0014118>)
- [2] Stefanidis, A., Vraga, E., Lamprianidis, G., Radzikowski, J., Delamater, P. L., & Jacobsen, K. H. (2017a). Zika in Twitter: Temporal Variations of Locations, Actors, and Concepts JMIR Public Health and Surveillance (URL: <https://publichealth.jmir.org/2017/2/e22/>)
- [3] Odium, M., & Yoon, S. (2015a). What can we learn about the Ebola outbreak from tweets? *American Journal of Infection Control*. (URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4591071/>)
- [4] Abd-Alrazaq, A., Alhuwail, D., Househ, M., Hamdi, M. & Shah Z. (2020) Top concerns of tweeters during the COVID-19 pandemic. *Infoveillance study, J. Med.Internet Res.* 22 (4) (URL: <http://dx.doi.org/10.2196/19016>)
- [5] Mackey, V. Purushothaman, J. Li, N. Shah, M. Nali, C. Bardier, B.Liang, M. Cai, R. Cuomo (2020). Machine learning to detect self-reporting of symptoms, testing access, and recovery associated with COVID-19 on Twitter: Retrospective big data invigilance study, JMIR Publ. Health Surveillance 6 (2) (URL: <http://dx.doi.org/10.2196/19509>)
- [6] Vayansky, I. & Kumar, A. P. S. (2020) A review of topic modeling methods, *Information Systems.* 94. (URL: <https://doi.org/10.1016/j.is.2020.101582>)
- [7] <https://www.sciencedirect.com/science/article/pii/S0306>