



**MASTER OF SCIENCE (MSc)**  
**AGRICULTURAL & FOOD DATA MANAGEMENT**  
**Institut Polytechnique UniLaSalle – Campus de Rouen. France**

Master's Thesis

"To build an interactive NLP dashboard tool to understand prevalent user perspectives around the concept of #organic as inferred from Instagram posts. The tool allows us to explore word frequency, posts frequency by unique word, word co-occurrence counts, and word embeddings, which we use to deep dive into the concept of 'organic food' as a test case."

**Student:** Ponnappa Appanna Machimanda, MSc AgFood Data management 2022

**Supervisors:**

Salima Taibi – Program supervisor, MSc AgFood Data management  
Alessandra Palmigiano – Full professor, VU Amsterdam

## Acknowledgements

I would like to express my heartfelt gratitude to my supervisors at the host institution, Dr Alessandra PALMIGIANO, and at the sending institution, Dr. Salima TAIBI, for giving me the opportunity to work on this research project. I have learned so much during this time which I believe will hold me in good stead throughout my future career.

Alessandra was gracious with her knowledge and expertise, providing me a direction to structure my research around. The discussions with her often opened my mind to new ways of looking at the data and capturing the intuition behind naturally occurring effects like the self-reinforcement effect, and the vector space semantics of lexical terms in language. She made me feel comfortable exploring unknown territory, which is common in a research project, while giving me a freedom to learn from my mistakes and deciding a suitable course of action at important junctures. I would also like to thank Nachoem Wijnberg for providing me an exposure to relevant projects in the field of NLP and offering more context into research in this field.

Now let me briefly describe my masters' program journey. After finding a passion for data analysis and introspecting into what domain I would like to work in, I realized that the idea of working with nature really appealed to me. I joined the masters' program at Unilasalle, Rouen, crafted by Salima Taibi and her team, which I found to have a unique and future relevant mix of data modules for agriculture. I am immensely grateful to the institution Unilasalle-Rouen, and Salima for accepting me as a student in this program.

Through the course of this masters' program, I was mentored by several notable professors who bolstered my understanding about statistical and programming concepts. I have mentioned a few people here but please note that there were also other professors and support staff in the program who significantly shaped my masters experience. Due to brevity, I do not mention them here. Needless to say, I whole heartedly appreciate the thought put into creating this masters' program and all the wonderful lecturers whose passion and expertise for their subject was shared with us students, eventually enabling a competence to carry out a project like this.

I would like to thank Jerome DANTAN, for providing guidance and coaching on the first python programs, and also being a responsible supervisor for our cohort. In all matters that involved administrative formalities during my stay in France, I would like to thank Emmanuelle GH for her tireless support despite having to attend to several international students. Rural sociology taught by Bill sparked interesting debates and a push to understanding perspectives of the different players in the Agricultural industry. The dashboard tool created in this thesis offers a window into the perspectives of the people and businesses as inferred from social media posts.

Along with an introduction to R programming, the foundational concepts about NLP (Natural Language Processing) methods using python, were taught by Redha MOULLA. It was indeed these notes and lectures on NLP that I first referred in the building of the tool.

I had the privilege of working with many (perhaps the most!) team-mates for the different projects in the program, all of whom have contributed to developing my skills in this field.

Hasnaa Latique, Uranie Jean-Louis, Robin Miranda, Elise BJ, Mouad Benkhouya, Yassine Ben-Mohamed, Marilena Charalambous, Fagr Gissan, and Malika Sekour. You guys were amazing to work with and I learned a lot from each of your unique strengths and personalities!

Last but not least, I would like to express gratitude to members of my family. I would like to thank my sister, Ashwini Ponnappa, and my parents, for encouraging me to pursue this masters' and supporting me throughout the program. I am deeply grateful and fortunate. During my time in the Netherlands, I was privileged to stay and explore the country with my cousin, Poonacha Medappa, who also gave me an immersive experience of Dutch culture. Much gratitude!

## Table of Contents

<b>Acknowledgements .....</b>	<b>2</b>
<b>Chapter 1. Introduction.....</b>	<b>6</b>
<b>Chapter 2: Literature review .....</b>	<b>7</b>
<b>2.1 LSA.....</b>	<b>8</b>
<b>2.2 Glove.....</b>	<b>9</b>
<b>2.3 Word2Vec .....</b>	<b>10</b>
<b>2.4 User types .....</b>	<b>10</b>
<b>3: Methodology .....</b>	<b>12</b>
<b>3.1 Data scraping:.....</b>	<b>12</b>
3.1.1 Selenium:.....	13
3.1.2 Steviesie data platform: .....	13
3.1.3 Phantombuster data platform: .....	13
<b>3.2 Data preprocessing: .....</b>	<b>16</b>
3.2.1 Dataframe_1.....	16
3.2.2 Dataframe_2.....	16
3.2.3 Dataframe_3.....	16
3.2.4 Language detection .....	17
3.2.5 Further filtering .....	18
3.2.6 Stemming and Lemmatization .....	18
3.2.7 Word Vocab.....	19
<b>3.3 Word Embeddings .....</b>	<b>19</b>
<b>3.3.1 DSM .....</b>	<b>19</b>
3.3.1.1 Co-occurrence matrix.....	19
3.3.1.2 SVD matrix .....	20
3.3.1.3 Plotting .....	20
<b>3.3.2 Word2Vec model .....</b>	<b>21</b>
3.3.2.1 Instantiation .....	22
3.3.2.2 TSNE transformation .....	22
3.3.2.3 Plotting .....	22
<b>4: Results and discussion.....</b>	<b>23</b>
<b>4.1 The top co-occurrences with the word “organic” .....</b>	<b>24</b>
<b>4.2 Deep analysis around the concept of “organic food”.....</b>	<b>26</b>
4.2.1 Instagram DataTable (DT-1) .....	26
4.2.2 Wordcloud (WC-1) – Word frequency .....	27
4.2.3 Z-score DataTable (DT-2).....	28
4.2.4 Wordcloud (WC-2) – Post frequency .....	28
4.2.5 Bar chart (BC-1) – Posts count.....	29
4.2.6 Bar chart (BC-2) – Co-occurrences count.....	29
4.2.7 Vector Exploration (VE) .....	30
4.2.8 Vector Math (VM) .....	31
4.2.9 Synthesis.....	32
<b>5: Conclusion .....</b>	<b>34</b>
<b>6: Bibliography .....</b>	<b>36</b>

<b>7: Appendix .....</b>	<b>39</b>
<b>7.1 Interactive Dashboard Tool .....</b>	<b>39</b>
<b>7.2 Development of code.....</b>	<b>40</b>
<b>7.3 Python Code – preprocessing and creating support files.....</b>	<b>41</b>
<b>7.4 Python Code – Plotly dash web app.....</b>	<b>52</b>
<b><i>Summary – English and French.....</i></b>	<b>66</b>

## Chapter 1. Introduction

Social media has made it possible for people to freely express their opinions and show their support or apprehension about certain topics. Groups of people that have similar opinions about a subject, tend to express themselves in similar ways using similar words. In this report we utilize social media data to understand how different groups of people use the word ‘organic’ in their relevant vocations. A dashboard tool is created for this purpose, and the entire code has been made available open source on the github repo: [https://github.com/Appymp/Agridata\\_mining](https://github.com/Appymp/Agridata_mining). The code has been developed in a block-wise, modular fashion and can be easily adapted for use in future projects.

This is not simply a positive or negative sentiment analysis but rather an approach to intuit a broader meaning. For example, as we will find out later in the results section, “organic” assumes certain connotations in the context of ‘cosmetics’ and a different one in the context of ‘food’. The groups of people creating this definition are the ones posting on social media.

In this study we have used data specifically from Instagram for posts where the user has added a hashtag ‘#organic’. So it is implied that the post has something to do with organic but not necessarily based on the same literal meaning. The text description entered by the user contains valuable information about what the user thinks “organic” is about, or what it is being portrayed as ‘organic’ by the user.

Broadly there are 2 categories of people that post content on Instagram:

- users that create posts for business purposes,
- users that do not have a business interest but more of passionate association.

We do not explore classification of the users in depth, but we do create an option for detecting a possibly influential post through a likes per minute count.

The dashboard tool we have created is geared towards an analyst who not only relies on the standard graph analysis approach of making sense from the data (usually sequentially), but who can also use the multiple graphs in a back-and-forth manner, to confirm and refute the assumptions being made. There is expected to be a learning curve where one can develop their own approach of using the tool as well. Subject matter and marketing experts in the domain can use this as a discussion and exploration tool to draw more general and intuitive inferences about their target market and audience.

The tool consists of word clouds for studying the frequency of words and frequency of posts-by-word, bar charts for co-occurrences and post counts, and word2vec based vector space graphs for checking “similarity” of a particular word through the context expressed in the posts. The data rendered in these graphs can be dynamically selected from the data table filtering and graph options:

- row selection and single or multiple column combinations of description, hashtags, usernames, caption mentions, weblinks
- word filtering of multiple columns

- influential post filtering with z-score slider from 0 to 4. Higher value indicates a higher likes count in a unit of time.
- zoom-in and out of the graphs and saving the images to file
- Sliders for setting the number of results to display in the bar charts and vector space graphs.
- Input boxes for words displayed as vectors and to check co-occurrences

## Chapter 2: Literature review

In order to obtain the meaning of the words we will create word vectors from the descriptions of the Instagram posts. Each word in the description is represented as a word embedding. Word embeddings are fundamentally a form of word representation that links the human understanding of knowledge meaningfully to the understanding of a machine. The representations can be a set of real numbers (a vector) (Selva Birunda & Kanniga Devi, 2021).

These vectors can be created using different methods ranging from the more transparent matrix factorization methods, to the more mysterious and better performing neural network models. In the neural network models like Word2vec, the actual algorithm is designed to locally predict the surrounding words and as a byproduct of this task, word embeddings are created which capture a lot of semantic information.

We explore some of these methods and the feasibility of implementing them in our use case. As pointed out by (Levy et al., 2015, p. 211) performance gains of word embeddings are due to certain system design choices and hyperparameter optimizations, rather than the embedding algorithms themselves. The embeddings hinge upon the distributional hypothesis by (Harris, 1954), where the degree of semantic similarity between two words (or other linguistic units) can be modeled as a function of the degree of overlap among their linguistic contexts.

The format of distributional representations vary depending on the specific aspects of meaning they are designed to model (Baroni & Lenci, 2010). At the core of these are the word co-occurrence counts. It is useful to identify these representations on the basis of what kind of similarity they are modelling. The concept of *relational similarity* and *attributional similarity* are significant in this regard (Gentner, 1983). As (Turney & Pantel, 2010) point out, the attributional similarity between two words a and b, depends on the degree of correspondence between the properties of a and b and the relational similarity the degree of correspondence depends on the relationship of the pairs a:b and c:d. Consider the example of *man* and *woman*. These have an attributional similarity because they are examples of humans. Whereas in the cases of *human : speak*, and *dog : bark* there is a relational similarity because they are examples of communication techniques. The relational similarity can be understood as “analogies”.

Distributional semantic models (DSMs) are particularly successful in tasks like synonym detection and concept categorization because these words share many properties (attributional similarity), such as synonyms and co-hyponyms (Baroni & Lenci, 2010). Their paper also describes how relational similarity is more suited to concepts related by hypernymy. Synonyms are words that have exactly or nearly the same meaning like

‘intelligent’ and ‘clever’. A hyponym is a subconcept of a broader term, and co-hyponyms are subconcepts of the same broader term like ‘rose’ and ‘daisy’ are both ‘flowers’. A hypernym is the broader meaning of a particular word, like ‘flower’ is the hypernym of ‘rose’.

In the following section we discuss the techniques for creating embeddings:

1. LSA
2. Glove
3. Word2Vec

## 2.1 LSA

One of the techniques of creating vectors uses a term-document matrix and is called LSA (Latent semantic analysis). Here the main objective is to identify the similarity in documents but it can also be used to find similarity between terms (words). As the paper by (Dumais, 2005) points out, there are mainly 4 steps in the analysis:

### 2.1.1 Term-document matrix:

This a matrix of counts with words as rows and the document or sentence as the column. The count of the words within the document or sentence are populated in the matrix. The order of the words does not matter and hence the terms “bag of words” is used to describe this method.

### 2.1.2 Transformed term-document matrix:

The raw word frequency counts are transformed cumulatively in a sublinear (log) function and inversely with the overall occurrence of the term in the document (inverse document frequency).

### 2.1.3 Dimension reduction:

A singular value decomposition (SVD) is performed to reduce the matrix to smaller dimensions by retaining the largest ‘k’ singular values and setting the remaining to 0. The vector space derived by the SVD now contains vectors of the words and documents.

### 2.1.4 Retrieval in reduced space:

Similarities are computed between words and documents in the reduced dimensional space rather than the original matrix. If a search is made to retrieve documents, then a new vector is formed using the constituent term vectors, and this is compared with the document vectors to find a match.

In essence this model captures information about the document as a whole and through this process it also gathers information about the different words. Each document has a vector defined by the words which have maximum representation in the particular document. The drawback of this as pointed by (Landauer & Dumais, 1997) is that the context window for the words is the size of the document (large), and the best performing algorithm was found to have a small context window of 5 with the center word as the target. This is the only pure “count based” approach we explore, which does not involve a neural network under the hood. This is widely considered as a global matrix factorization method.

We do not implement this model because we are not interested in identifying what a post is specifically about. Rather, we want to understand what a specific word means as a representation of other words. So, in our case a word-word co-occurrence matrix would be apt. However, we carry forward the idea of using a context window of 5 words in our co-occurrence matrix.

## 2.2 Glove

In the previous method, LSA, word analogies are not possible due to the typically large document size and loss of local context information. The Glove method, short for global vectors, introduced by the team at Stanford (Pennington et al., 2014) tries to leverage the global document statistics like LSA but instead of term-document matrix it relies on the use of a word-word co-occurrence matrix:

- The method hinges upon the idea of comparative probabilities as in the paper by the Stanford team. Consider a word ‘k’=[solid, gas, water] and its co-occurrence with 2 other words ‘i’=steam and ‘j’=ice. If a word is similar to both words like ‘water’:>{‘ice’, ‘steam’} or dissimilar to both words like ‘water’:>{‘dog’, ‘cat’} would result in both high probabilities or both low probabilities which would approximately equal to the value 1. Whereas for relations like ‘solid’:>{‘ice’, ‘steam’} or ‘gas’:>{‘ice’, ‘steam’} we can see that probability of one of the pairs in each of these relationships will be high resulting in a value greater than one or lesser than 1.
- Words with no discerning ability like the one which are used in all word contexts ‘a’, ‘the’ etc, called “stopwords”, and also other words which are not relevant to the target word, are identified with the use of global statistics.
- The main issues in considering the co-occurrence matrix for both global and local statistics are:
  - There is no direct equation to represent the relation between the 3 word vectors and their relative probabilities of occurrence. Some function has to use the 3 words as arguments. This would be important in the calculation of analogies.
  - The word vectors are high dimensional vectors whereas the probabilities are scalars. This has issues in computation.
  - Since are 3 words involved, the loss function computation for this is tricky.

In comparison with the LSA model, this is similar in that it is also a matrix of counts. However, there are no transformations which takes place on the matrix. These counts are used as is. The complexity reduction is mainly executed by the logic in the weighting scheme. SVD as used in LSA sees decreased performance with a larger corpus.

In our project, some of the main benefits of the global document statistics like “stopwords” are addressed by their explicit treatment and removal using existing stopword indexes. Also, since the language use in social media tends to be more dynamic, the local statistics may be the foremost consideration rather than global. New posts are constantly coming in, and an ideal data pipeline must include a better alternative than repeatedly creating computationally expensive global co-occurrence matrixes.

## 2.3 Word2Vec

Finally, we explore Word2vec introduced by ex-google's (Mikolov et al., 2013) which is the model we implement in this project. We use the word2vec implementation from Gensim package (*Models.Word2vec – Word2vec Embeddings — Gensim*, n.d.). New posts (text information) can be added to an already trained model to update the weights and create new embeddings. Since this is based on a neural network architecture, the dimensions representing each vector are trained in a linear time complexity (very fast). Listed below are some of the significant details as relevant in our use case.

- We use the SGNS (skip gram negative sampling) algorithm, setting the window size to 5 and the vector size to 300 dimensions. In the study behind the glove method (Pennington et al., 2014), the vector size and window size were found optimal for values of 300 and 5 respectively (figure below).

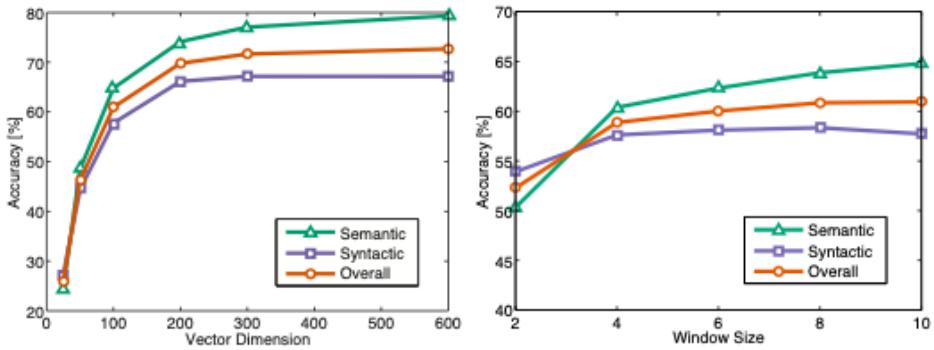


Figure 1: Vector dimension and window size selection (Pennington et al., 2014)

- The skipgram model is tuned to predict the context words from the center word. This is found to be better at creating word embeddings for similarity than the CBOW (Continuous bag of words) model which predicts the center word based on the context words.
- In this model, the words in a context window are weighted by probability of occurrence as opposed to the method in glove where it is harmonic function (Levy et al., 2015). It is also faster to train than glove.
- Words larger than some threshold are randomly removed through a process called subsampling. In this way, stopwords are also implicitly removed.

We will use the word2vec model to identify most similar words and also visualize the words in a 2-dimensional space. The closer the words are in this plot, the more semantically similar they are assumed to be. Here we will have to use our discretion in drawing conclusions from potentially "similar" words because they may be similar in aspects other than synonymity as well.

## 2.4 User types

Different groups of people use different language styles to describe their content. In real life we see an obvious distinction in the way a teenager speaks versus the way an individual in his late 20's and older would speak.

A consumer may express an interest or passion towards a certain idea through the usage of certain hashtags and the post description (text). Whereas an influencer or business may use the same hashtags along with a call to engage, like requests to follow, share etc.

An influencer could be influential purely due to celebrity status gained from another avenue like movies, sport etc. Parts of their fanbase can be quick to resonate with certain values and ideas expressed by them. An influencer could also be someone who offers credibility to the post because of a strong track record of selection through review and study, in the specific domain.

On reviewing the existing literature on user classification based on social media posts, we find an in-depth study based on health related social media post data. (Rivas et al., 2020) derive classification models for age, gender, ethnicity, and location based on data from twitter, google +, and specific patient drug forums. These involve manual classification techniques. We do not pursue classification of the posts as the main objective of this study, but we include an option in the dashboard tool to pursue manual classification in the eventuality that we train a model based on user classified data. However, an indicator metric of influential posts is created in this project called likes\_per\_minute, which is considered in a bivariate setting along with time elapsed, to filter out (influential) posts which are greater than a certain z-score threshold. Here it was found important to factor time\_elapsed with the likes\_per\_minute because it seems to have an effect on it.

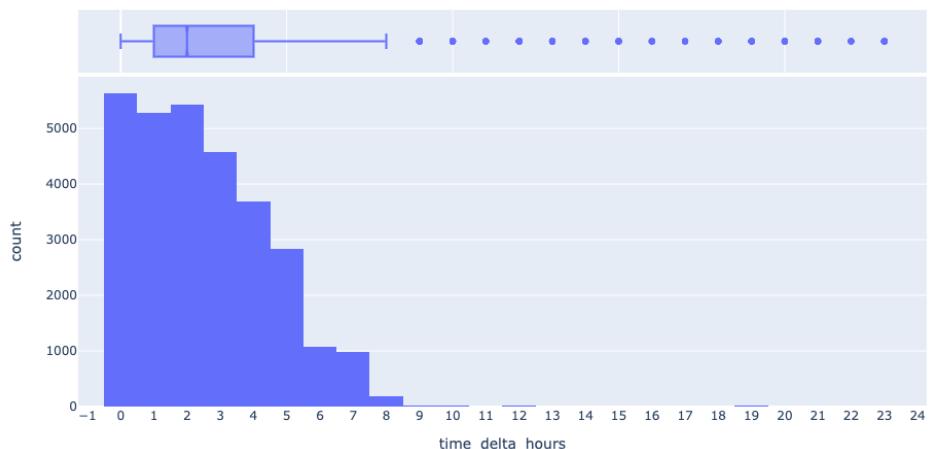


Figure 1: Count of posts by elapsed time delta

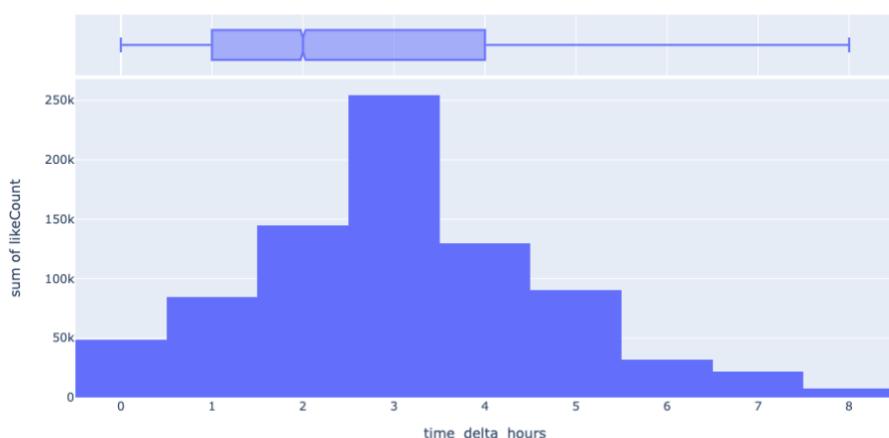


Figure 3: Sum of likes by elapsed time delta

In figure 3, we can see that the sum of likes increases till time\_delta of 3 hours and after that it starts reducing. The reduction can be explained by decreased number of posts. The increase seems a little strange because the count of posts (figure 2) falls approximately in the same range till 3 hours, and in fact dipping in count of posts from the 3 hour time delta mark.

When we look at figure 4, we see that the sum of likes is pronounced for a time delta of 3 hours because of a very popular post which has over 150k likes. To reduce the effect of such a large outlier on the selection of “influential” posts in the entire corpus, we can factor the time delta in bivariate relationship. So this balances the effect of the scraping time on detecting influential posts. However, since we have categorized the posts in terms of hours, the time\_elapsed resolution of a post for optimal segregation using the bivariate z-score, is 1 hour. Posts which are less than an hour may be more sensitive to time\_elapsed rather than actual “influence”.

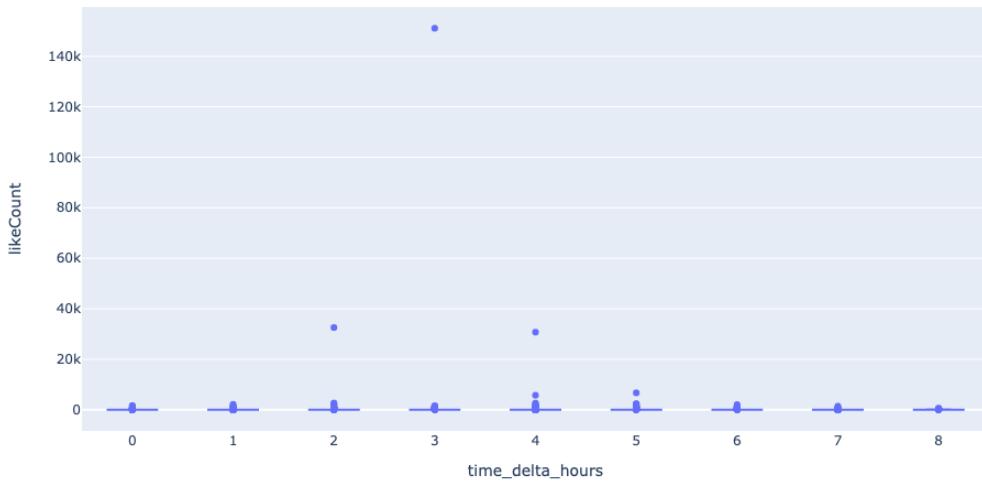


Figure 4: Outliers in like count by time delta

### 3: Methodology

In this section, we will discuss the bulk of the work done in terms of the python coding for cleaning the data, the data pre-processing steps, and the method of selecting and tuning the models.

#### 3.1 Data scraping:

Sourcing the data was the first step of the research. Instagram currently does not have a publicly accessible data API from where we can query the data. So other methods of data collection were tried and the best method was chosen for the collecting the large data corpus. Also there were found to be no existing banks of data containing the Instagram post level data that we required.

In our project we wanted to obtain recent data containing the Instagram post descriptions for posts containing a particular hashtag (#organic). A post would normally contain this

hashtag if the image or video is relevant to what the #hashtag means. The meaning of the hashtag is subject to the users perception of what it means. Since language is constantly evolving, although a word might have a limited scope of understanding in its initial usage, the meaning can grow to accommodate a broader scope as the natural usage by the user evolves. Here we list out the different methods attempted and describe in more detail the method which was finally used.

### **3.1.1 Selenium:**

First, coding a scraping bot using Selenium python package was considered. This would interact with the objects of the web page and execute the manual tasks of clicking and viewing the data (Python Simplified, 2020). The data from Instagram has to first load in the browser, by viewing it in page, and then the required objects can be populated into a data table. This process was found to be too time consuming to obtain a large text corpus for our research and would directly violate the terms of service of Instagram which forbids any bots on the platform.

### **3.1.2 Steviesie data platform:**

Here the user has to access the pages manually on a web browser (like normal) and make sure to scroll through the required posts. Then, using the inspect tool of the web browser (Google Chrome in our case), we download the .HAR file which contains a mirror representation of all the data objects that has been displayed on the screen. This .HAR file is then uploaded to the Steviesie data platform (*Scrape Data Legally / No-Code Data API & Website Scraper / Steviesie Data*, n.d.), which parses the data and creates data tables out of similar types of objects. So if we require the “post” data, then we just need to select the relevant output from the parsed files on the platform, and then download the file in .csv format. However, even this method was found to be time consuming because of the same buffering bottlenecks of the previous method, but with an added element of manually having to scroll through the posts.

### **3.1.3 Phantombuster data platform:**

Finally, a commercial method was found which provided a fully automated scraping solution for Instagram data. Selected hashtags were queried and on each run about 5000 posts were retrieved in a span of about 15 minutes. Observations were made with regard to the data scraping velocity of this tool.

Hashtag popularity on the platform (Instagram), affects the number of posts retrieved. For example, 3 hashtags were scraped at different run times to check the span of time covered by the retrieved posts. The table below describes the time of each of a few runs, which we then use to estimate the ideal automation parameter for the scraping bot.

Table 1: Data scraping time overlap

Hashtag   Date	Time stamp of first and last post scraped (hh:mm)				
	11Aug_1220	11Aug_1038	11Aug_0020	10Aug_1435	9Aug_1914
#organic	10:19 to 3:18	8:38 to 00:43	22:19 to 16:50	12:35 to 13:55	17:14 to 14:36
#sustainable	10:31 to 8:17	8:49 to 16:15(pd*)	22:31 to 19:05	12:43 to 20:21	17:16 to 7:01
#fairtrade	-	9:02 to 10:31(pd*)	22:31 to 10:33	12:55 to 10:33	9-17:26 to 6-15:53

\*pd-previous day

In the above table the column headers refer to the date and time of running the query and the rows represent the different hashtags. From the above table we can infer the following:

- The more popular a hashtag is, i.e it has more posts, then there are more posts scraped within a certain unit of time. The spans of time for 5000 posts of each of the hashtags are:
  - #organic 7:00 hrs, 7:55 hrs, 5:29hrs
  - #sustainable 16:34 hrs, 16:23 hrs
  - #fairtrade :~106:00 hrs but only up to 3 days of previous data available.
- As seen above #organic is most popular of the selected hashtags. Since each run can pull 5000 posts, how can we automate the pull requests so that we can cover all times of the day:
  - If we set the bot to automatically query the data 4 times a day, then every 6 hours the new data is pulled. Since each query covers around 7 hrs of #organic posts, we will have a seamless pull of data across the day, with possibly a 1 hour duplication due to overlap.
- # sustainable contains 5000 posts across 16 hrs. So it is roughly 2.5 times less popular than #organic in terms of posts per unit time. #fairtrade however has the lowest number of posts per unit time and hits the Instagram backtrack query limit of 3 days before 5000 posts are pulled. Total posts containing hashtags at time of data collection:
  - #organic : 52.6M
  - #sustainable : 11.4M
  - #fairtrade : 3.4M

#### Conflict of Instagram session while data is being pulled:

While the phantombuster data scraping tool is carrying out a scheduled query, if the user (me) accesses Instagram on another device (like a smartphone), then Instagram blocks the user ID with a warning. This is because Instagram does not encourage running a bot in a session. However, the data itself is openly available if a user chooses to manually access it. So it is technically not private data we are scraping. Instagram allows for an API to fetch the data if there is a commercial interest so phantom buster is a noted scraping player which is sometimes advertised on Instagram itself. Since the dashboard created here is not yet a commercial product, usage of this tool was halted after the first warning without pursuing any extended data collection.

Data fields contained in the queried data:

The below table describes the different fields available and provides an example of the data. Explanations of the field types are self explanatory but some clarifications are listed below.

The data table produced from the phantom buster tool consists of the following fields:

Table 2: Example data fields scraped from Phantombuster

Field	Example
profileUrl	<a href="https://www.instagram.com/myhomefarming">https://www.instagram.com/myhomefarming</a>
username	myhomefarming
fullName	Home Farming
commentCount	0
likeCount	1
pubDate	2021-08-10T12:34:56.000Z
description	Planting sweet potatoes using a rice bag. In the future, I will show you how I do it step by step. It is easy for home farming.  #urbanfarming #gardening #hydroponics #urbangarden #urbanfarm #growyourownfood #organic #hydroponic #urbangardening #organicgardening #agriculture #microgreens #urbanfarmer #homegrown #vegetables #homefarming #sayurorganik #sayur #kebun #kebunrumah #healthyeating #sustainablelifestyle #vegetablegardens #sweetpotatoes
imgUrl	<a href="https://scontent-lhr8-1.cdninstagram.com/v/t51.2885-15/e15/234902054_1511876039156842_3394796396722981677_n.jpg?_nc_ht=scontent-lhr8-1.cdninstagram.com&amp;_nc_cat=110&amp;_nc_ohc=Jo3LSzWF7coAX-T5FL6&amp;edm=ABZsPhsBAAAA&amp;ccb=7-4&amp;oh=e026bdf0c9656b16fbca85b7e76b6fb3&amp;oe=61186640&amp;_nc_sid=4efc9f&amp;ig_cache_key=MjYzNzIwMTA5NTEwNzAwMDQ0Ng%3D%3D.2-ccb7-4">https://scontent-lhr8-1.cdninstagram.com/v/t51.2885-15/e15/234902054_1511876039156842_3394796396722981677_n.jpg?_nc_ht=scontent-lhr8-1.cdninstagram.com&amp;_nc_cat=110&amp;_nc_ohc=Jo3LSzWF7coAX-T5FL6&amp;edm=ABZsPhsBAAAA&amp;ccb=7-4&amp;oh=e026bdf0c9656b16fbca85b7e76b6fb3&amp;oe=61186640&amp;_nc_sid=4efc9f&amp;ig_cache_key=MjYzNzIwMTA5NTEwNzAwMDQ0Ng%3D%3D.2-ccb7-4</a>
postId	2637201095107000300
ownerId	48877898830
type	Video
query	#organic
timestamp	2021-08-10T12:35:19.268Z
videoUrl	<a href="https://scontent-lhr8-1.cdninstagram.com/v/t50.2886-16/234259757_428427111750827_6156703384781018693_n.mp4?efg=eyJxZV9ncm91cHMiOjIjbXCJpZ19wcm9ncmVzc2l2ZV91cmxnZW4ucHJvZHVjdF90eXBILmZIZWRcl0ifQ&amp;_nc_ht=scontent-lhr8-1.cdninstagram.com&amp;_nc_cat=111&amp;_nc_ohc=qitd3o_ju5YAX8wm9p0&amp;edm=ABZsPhsBAAAA&amp;ccb=7-4&amp;oe=61149592&amp;oh=073607c29be71b62078e62757cdfce00&amp;_nc_sid=4efc9f">https://scontent-lhr8-1.cdninstagram.com/v/t50.2886-16/234259757_428427111750827_6156703384781018693_n.mp4?efg=eyJxZV9ncm91cHMiOjIjbXCJpZ19wcm9ncmVzc2l2ZV91cmxnZW4ucHJvZHVjdF90eXBILmZIZWRcl0ifQ&amp;_nc_ht=scontent-lhr8-1.cdninstagram.com&amp;_nc_cat=111&amp;_nc_ohc=qitd3o_ju5YAX8wm9p0&amp;edm=ABZsPhsBAAAA&amp;ccb=7-4&amp;oe=61149592&amp;oh=073607c29be71b62078e62757cdfce00&amp;_nc_sid=4efc9f</a>
viewCount	5

- Each row of the data table refers to a post. The post can either be an image or a video. If it is a video, then there is an additional field called ‘viewCount’ which may be populated.
- ‘commentCount’ and ‘likeCount’ are engagement metrics. This can indicate the popularity of the post with other users. A user can either write a comment on the post or hit the “like” button for the post.
- ‘pubDate’ refers to the date and time that the post was posted.

- ‘description’ is the text field which is what we use to infer embeddings in our project.

## 3.2 Data preprocessing:

This section deals with all necessary preprocessing in order to get the data to a form in which we can carry out analysis.

### 3.2.1 Dataframe\_1

The first dataframe is created by combining the outputs of the individual query runs from the phantombuster data tool. This dataframe includes all the different hashtags as well, which we had used for testing the time span of posts (including #fairtrade #sustainable). The python code (appendix python code line 22) takes these different files and concatenates them along the same columns. If there are no new individual queries and the old combined dataframe is the latest, then the old dataframe is loaded into the program.

- There were a total of 87,423 rows in this dataframe.
- Includes the hashtags #organic, #sustainability, #fairtrade.
- A total of 16 datasets were combined.
- This raw dataframe included 351 empty rows. So there are a total of 87072 non-empty rows in the raw dataframe.

### 3.2.2 Dataframe\_2

From Dataframe\_1 the description column is then treated to extract relevant data into its unique columns. This creates the Dataframe\_2. The new columns created are:

- hashtags – consists of the hashtags. Some hashtags are not preceded by a space. So it is first treated to include a space so that text splitting functions can be properly applied.
- cap\_mentions – this consists of the users that have been mentioned preceded by an ‘@’ character.
- web\_links – any websites mentioned in the description will be populated in this field.
- clean\_captions – this consists of the plain text after all the previous types of text have been extracted into their relevant columns.

### 3.2.3 Dataframe\_3

Next, the clean\_captions column which contained just the text data was processed with the following transformations:

- font\_uniformity – the encoding of different fonts causes mismatch of word representations. So the fonts have to be made uniform. ‘caption\_processed’ column.
- convert\_lower\_case – Next all the words are converted to lower case. This ensures that the words starting with a capital letter are represented properly. ‘caption\_processed\_2’ column.
- remove\_punctuation – All special characters are then removed. This include "!"#\$%&()\*+- .•/:;<=>?@[{}]^`{\}`\n". ‘caption\_processed\_3’ column.
- diff\_encodings – Sometimes the punctuation marks like the triple dot (...) and inverted commas (‘) are of a different encoding. These had to be removed separately. ‘caption\_processed\_4’ column.
- Emojis – after carrying out the above steps independently, emojis still remained. These were removed using the “preprocessor” library before removing stopwords.

For identifying the different treatments that were required on the data, the dataframe table had to be explored in the table form. Exemplars of a certain row which had unique fonts and punctuations were selected and viewed as test cases to check if a particular transformation was working. Shown below in figure 5, are the columns after the relevant preprocessing steps.

	clean_captions	caption_processed	caption_processed_2	caption_processed_3	caption_processed_4
11	Health is precious - protect it. For more information contact: Shatayu Ayurved - 9826092380. . .	Health is precious - protect it. For more information contact: Shatayu Ayurved - 9826092380. . .	health is precious - protect it. for more information contact: shatayu ayurved - 9826092380. . .	health is precious protect it for more information contact shatayu ayurved 9826092380	health is precious protect it for more information contact shatayu ayurved 9826092380
13	Shatkratu is one of the 10 sacred flowers (Dashpushpam) in Kerala, also called "Love in a Puff". Enhances hair growth, thickening & reduces hair fall strengthening the hair roots. . .	Shatkratu is one of the 10 sacred flowers (Dashpushpam) in Kerala, also called "Love in a Puff". Enhances hair growth, thickening & reduces hair fall strengthening the hair roots. . .	shatkratu is one of the 10 sacred flowers (dashpushpam) in kerala, also called "love in a puff". enhances hair growth, thickening & reduces hair fall strengthening the hair roots. . .	shatkratu is one of the 10 sacred flowers dashpushpam in kerala also called love in a puff enhances hair growth thickening amp reduces hair fall strengthening the hair roots	shatkratu is one of the 10 sacred flowers dashpushpam in kerala also called love in a puff enhances hair growth thickening amp reduces hair fall strengthening the hair roots
23	She's here! Judith is here! She has cucumbers, MELONS, courgettes, herbs, chard, plums, beetroot, MELONS... did I mention MELONS? I'm very excited about this little collab. She's a star and I can't wait to see what she brings to She's already sold a big bag full to a regular so you'd better get down here quickly!	She's here! Judith is here! She has cucumbers, MELONS, courgettes, herbs, chard, plums, beetroot, MELONS... did I mention MELONS? I'm very excited about this little collab. She's a star and I can't wait to see what she brings to She's already sold a big bag full to a regular so you'd better get down here quickly!	she's here! judith is here! she has cucumbers, melons, courgettes, herbs, chard, plums, beetroot, melons... did i mention melons!? i'm very excited about this little collab she's a star and i can't wait to see what she brings to she's already sold a big bag full to a regular so you'd better get down here quickly!	she's here judith is here she has cucumbers melons courgettes herbs chard plums beetroot melons did i mention melons im very excited about this little collab she's a star and i can't wait to see what she brings to shes already sold a big bag full to a regular so you'd better get down here quickly	shes here judith is here she has cucumbers melons courgettes herbs chard plums beetroot melons did i mention melons im very excited about this little collab she's a star and i cant wait to see what she brings to shes already sold a big bag full to a regular so youd better get down here quickly
106	Brownie balayage .	Brownie balayage .	brownie balayage .	brownie balayage	brownie balayage

Figure 5: Caption processing of exemplars

### 3.2.4 Language detection

A language detection library was used to identify the language of a particular caption. Here an offline model called “langdetect” (Danilak, n.d.) was used and it took 11:30 mins for the classification of 87072 posts captions on a macbook air 2017 laptop system. Other models involved an online api method of language classification (textblob which uses a google api) but seemed to throw a “too many requests” error.

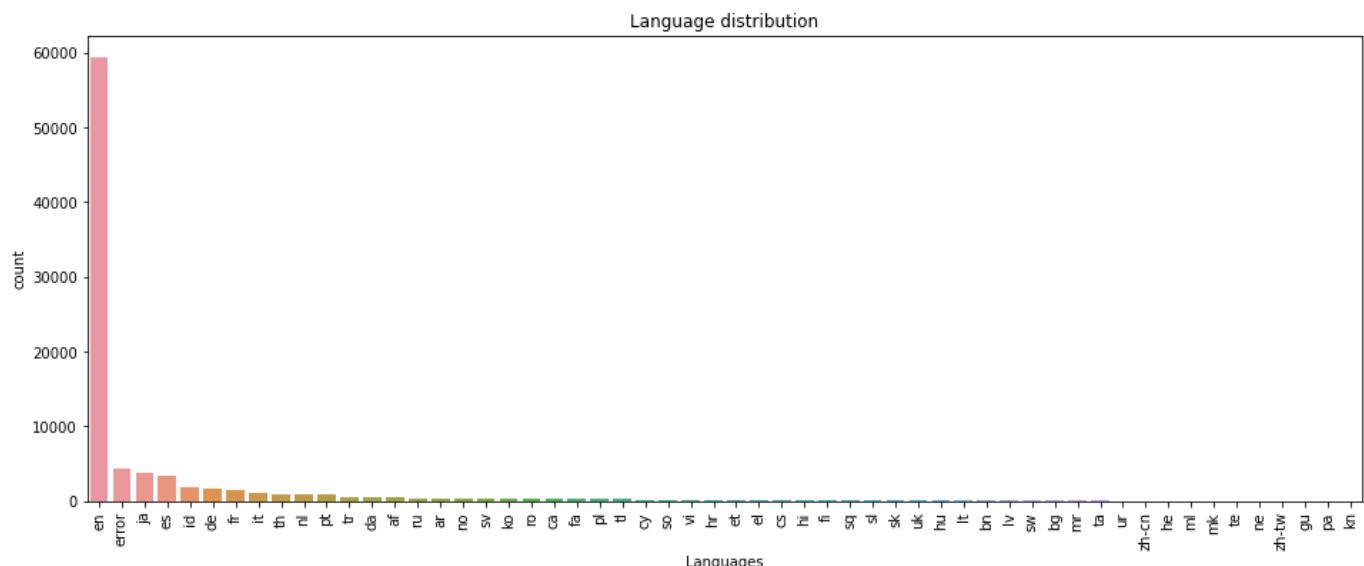


Figure 6: Language countplot of initial data

As seen in the above countplot, figure 6, most of the posts are in the English language. The second highest categorization are the errors which are mainly due to there being no data in ‘clean\_captions’, after separating all other textual data like ‘hashtags’, ‘web\_links’ etc. into their respective columns. The errors also arise if there are only *emojis* left in description after all other

text data has been separated into their columns. These error rows may contain useful hashtags, but since they do not make sense sequentially (like descriptions using a certain language), they can be safely discarded.

- The languages short forms are as per the list of ISO 639-1 codes (“List of ISO 639-1 Codes,” 2022).
- We can see that the top languages on this dataset apart from English are: Japanese, Spanish, Indonesian, German, French, Italian, Thai, Dutch, Portuguese, and Turkish.
- The languages give us an indicator of which countries are more interested in the concept defined by these hashtags (#organic, #sustainability, #fairtrade). Since English can be considered as a global language, the only reliable inferences we can make are that
  - Users in Japan and Spain have high interest.
- For the purpose of using the text to interpret meanings, we will have to use a single language. The language selected was English which formed the bulk of posts. The dataframe was filtered to only English posts. This leaves us with 59374 rows of data. This is a drop of 33% from the preprocessed “Dataframe\_3”.

### 3.2.5 Further filtering

After the above transformations, it was decided to focus on only 1 type of hashtag to maintain a uniformity and depth of inference. The dataframe was filtered to only #organic which makes up most of the dataset.

- #organic posts made up 40198 posts. This is again a drop of 33% from the “English” only posts.
- Finally, duplicate posts (based on postId) were removed which consisted of about 5067 posts. After that duplicate posts based on same description were eliminated which amounted to another 4891 posts. The remaining dataset contains 30240 rows.

### 3.2.6 Stemming and Lemmatization

Stemming is the process by which a particular word is reduced to its root word (“Introduction to Stemming,” 2018). For example the words “chocolates”, “chocolatey”, “choco”, are reduced to the root word “chocolate”. We should keep in mind the errors of stemming that can occur, namely over-stemming and under-stemming (Heidenreich, 2018).

- Over-stemming occurs when too much of a word is stemmed to a nonsensical term or different words being reduced to the same stem. Consider the words university, universal, universities, and universe; if these are reduced to the stem term “univers”, then over stemming has occurred.
- Under stemming is the opposite issue where words which are actually forms of each other do not resolve to the same stem.
- A Porter stemmer from NLTK package was used on this dataset. This is considered as a gentler algorithm which offers good reproducibility and dates back to the 1980’s.

Lemmatization is a more calculated process compared to stemming. Here the words are reduced based on their dictionary forms rather than by their common starting letters. For example, words like “is” and “are” can be reduced to “be”.

- Lemmatizers require a lot more knowledge about the structure of the language and hence tends to be more intensive.
- The Wordnet lemmatizer from NLTK was used on this dataset.

### 3.2.7 Word Vocab

All the above steps were carried out to bring the text data into a tokenized form with which word co-occurrences and word embeddings can be created. The number of unique words (vocab) in each of the different processing steps are:

- ‘cp4\_no\_emojis’ contained 46066 unique words after general preprocessing of converting to lower case, removing punctuation, and removing emojis.
- ‘cp4\_rm\_sw’ is the same as the previous column with the added treatment of removing the stopwords. 45925 unique words remained.
- ‘rm\_sw\_stem’, after running Porter stemming on ‘cp4\_rm\_sw’ reduced the vocabulary to 34051 unique words.
- ‘rm\_sw\_lemt’, after running Wordnet lemmatizer on ‘cp4\_rm\_sw’ resulted in an output of 41526 unique words.

Lemmatized text is more responsive to different part of speech, and hence is the more nuanced algorithm we have considered. Although the goal was to reduce the vocab to a smaller size for easier computations, we have to still use the lemmatized text column (with more words than stemming algo) as the input to our Word2vec algorithm.

## 3.3 Word Embeddings

We explore 2 different methods of making word embeddings. The first method uses a traditional distributed semantic model (DSM), and in the second method “word2vec” they are formed using a neural network. We end up choose the word2vec model over the DSM because they significantly outperform the DSM model (*On Word Embeddings - Part 3*, 2016).

The *vocabulary* is a list of unique words of text from the dataset. This corpus of text is obtained from the lemmatized column ‘rm\_sw\_lemt’ which also omits stopwords and contains other preprocessing as described in the previous section.

Stopwords removal is an important aspect. *Stopwords* are a set of commonly used words in a language. The reason why stop words are critical to many applications is that, if we remove the words that are very commonly used in a given language, we can focus on the important words instead (Ganesan, 2014), and reduce computations.

Each word in the vocabulary (vocab), is represented as a vector and then reduced to 2 dimensions to give us a visual representation on a 2D plot.

### 3.3.1 DSM

#### 3.3.1.1 Co-occurrence matrix

With the token word at the  $i^{th}$  position, a rolling bidirectional window is used to define the co-occurring words with 2 words in the front and 2 words behind the token word. The token word is iteratively selected from each word of the row, and is paired in a tuple with the surrounding words within a 5 word window. Then the counts of co-occurrence are updated as the value to the unique (word,word) tuple pair key of a dictionary.

- The vocab consists of 41526 words resulting in a co-occurrence array of shape  $41526 \times 41526$ . So each word is represented as a vector of 41526 dimensions.
- This is a sparse matrix (many redundant 0's) which makes it computationally very expensive.
- These vectors must be reduced to 2 dimensions while preserving the relationship of the words in the multidimensional space.

- For this an SVD matrix factorization can be performed to identify the most significant contributors to a vector's identity and represent it in 2 dimensions.

### 3.3.1.2 SVD matrix

To address the computational cost of working with sparse matrixes of very large dimensions, we can use the COO formatting method in sklearn library (*Scipy.Sparse.Coo\_matrix — SciPy v1.7.1 Manual*, n.d.). Without this it is not feasible to perform matrix factorization operations. The advantages of using this format are:

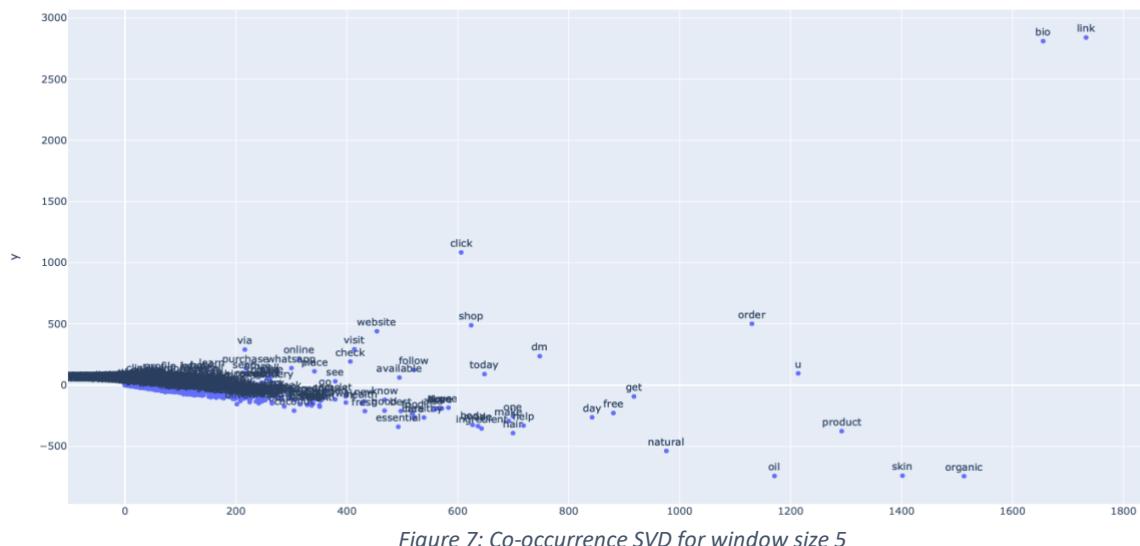
- facilitates fast conversion among sparse formats
- permits duplicate entries
- very fast conversion to and from CSR/CSC formats (compressed sparse row/column)

This matrix then has to be upcast to a floating point format before carrying out the SVD.

We use the truncated SVD method from the sklearn library to decompose the multidimensional matrix into 2 dimensions. This method is also called LSA (latent semantic analysis). Contrary to PCA, this estimator does not center the data before computing the singular value decomposition. This means it can work with sparse matrices efficiently (*Sklearn.Decomposition.TruncatedSVD*, n.d.). The `fit_transform` method is applied on the preprocessed co-occurrence array (COO formatting) and the output is an array with 2 columns and the number of rows defined by the vocab – shape (41526,2).

### 3.3.1.3 Plotting

The co-occurrence matrix contains the index and columns of the vocab in sorted order. This means that when we obtain the SVD matrix the index is in the sorted order here as well. Since the SVD matrix does not contain the actual word, the vocabulary to plot must be referenced by the index value. So for example if we want to plot a word which is in the 5<sup>th</sup>



As seen in the below plot, most of the words are very densely concentrated on the left-hand side and it funnels out towards the right. The right hand side, which have words that are much further away in terms of distance, is indicative of the words which occur most from the co-occurrence values.

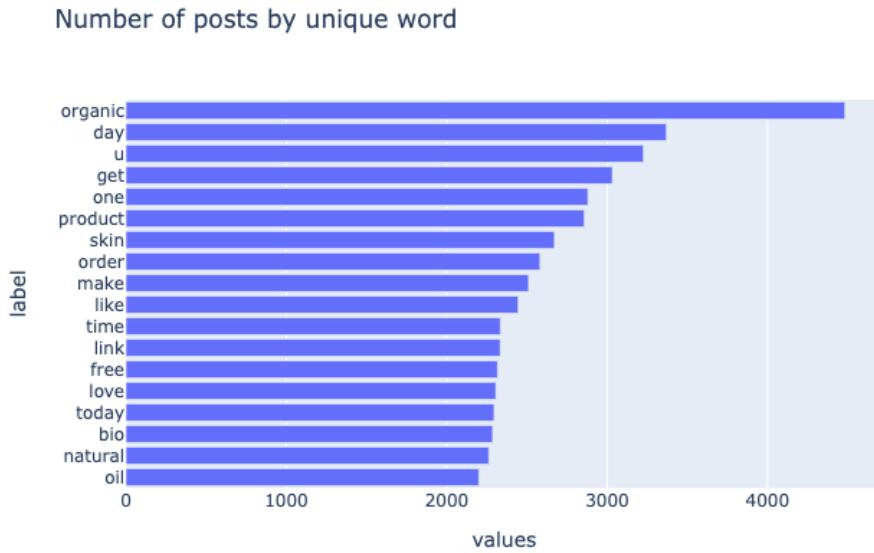


Figure 8: Count of posts by word

In the above bar chart, figure 8, the top results are the same words which are spread out towards the right-hand side of the plot in figure 7. The words seem to have enough information to assume a more specific position in the 2-dimensional representation.

Since most of the words are not adequately separate from each other, we cannot infer much from their 2D representations. It was noted that the model could have been improved by lowering the weights of words further away in the window, setting activation functions, and performing more mathematically heavy processing, in order for the vectors to gain better real-world meanings. However, we did not investigate further into this area because the method was computationally expensive.

### 3.3.2 Word2Vec model

To overcome the lack of context association between the words in the previous model, we use Word2Vec. Here the algorithm uses a large amount of unannotated plain text to learn the relationship between the words automatically (*Making Sense of Word2vec / RARE Technologies*, n.d., p. 2). The output are vectors, one vector per word, with remarkable linear relationships that allow us to do things like  $\text{vec}(\text{"king"}) - \text{vec}(\text{"man"}) + \text{vec}(\text{"woman"}) \approx \text{vec}(\text{"queen"})$ . This seems intuitive to us in our normal usage of these terms in everyday language. The relationships depend on the training corpus. The example mentioned here pertains to a dataset where the context of king and queen have been sufficiently represented with different context words so that it intuitively represents the relationship.

In our dataset however, the words king and queen are not used in normal contexts. So this does not yield the same result. This is the main idea we will leverage in using this model for the Instagram social media data. We use this to identify the underlying meaning of a

particular concept based on the usage of language in these contexts. Since our entire text corpus consists of small summarized pieces of text (descriptions) based on some notion of “organic”, we can assume that the relationships and the implied meanings, apply specifically within a particular domain of selection.

### 3.3.2.1 Instantiation

We will use the word2vec implementation from the gensim package. One of the benefits of this implementation is that it uses optimized C routines, data streaming and Pythonic interfaces. The training is streamed, so ``sentences`` can be an iterable, reading input data from the disk or network on-the-fly, without loading your entire corpus into RAM (*Gensim*, n.d.).

The code was implemented in a Spyder python console and the model was set with the following hyperparameters:

- The skip gram model was used. Here the context words within a window are probabilistically predicted based on the center word.
- vector size is set to 300. This means that each word is represented as a vector with 300 dimensions.
- Window size was set to 5.

### 3.3.2.2 TSNE transformation

With the model instantiated, a TSNE transformation is performed on the vectors. TSNE stands for T-distributed stochastic neighbor embedding. The algorithm maps multi-dimensional data to two or more dimensions, where points which were initially far from each other are also located far away, and close points are also converted to close ones (Smetanin, 2018). This is analogous to the SVD transformation in the DSM method and benefits from smaller dimensions since vector size is 300 dimensions (2 powers of magnitude lesser) . The hyperparameters for the TSNE model were set to the following:

- Perplexity set to 40. This value is effectively the number of nearest neighbors (*Sklearn.Manifold.TSNE*, n.d.).
- n\_components is set to 2. This is the number of dimensions we want to reduce to.
- initialization of the embedding is set to PCA.
- n\_iterations is set to 2500. This is the maximum number of iterations for the optimization
- random\_state set to an arbitrary number 23. Setting a value for this allows for reproducible results over multiple function calls.

### 3.3.2.3 Plotting

After carrying out the TSNE transformation, a dataframe is created and saved to file. The dataframe contains 3 columns, the word and columns representing the x and y axis. These dataframes are then readily loaded into the app during runtime, for exploring the word vectors on a graph. We can explore similar words and their placement from each other, and allows us to infer if a high-level concept is reflected by the directionality between the vectors.

Shown below is the plot of the full vocab. We can already notice a much more uniform distribution of the words, in a fairly symmetrical manner about the x and y axis. We can assume that the uniform and spread out representation, has maximized the information encoded in the 2D. Hence, we can use this to draw inference based on direction and distance.

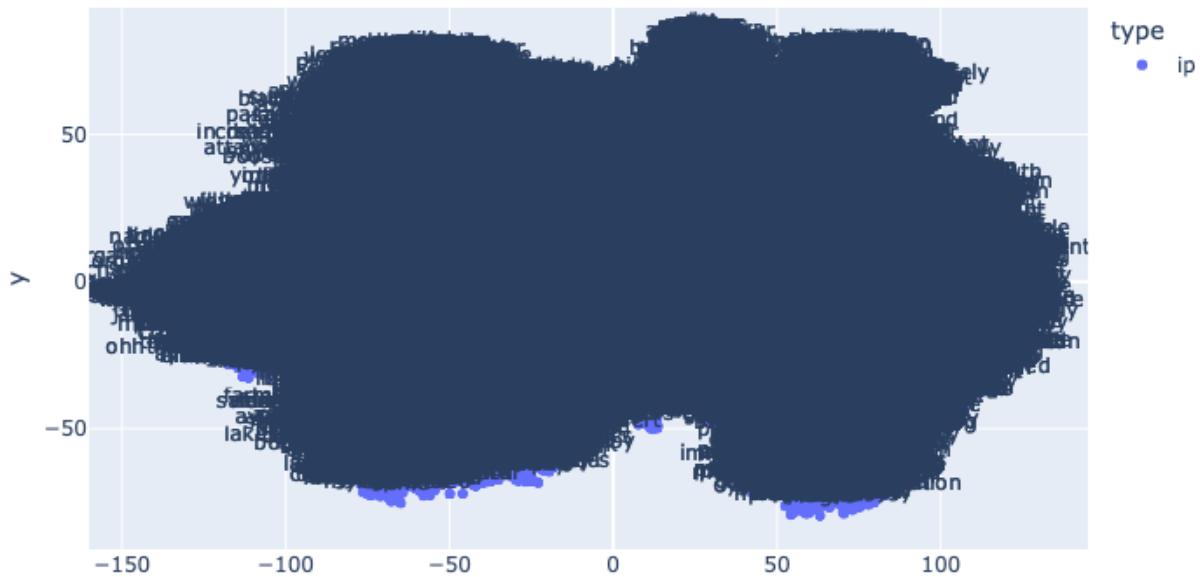


Figure 9: Word2vec represented in 2 dimensions ( 300 dimensions, windows size 5)

## **4: Results and discussion**

To use the dashboard tool, we can start by filtering out the rows of the data table based on a target word. The description column and the hashtag column are the main columns with which we can filter by text. Often, the number of row results while filtering for a word, is more for hashtags rather than the descriptions. This can generally be attributed to the fact that although the word is not used in the description, the user may have used the hashtag because they believe that the post is related to the target word concept in some way. The filtered rows text can be rendered in a wordcloud to allow us to abstract a high level generalization about the text in the filtered selection.

Next, we can use any of the other graphs to gather our intuition about the particular concept. The available visualizations include barcharts for post count by word and co-occurrences, and vector space graphs for checking the words which are close (most similar) to a particular word.

We can also perform simple mathematical operations like addition and subtraction on the word vectors, and they can point towards related terms or concepts. This allows us to play with different combinations of words relying on what feels intuitive to add or subtract or both.

- If word vectors are added it gives a resulting word which is similar to both words, like a middle word. It can be intuited as an average.

- Whereas if we subtract and add a word, then we can determine analogous relationships. For example in an analogy like: king is to man what woman is to ‘\_\_\_\_\_’, we can type it in the form ‘king’ - ‘man’+ ‘woman’ = ‘queen’.

We can use any domain of interest like “textile” or “pharma” and add “organic”, then it should give the word which is closest to the meaning of the summed words, based on the text descriptions in the corpus. Example:

- Consumer + waste = footprint
- Plastic + organic = fairtade
- Textile + organic = fairtrade

The count of co-occurrences plot can be a good place to start for analysing generalizations about our target concept “organic”. The bar chart will display the co-occurring words based on the word typed into the input box. If required, we can use the slider to adjust the number of words to display. Note that some of the words might be omitted because the sidebar of the plot cannot fit it. In these cases, we can zoom into that particular section using the controls on the top right corner of the plot till all the words are displayed on the sidebar.

#### 4.1 The top co-occurrences with the word “organic”

If we type “organic” in the input box of the co-occurrence plot, we get the co-occurring words in descending order. The top words and their associated counts are described in the figure below various high-level insights can be derived.

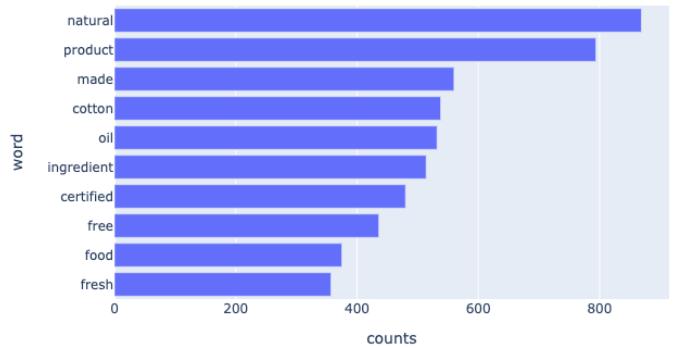


Figure 10: Top co-occurrences in corpus, input word – “organic”

- It has the highest co-occurrence with ‘natural’. As opposed to ‘artificial’. Synonymous.
- It is used most in the context of ‘products’. And this is most associated with the verbs ‘use’ and ‘made’; and the *domain* of ‘skin care’ which we will not further pursue at this point.
- The premier *material* in the ranks is ‘cotton’, which is associated the strongest (most number of co-occurrences) with ‘certification’. The other associations of cotton are:
  - ‘soft’, ‘fabric’ – descriptive words. Cotton is soft, and it is a fabric.

- Interestingly there are also high associations of ‘fabric’ usage with the cannabis industry ‘pot’ and ‘hemp’. On filtering the dash table to “fabric pot” we can get a better picture of why; These are posts on indoor cannabis cultivation using fabric pots.
- ‘bag’- could this be cotton bags? Yes because when we look at the associations for ‘bag’, the foremost *material* is ‘plastic’ followed by ‘cotton’.
- Also we have ‘free’. Filtering the “hashtags” column to #organiccotton and further filtering the “description” column to ‘free’ on the DashTable we can see the context of ‘plastic’ appear in a large number of posts (barchart on the left).
- We could go deeper on the different associations with ‘cotton’ based on our business interests, using the other charts as well.
- Notice that till now we have only used the bar charts for *co-occurrences* and *posts by unique word* respectively.
- Let us now look at ‘oil’. Here ‘essential’ seems to be the dominant co-occurrence. Let us sum these vectors in the *vector math* graph to understand more.
  - Oil + essential = jojoba, carrier, marula, frankincense etc. These are the types of essential oils.
  - On the unique word posts bar chart, ‘skin’ seems to be the premier use case followed by use as an ‘ingredient’.
  - ‘olive’ and ‘coconut’ oils seem to be the most popular *types* in terms of co-occurrences. But these are not ‘essential’ oils because these do not feature in the previous equation.
    - From the dashtable filtering, we can see that ‘coconut oil’ is mainly mentioned in posts for skincare. Infact ‘ingredient’ is most co-occurring with ‘skin’ as a domain of use.
    - Whereas for ‘olive oil’ the wordcloud and bar chart of post mentions suggest application with other cooking ingredients and food.
  - We can similarly explore deeper for ‘oils’ using its other co-occurrences as a start point. Words like ‘hair’, ‘skin’, ‘body’ suggest applications in the beauty industry.
- For ‘ingredient’, if we enter the vector math formula ‘organic’ + ‘ingredient’, we get the resulting words that suggest concepts like ‘purest’, ‘artisanal’, ‘nasties’, ‘leaping’, and ‘unrefined’.
  - When exploring ‘nasties’ on the dashtable we notice that it is better represented with its prefix word ‘no’ as in ‘no nasties’. We miss this in the visualizations because they focus on more significant words. ‘leaping’ is associated with ‘leaping bunny’ which is a certification against animal testing, used for cosmetic and household products.
- In the co-occurrences for ‘certified’, we see that ‘cotton’ has a strong association. When we execute the formula ‘cotton’ + ‘organic’ we get the results like ‘gots’, ‘oeko’, ‘softest’, ‘fsc’, ‘ecocert’. These are all textile certifications, except ‘softest’ which is more of an implied characteristic of organic cotton.
- The word ‘free’ co-occurs most with ‘gluten’, ‘shipping’, ‘cruelty’, ‘chemical’, ‘sugar’, ‘dairy’, ‘paraben’, ‘plastic’. The context of free seems pretty broad. For one it is a concept of something without charge, and on the other it indicates the lack of a

certain harmful substance. We can use the vector math graph to identify concepts that resonate in specific domains.

#### **4.2 Deep analysis around the concept of “organic food”**

#### **4.2.1 Instagram DataTable (DT-1)**

- We can use the first datatable on the dashboard “Instagram DataTable” and filter out the posts containing the keyword “food”. These should automatically include the idea of “organic” because the entire dataset consists only of “organic” posts.
    - Under the ‘description’ column header, there is a space to type an input query. If we type “food” and hit enter, the table gets filtered to post descriptions which contain the word.
    - On filtering the table, we see that the number of pages reduces to 185. Each page contains 10 posts, so we can infer that there are more than 1840 posts about food.
    - We can read the details about each post like ‘description’, ‘hashtags’, ‘cap\_mentions’ – which means other users who have been tagged in the post, and ‘web\_links’ – which are the websites mentioned in the post.
    - The tool is designed so that we can make general inferences about the data. So now with the table filtered to ‘food’, we can click on the *sel\_all* button to select these filtered rows to be rendered in the wordcloud (WC-1).
    - By default, the ‘description’ column is highlighted. Using the dropdown, a different column or multiple columns can be selected. The selected rows and columns are visualized using WC-1 beside the Instagram DataTable. This DataTable only feeds into *WC-1*.

Instagram DataTable					
			Sel_all	Des_all	
x	description				x ▾
#ix	#description	#hashtags	#cap_mentions	#web_links	
filter data	food				
<input checked="" type="checkbox"/> 25	ginger is commonly used as a flavoring agent and spice in foods and drinks :: but this herb has great health benefits and here are some of them :: it is a powerful antioxidant :: ginger has strong anti inflammatory properties :: its an antidote for stomach issues ::	#herb #medicalherbs #organic #giftofnature #herbalism #herbera #herbera_herbs			
<input checked="" type="checkbox"/> 36	this is why i love our food coop the ingredient love ❤ what are you putting love into today	#1 #love #realfood #organic #putloveintoallyoudo			
	want to join us in building a better food system there are plenty of ways to get involved including 🌱 volunteer for the				

*Figure 11: Dashtable DT-1. Allows for filtering and selection of rows and columns*

## 4.2.2 Wordcloud (WC-1) – Word frequency

- The wordcloud offers an intuitive way to see significant words and themes in a large corpus of text (*WordCloud for Python Documentation — Wordcloud 1.8.1 Documentation*, n.d.). The parameters in the word frequency wordcloud were set to the following:
    - Meaningless words are not considered. English stopwords and irrelevant website link words like “https”, “http”, “www” etc were added to the ignore list and not rendered.
    - Collocations is set to true. This means that words which appear together “bigrams” are also identified. These bigrams are only significant for the ‘description’ column since in the other columns there is no inherent meaning of words which appear successively together.
    - The columns that can be visualized apart from the ‘description’ and ‘hashtag’ columns are ‘cap\_mentions’ and ‘web\_links’.
  - Let us explore the columns one at a time to see what trends emerge. It is useful to compare the wordclouds with a baseline of the entire corpus, figure 12, to identify what is changing.
    - The larger the word, the greater is its frequency of occurrence in the corpus.
    - Figure 13, on the right, is for posts filtered to contain “food”.
    - We notice that for ‘food’, some of the words are common from the full corpus WC on the left like: ‘**organic**’, ‘**make**’, ‘**get**’, ‘**one**’.
    - Apart from these, the words which seem pronounced are ‘**body**’, ‘**also**’, ‘**health**’. So we know that the word ‘food’ is associated with the concept of ‘body’ and ‘health’.
    - Using the zoom option of WC-1 on the dashboard, we can explore the significance of other words at smaller sizes.
    - The insights drawn from this are broad and feel based. At this stage, the analysis is not empirical or conclusive, but gives us a direction to approach

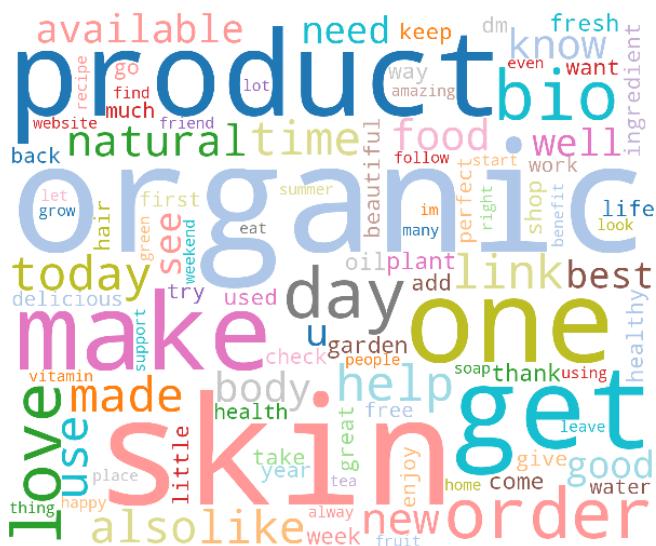


Figure 12: WC-1 based on entire corpus. Word frequency.



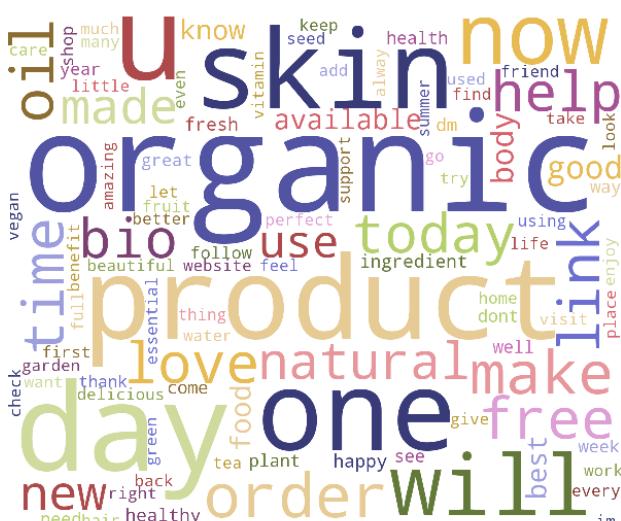
*Figure 13: WC-1 based on "food" filtered posts. Word frequency.*

#### **4.2.3 Z-score DataTable (DT-2)**

- This table looks and functions similar to the previous one with some additional features:
    - Z-score slider from 0-4. Default value is 0, in which case the table contains the same number of rows as DT-1(all rows). The Z-score is based on bivariate relationship between the ‘likes\_per\_min’ and ‘time\_delta\_hours’. These columns are selected to normalize the effect of time elapsed on the ‘likes\_per\_min’ values. The outliers indicate high ‘likes\_per\_min’ values suggesting that they may be influential.
      - A value of 3 filters to the top 15 influential posts.
    - We can enter values in the ‘Inf\_non’ and ‘Domain’ columns for manual classification of the posts as “influencer” or “not” and the respective domain like “textile”, “pharma”, “food” etc. This can then be exported for further processing. That is left for future work and we will not discuss it here.
    - The columns are different. Instead of ‘cap\_mentions’ and ‘web\_links’, we have ‘username’ and ‘fullname’.
    - The row and column selections of this table feed into 2 plots: wordcloud *WC-2* and the barchart *BC-1*.

#### 4.2.4 Wordcloud (WC-2) – Post frequency

- Next we can add the insights from the post frequency wordcloud. This is a wordcloud where the size of the word is dependent on the number of posts that a word occurs in. So this wordcloud controls for over emphasis of a word due to there being large amount of text about it.
  - The wordcloud parameters are set to the following:
    - English stopwords are not displayed.
    - Only unique words are considered for each row. So that means collocations also become irrelevant and is set to False.
    - The special columns which can be visualized here are ‘username’ and ‘fullname’.
    - We can render outliers in terms of the bivariate relationship between ‘likes per minute’ and ‘time delta hours’, using the 4 point z-score slider.



*Figure 14: WC-2 based on entire corpus. Post frequency.*



Figure 15: WC-2 based on “food”. Post frequency.

- In the above figures 14 and 15, we can see the number of posts for each word of the ‘food’ filtered rows. We can see that **‘health’ and ‘make’ are the pronounced words and ‘body’ does not feature as prominently (as in the word frequency WC)** suggesting that there are few posts with lot of text containing ‘body’ in the context of ‘food’.

#### 4.2.5 Bar chart (BC-1) – Posts count

- From the table DT-2, we can select the ‘hashtags’ column along with ‘description’. This gives us the distribution of words and hashtags in the posts for ‘food’.
- We can see that hashtags offer more grouping possibilities because the users themselves try to connect the relevance of the content to standard concepts. There are much more posts which contain the hashtag “#organic”, rather than the word.

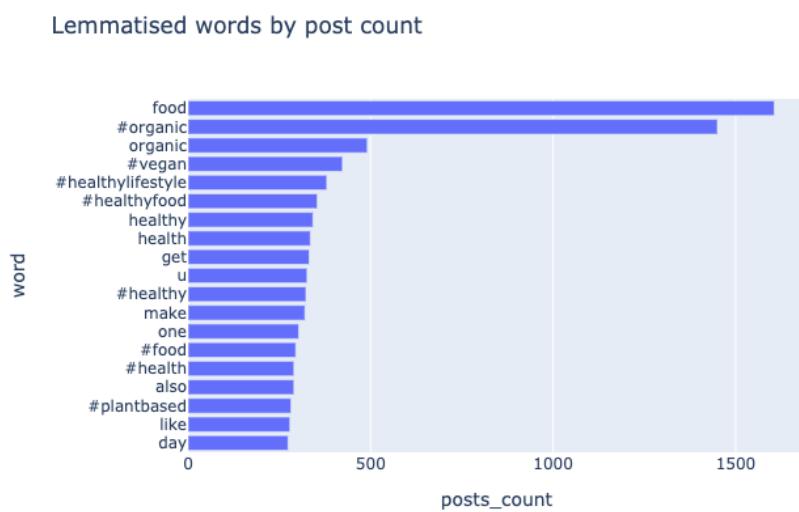


Figure 16: BC-1 number of posts by post count.

- In figure 16, apart from ‘food’ and ‘organic’, the concepts which stand out are: ‘vegan’, ‘healthylifestyle’, ‘healthyfood’, concepts around ‘health’, and ‘plantbased’.

#### 4.2.6 Bar chart (BC-2) – Co-occurrences count

- The co-occurrence count bar chart is particularly useful to understand in what specific context, within a window of 5 words, is the center word used. In the context of ‘food’ posts it is particularly useful in detecting verbs.
- For example, in the previous chart figure 16, we can see words like ‘get’ and ‘u’. When we type ‘u’ in the input box of BC-2 as seen in figure 17, the top co-occurrence is with ‘dm’. ‘dm’ means direct message in Instagram. And the 2<sup>nd</sup> most co-occurring word is ‘follow’. These are *terms used to solicit engagement with the customers, typically used by businesses and influencers*.

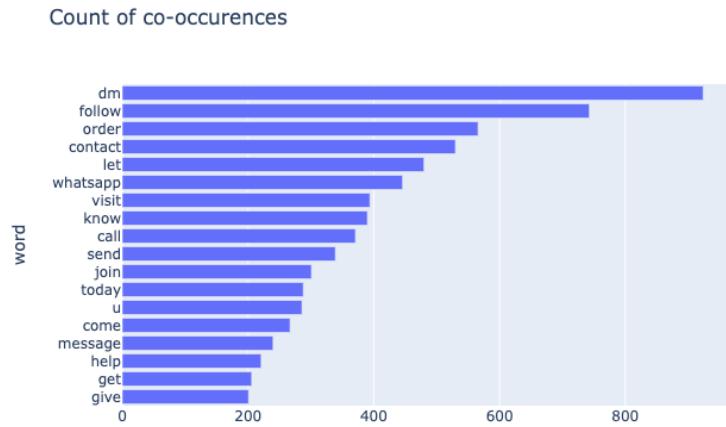


Figure 17: BC-2 word co-occurrences based on input word 'u'.

- On filtering out the description column of DT-2 to ‘dm’, we see that it is most used in the context of ‘us’. Examples ‘dm us’, or ‘follow us’, ‘call us’ etc. Lemmatizing the word ‘us’ has resulted in it being cut short to ‘u’ because it thinks it is a plural.

#### 4.2.7 Vector Exploration (VE)

- The closer the words are on this 2 dimensional space, the closer the points are in the 300 dimensional vector space. Starting with the slider at the 0 position we can increment it in steps of 1 to get the most similar words to the input word.
- The ‘Clr’ – clear button removes all input words. We can then type a word or multiple words (separated by comma ‘,’ ) in the input box.



Figure 18: VE Vector exploration of word2vec words. Most similar to ‘food’ in the corpus

Figure 19: VE Vector exploration of word2vec words. Most similar to ‘make’ in the corpus

- The figure 18 contains the words which are most similar to ‘food’ from the word2vec vector space. The directions and grouping of the words are indicative of different contexts as well. These are not necessarily semantically similar, but they are used in the context of food. These words are themselves encoded with some meaning about food.
- The graph can be used to check which words are most similar to the word typed in the input box. If we type ‘get’ which is another significant word in ‘food’ posts, the top similar words suggest a sense of urgency and rush: ‘grab’, ‘asap’, ‘rush’. This is again indicative of usage in marketing strategies which seek to create a sense of urgency for sales.
- For the word ‘make’, figure 19, the most similar word is ‘swap’. This is followed after a word by ‘prepare’. So ‘make’ could be used in both the contexts, making food (prepare) and making changes(swap) to food consumed. In what context are the food related posts mostly about? We can use the Vector math graph to clarify.

#### 4.2.8 Vector Math (VM)

- The Vector Math graph display requires that at least 3 out of the 4 input boxes contain a valid word from the vocab, in order to display a result.
- Let us type ‘food’ in one input and ‘make’ in the other input beside the ‘+’ sign to signify summation. Addition in the word2vec model outputs a word which is close to both concepts ie similar to both words.

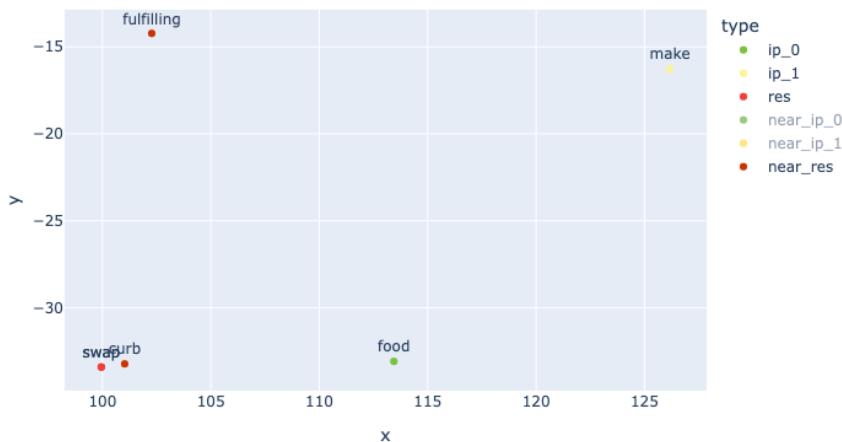


Figure 20: VM Vector Math graph for top 3 most similar words to ‘food’ and ‘make’.

- As seen in figure 20, the most similar word is ‘curb’ followed by ‘swap’. These suggest making changes.
- The next most similar word is ‘fulfilling’ which is popular in garden grown pickles and salads as inferred by filtering DT-1 and observing the render of WC-1.

- Exploring the co-occurrences for ‘make’ in BC-2, we see how it is used in conjunction with different words which can assume totally different meanings. In this sense it can be considered as a stop word.
- Subtraction implies that the result must be similar to the first word and dissimilar to the negative word. This is useful for controlling against a certain concept.
  - For example, when we type ‘food’ + ‘organic’, the most similar word is ‘carna4’.
  - When we check ‘carna4’ using the wordclouds WC-1 and WC-2, we understand that this is a pet food brand which seems to use many words related to organic, in different contexts. So, the model has learned this as the most similar concept to ‘organic’ and ‘food’.
  - We can control for this by subtracting ‘pet’ in the input box. This gives us results which are likely to be related to human consumption. The results are ‘barley’, ‘hearty’, ‘raw’, ‘millet’.

#### 4.2.9 Synthesis

- From WC-1 we were able to identify more potential stopwords (does not add discerning meaning to the ‘food’ data) words like ‘organic’, ‘make’, ‘get’, ‘one’. Controlling for these stopwords, we find that words like ‘body’ and ‘health’ were used frequently.
- In WC-2 we saw that ‘health’ and ‘make’ is pronounced but ‘body’ is not emphasized. Since ‘make’ can be considered a stopword, ‘health’ is the major concept related to ‘food’ when it comes to ‘organic food’.
- Filtering DT-1 and DT-2 to ‘body’ under the description column and ‘food’ under the hashtag column:
  - From the hashtags, it appears that the word ‘body’ in the context of food is associated with ‘healthyfood’ and ‘healthylifestyle’.
  - From WC-1 and WC-2 we can see that ‘help’ is a prominent word in this selection of data. Reading some of the posts we can see that these are posts that ‘help’ the ‘body’ in some way or the other.
- Lateral thinking approaches:
  - Let us think about what organic food could also mean. It must relate to organic methods of growing right? But why should there be significance about organic methods of growing? Ideas like ‘biodiversity’ and ‘sustainability’ come to mind.
  - Filtering DT-2 description column to ‘biodiversity’, only about 25% of the biodiversity posts speak about food, as observed in figure 21.

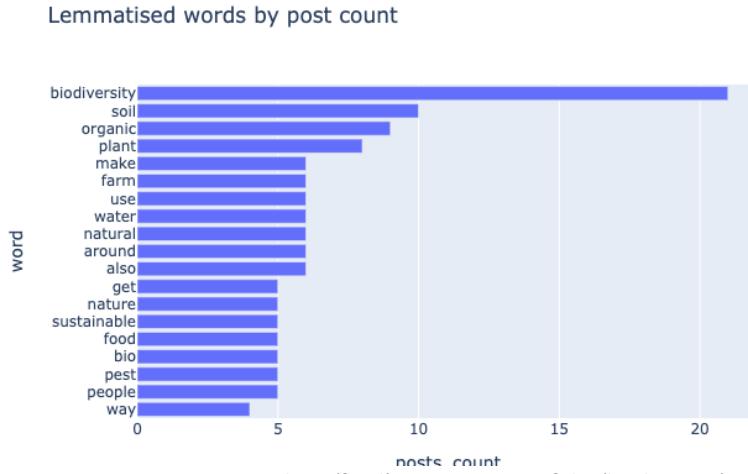


Figure 21: BC-1 where 'food' posts are 25% of the 'biodiversity' posts

- We can also see that terms related to farming are more prevalent in these posts: '**soil**', '**plant**', '**farm**'.
- It must also be noted that there are very few posts, only 21, about 'biodiversity'. So we must be cautious about drawing generalizations.
  - When we visualize the hashtags column for this same filtering, we see that they are mainly about **selling wines**.
  - This is a potential insight for organic food companies to hire biodiversity consultants who are working in the wine industry and the social media teams. They can enable a transformation in the meaning of 'organic food' not just being related to 'health' but also important for 'biodiversity' preservation.
- Filtering DT-2 to 'sustainable' and observing WC-2 and BR-1.
  - #sustainableliving and #sustainablefashion seem to be the notable themes.
  - When we visualize the hashtags specifically associated with 'food', we still see #sustainableliving to be in the top with other notable hashtags being #farmtotable, #sustainableliving, #smallbusinesses and #organicgardening.
  - From this we can say that the main idea of 'sustainable' within an organic and food context, is about a **way of life 'living'**.
- Filtering DT-2 to 'sustainability', we observe that #organicgardening is a prominent theme, which again suggests a way of living.
- Using the VE graph, we got a better understanding of what is meant by 'get' and 'make'
- In the VM graph we could clarify what meaning of 'make' is more prominent. And also we could get a list of what foods are most organic with terms like 'barley', 'hearty', 'raw', 'buckwheat'.
- To sum it up in a small paragraph:

"From the social media posts, we can infer that "Organic food" caters towards people interested in healthy lifestyle. The products generally claim to help the body in some way."

They encourage making changes to the consumed food towards healthier alternatives. The main foods which are organic include crops like ‘barley’ and ‘buckwheat’. Influencers and businesses presence is at least 16% with the posts containing a call to engage (‘get’, ‘u’)

## 5: Conclusion

The aim of this project was to create a tool to infer word meanings using mathematical methods. Although mathematical operations can be performed with great success on numerical tasks, language cannot be fully reduced to numbers. Extracting insights with these models should be considered as an art. The task of translating words into numbers and making use of computers to analyze them, has only seen rapid progress since the past decade, and for this, implementing neural network models has been a game changer in the field. Statistics is still at the base of these methods, and most of the recent developments in the state of the art have been achieved by finetuning the parameters.

A word assumes different meanings depending on how they are used by the people. A word like “organic” assumes different meanings depending on which context it is used in.

Although originally used in the context of farming, the sense of the word has assumed many connotations which intuitively make sense to us when used in those different contexts. In light of this, we had rather not relegate this intuitive sense making ability in the face of an empirical model, when it comes to language and word use.

Ideally, we should employ both an intuitive and empirical sense making approach. This tool attempts to offer such a functionality for the analyst or researcher. You may notice that the wordcloud immediately (perhaps “intuitively”) offers a deep insight into the nature of the posts without any numbers representing them. Whereas the bar charts add a numerical dimension to this which is relatively simple for most people to understand. The vector space graphs however, are a little more complex and may not be so intuitive for us to make sense of it right away. The results seem valid enough in most cases, but some results might not make sense, especially the ones with limited textual information.

In these vector space models, if the training dataset is a large corpus of text containing everyday use of language, then the results may make more sense because it has encoded most of the meaning behind it, in the way that we are used to using them. These results make intuitive sense to us only if we ourselves have been exposed to, and participate in, such everyday language use. I propose that this tool can uncover the evolving language use within a certain domain (example “organic”), which may not immediately feel intuitive, but over time gains its truth as such, as it is adopted by the society at large. Training the models on larger corpus of data will lead to better results.

The vector space graphs, might feel more useful for someone who has some (logical) understanding about the underlying statistics and the way it influences the results. The ideal candidate to use this tool could be a linguistics major with a statistical background and an interest in marketing.

It is worth noting that gaining insights from data depends on both, the methods, and the ability of the analyst to understand the methods. Hence, we have not considered more complex neural network models like transformers. This tool would be great for projects

involving social media market research and understanding the language use and motivations of the target users. In the coming months and years, we can expect to see further leaps in this field of abstracting information from words using computers, and it will become increasingly important for the data analysts to have some fundamental idea of the logic behind the methods in order to start using them, until they eventually feel intuitive.

## 6: Bibliography

- Baroni, M., & Lenci, A. (2010). Distributional Memory: A General Framework for Corpus-Based Semantics. *Computational Linguistics*, 36(4), 673–721. [https://doi.org/10.1162/coli\\_a\\_00016](https://doi.org/10.1162/coli_a_00016)
- Danilak, M. M. (n.d.). *langdetect: Language detection library ported from Google's language-detection*. (1.0.9) [Python; OS Independent]. Retrieved January 11, 2022, from <https://github.com/Mimino666/langdetect>
- Dumais, S. T. (2005). Latent semantic analysis. *Annual Review of Information Science and Technology*, 38(1), 188–230. <https://doi.org/10.1002aris.1440380105>
- Ganesan, K. (2014, October 19). What are Stop Words? *Kavita Ganesan, PhD*. <https://kavita-ganesan.com/what-are-stop-words/>
- Gensim: Topic modelling for humans*. (n.d.). Retrieved January 17, 2022, from <https://radimrehurek.com/gensim/models/word2vec.html>
- Gentner, D. (1983). Structure-Mapping: A Theoretical Framework for Analogy\*. *Cognitive Science*, 7(2), 155–170. [https://doi.org/10.1207/s15516709cog0702\\_3](https://doi.org/10.1207/s15516709cog0702_3)
- Harris, Z. S. (1954). Distributional Structure. *WORD*, 10(2–3), 146–162. <https://doi.org/10.1080/00437956.1954.11659520>
- Heidenreich, H. (2018, December 21). *Stemming? Lemmatization? What? Medium*. <https://towardsdatascience.com/stemming-lemmatization-what-ba782b7c0bd8>
- Introduction to Stemming. (2018, June 20). *GeeksforGeeks*. <https://www.geeksforgeeks.org/introduction-to-stemming/>
- Landauer, T. K., & Dumais, S. T. (1997). A solution to Plato's problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*, 104(2), 211–240. <https://doi.org/10.1037/0033-295X.104.2.211>
- Levy, O., Goldberg, Y., & Dagan, I. (2015). Improving Distributional Similarity with Lessons Learned from Word Embeddings. *Transactions of the Association for Computational Linguistics*, 3, 211–225. [https://doi.org/10.1162/tacl\\_a\\_00134](https://doi.org/10.1162/tacl_a_00134)

List of ISO 639-1 codes. (2022). In *Wikipedia*.

[https://en.wikipedia.org/w/index.php?title=List\\_of\\_ISO\\_639-1\\_codes&oldid=1064376433](https://en.wikipedia.org/w/index.php?title=List_of_ISO_639-1_codes&oldid=1064376433)

*Making sense of word2vec / RARE Technologies*. (n.d.). Retrieved January 16, 2022, from

<https://rare-technologies.com/making-sense-of-word2vec/>

Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient Estimation of Word Representations in Vector Space. *ArXiv:1301.3781 [Cs]*. <http://arxiv.org/abs/1301.3781>

*models.word2vec – Word2vec embeddings—Gensim*. (n.d.). Retrieved January 30, 2022, from

<https://radimrehurek.com/gensim/models/word2vec.html>

*On word embeddings - Part 3: The secret ingredients of word2vec*. (2016, September 24). Sebastian Ruder. <https://ruder.io/secret-word2vec/>

Pennington, J., Socher, R., & Manning, C. (2014). Glove: Global Vectors for Word Representation.

*Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1532–1543. <https://doi.org/10.3115/v1/D14-1162>

Python Simplified. (2020, October 10). *Web Scraping Instagram with Selenium*.

<https://www.youtube.com/watch?v=iJGvYBH9mcY>

Rivas, R., Sadah, S. A., Guo, Y., & Hristidis, V. (2020). Classification of Health-Related Social Media Posts: Evaluation of Post Content-Classifier Models and Analysis of User Demographics. *JMIR Public Health and Surveillance*, 6(2), e14952. <https://doi.org/10.2196/14952>

*scipy.sparse.coo\_matrix—SciPy v1.7.1 Manual*. (n.d.). Retrieved January 15, 2022, from

[https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.coo\\_matrix.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.coo_matrix.html)

Scrape Data Legally | No-Code Data API & Website Scraper | Stevesie Data. (n.d.). Retrieved January 9, 2022, from <https://stevesie.com/>

Selva Birunda, S., & Kanniga Devi, R. (2021). A Review on Word Embedding Techniques for Text Classification. In J. S. Raj, A. M. Iliyasu, R. Bestak, & Z. A. Baig (Eds.), *Innovative Data Communication Technologies and Application* (Vol. 59, pp. 267–281). Springer Singapore.

[https://doi.org/10.1007/978-981-15-9651-3\\_23](https://doi.org/10.1007/978-981-15-9651-3_23)

- Sklearn.decomposition.TruncatedSVD.* (n.d.). Scikit-Learn. Retrieved January 16, 2022, from  
<https://scikit-learn/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html>
- Sklearn.manifold.TSNE.* (n.d.). Scikit-Learn. Retrieved January 19, 2022, from <https://scikit-learn/stable/modules/generated/sklearn.manifold.TSNE.html>
- Smetanin, S. (2018, November 16). *Google News and Leo Tolstoy: Visualizing Word2Vec Word Embeddings with t-SNE*. Medium. <https://towardsdatascience.com/google-news-and-leo-tolstoy-visualizing-word2vec-word-embeddings-with-t-sne-11558d8bd4d>
- Turney, P. D., & Pantel, P. (2010). From Frequency to Meaning: Vector Space Models of Semantics. *Journal of Artificial Intelligence Research*, 37, 141–188. <https://doi.org/10.1613/jair.2934>
- WordCloud for Python documentation—Wordcloud 1.8.1 documentation.* (n.d.). Retrieved January 24, 2022, from [https://amueller.github.io/word\\_cloud/](https://amueller.github.io/word_cloud/)

## 7: Appendix

### 7.1 Interactive Dashboard Tool

This is how the tool appears when scrolled down on the web page.

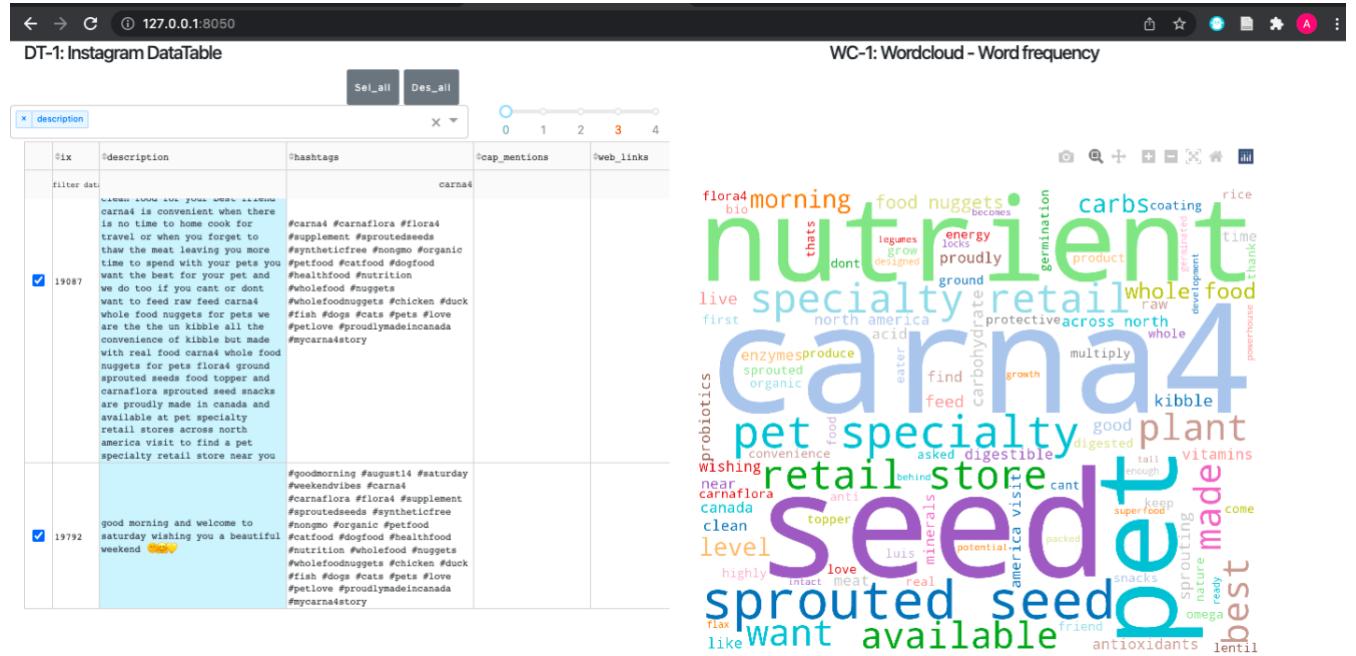


Figure 22: Dashtable DT-1, and Wordcloud WC-1 of the dashboard

DT-2: Z-Score DataTable

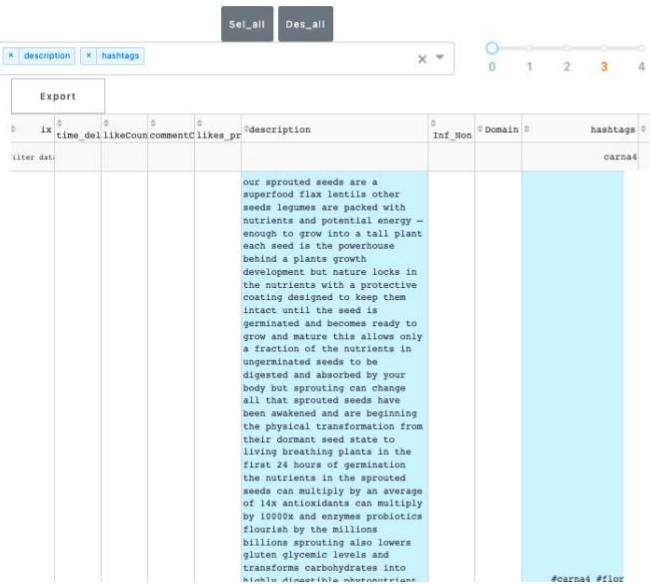


Figure 23: Dashtable DT-2, and Wordcloud WC-2 of the dashboard, on scrolling down

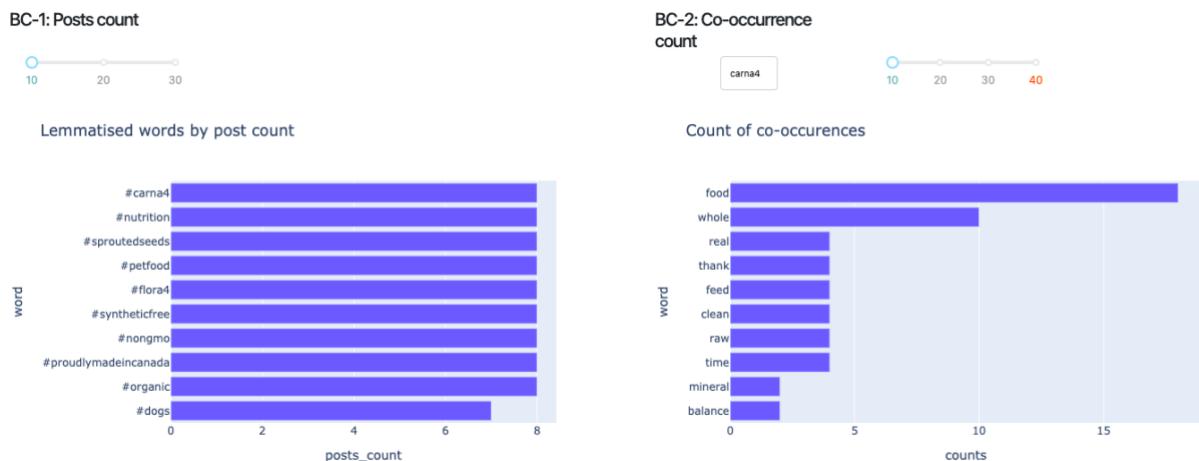


Figure 24: Barchart BC-1, and Barchart BC-2 of the dashboard, on scrolling down further



Figure 25: Vector Exploration VE graph, and Vector Math VM graph of the dashboard, towards the bottom end of the page

## 7.2 Development of code

SI No.	Development Version name	Date	Summary
1	Thesis_proj_v2.ipynb	20 July 2021	Exploring data sources. Data preprocessing.
2	Thesis_proj_v3.py	2 Sept 2021	App development started. Wordcloud.
3	Thesis_proj_v3.2.py	10 Sept 2021	Data table interactivity.
4	Thesis_proj_v3.3.py	16 Sept 2021	Stabilized app code. Code cleanup.
5	Thesis_proj_v3.4.py	27 Sept 2021	Word co-occurrence matrixes.
6	Thesis_proj_v3.5.py	26 Oct 2021	Word2vec math interactivity.
7	Thesis_proj_v4.py	10 Jan 2021	Final Explorations.

### 7.3 Python Code – preprocessing and creating support files

```

1 # In[1]:
2 print("Starting program..")
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 import unicodedata
9 from langdetect import detect
10
11 from wordcloud import WordCloud, STOPWORDS
12 from nltk.corpus import stopwords
13
14 import warnings
15 warnings.filterwarnings('ignore')
16
17 import os
18 import os.path
19
20 from datetime import datetime
21 app_launch_start=datetime.now() #Set start time for program start
22 # In[2]: Combine datasets; final_df.to_csv("combined_df.csv")
23 ##Collect multiple datasets saved in "datasets" folder and create 1 dataframe
24
25 print("\nReading datasets in the folder :")
26 start=datetime.now()
27 working_dir = "datasets" #folder where the datasets exist
28 all_files=[]
29 for root, dirs, files in os.walk(working_dir):
30     file_list = []
31     for filename in files:
32         if filename.endswith('.csv'):
33             # print(filename)
34             file_list.append(os.path.join(root, filename))
35     all_files+=file_list
36 print(all_files) #list of file names
37 print("Number of datasets considered:", len(all_files))
38
39
40 read_df_sh=[]
41 df_list=[]
42 for file in all_files:
43     read_df= pd.read_csv(file)
44     read_df_sh.append(read_df.shape)
45     df_list.append(read_df)
46 print("\nList of df shapes are: ",read_df_sh)
47
48 df_shape_sum=sum(i for i, j in read_df_sh)
49 print("\nSum of rows of multiple datasets", df_shape_sum)
50
51 #Check if the existing dataframe already includes all sub datasets
52 existing_comb_df=pd.read_csv('combined_df.csv')
53 existing_comb_df_rows=existing_comb_df.shape[0]
54 print("Sum of rows of existing combined dataset", existing_comb_df_rows)
55
56 if existing_comb_df_rows<df_shape_sum:
57     print("\nNew dataframe exists in the dataset folder..")
58     print("Create new combined dataframe")
59
60
61 if df_list:
62     final_df = pd.concat(df_list, ignore_index=True) |
63     final_df.to_csv("combined_df.csv")
64

```

```

55
56     else:
57         print("\nAll datasets already included in combined dataset")
58
59 t=datetime.now() - start #datetime object
60 s=str(t) #string object
61 print("Execution time ", s[:-5])
62
63 # In[3]:
64 ##Drop empty rows and create an SI number column
65 print("\nDrop empty rows and add SI number column..")
66 start=datetime.now()
67 comb_df=pd.read_csv('combined_df.csv')
68 print(comb_df.shape)
69 # comb_df.info() #Check where the empty rows are
70 comb_df = comb_df.dropna(axis=0, subset=['description']) #Drop empty rows in de
71 print("Redundant rows removed")
72 print(comb_df.shape)
73
74 #Hardcode index values
75 comb_df.reset_index(inplace= True,drop=True)
76 comb_df.reset_index(inplace= True)
77
78 comb_df.rename(columns={"index": "ix"},inplace=True) #to ref after dropping rows
79 comb_df['ix']=comb_df['ix'].astype(str) #convert to string type
80
81 t=datetime.now() - start
82 s=str(t)
83 print("Execution time ", s[:-5])
84
85 # In[4]:
86 ##Functions for splitting into new columns
87 print("\nLoading column splitting functions..")
88
89 def desc_cleaning(org_str): #some hashtags are not separated by a space.
90     orig_str=str(org_str)
91     new_str = orig_str.replace('#', "# ")
92     new_str2 = new_str.replace(' ', " ")
93     return new_str2
94
95 def desc_splitting(df):
96     df_upd=df
97     df_upd['clean_descriptions']=""
98     df_upd['hashtags']=""
99     df_upd['cap_mentions']=""
100    df_upd['web_links']=""
101
102    for ind, row in df_upd.iterrows():
103        X=str(row['new_desc']).split()#ensure description column is cleaned prior
104        X_cc=[]
105        X_hstgs=[]
106        X_cms=[]
107        X_http=[]
108
109        for x in X:
110            if not (x.startswith('#', '@', 'http')) or ('www.' in x) or ('.com' in x)) :
111                X_cc.append(x)
112
113            if x.startswith('#'):
114                X_hstgs.append(x)
115
116            if x.startswith('@'):
117                X_cms.append(x)
118
119            if (x.startswith('http') or ('www.' in x) or ('.com' in x)):
120
121
122
123
124
125
126
127
128
129

```

```

130         X_http.append(x)
131
132     df_upd.at[ind, 'clean_descriptions'] = ' '.join(X_cc)
133     df_upd.at[ind, 'hashtags'] = ' '.join(X_hstgs)
134     df_upd.at[ind, 'cap_mentions'] = ' '.join(X_cms)
135     df_upd.at[ind, 'web_links'] = ' '.join(X_http)
136
137     return df_upd
138
139
140 def new_cols(df): ##wrapping above functions
141     df['new_desc'] = df['description'].apply(lambda x : desc_cleaning(x))
142     df_updated=desc_splitting(df)
143     return df_updated
144
145 print("Loaded")
146
147
148 # In[5]:
149 ##New columns dataset "df_updated"
150 print("Applying splitting functions...")
151 start=datetime.now()
152
153 df_updated=new_cols(comb_df)
154 print("New columns created from splitting func")
155
156 t=datetime.now() - start
157 s=str(t)
158 print("Execution time ", s[:-5])
159
160
161 # In[6]: Text preprocessing; fonts, lower, punct, char_encodings
162 ##Tried to make granular transformations for specificity of processing.
163
164 print("\nLoading text preprocessing functions..")
165
166 def font_uniformity(x):
167     return unicodedata.normalize('NFKC', x)
168
169 def convert_lower_case(s):
170     return np.char.lower(s)
171
172 def remove_punctuation(data):
173     symbols = "#$%&()*+-./';<>?@{}^_{})~\n"
174     for i in range(len(symbols)):
175         data = np.char.replace(data, symbols[i], ' ')
176
177     data = np.char.replace(data, " ", " ")
178     data = np.char.replace(data, ',', ' ')
179     return data
180
181 def diff_encodings(s):
182     s = np.char.replace(s, "„", " ")
183     s = np.char.replace(s, "„", " ")
184     return s
185
186 print("Loaded")
187 print("Applying text preprocessing...")
188 start=datetime.now()
189
190 df_updated['caption_processed']=df_updated['clean_descriptions'].apply(lambda x: font_uniformity(x))
191 df_updated['caption_processed_2']=df_updated['caption_processed'].apply(lambda x: convert_lower_ca
192 df_updated['caption_processed_3']=df_updated['caption_processed_2'].apply(lambda x: remove_punctua
193 df_updated['caption_processed_4']=df_updated['caption_processed_3'].apply(lambda x: diff_encodings
194
195 print("Text preprocessing done")
196 t=datetime.now() - start
197 s=str(t)
198 print("Execution time ", s[:-5])
199
200 # In[8]: "views" and "exemplars" - testing
201 ##Define view frame to view previous processing steps with exemplars
202 view1=['description','new_desc','clean_descriptions' ]
203 view2=[]
204 view2=['clean_descriptions','caption_processed','caption_processed_2', 'caption_processed_3',
205         'caption_processed_4']
206
207 #define the rows to display-exemplars of diff test cases
208 ex_row_list=[11,13,23,106] #subtract 2 because index is reset.
209
210 #df_updated[view2].loc[ex_row_list]
211 #df_updated[view2].head(30)
212
213 # In[9]: Language detection to_csv("App_dataframe.csv")
214 # Visualise how many languages are there:
215 print("\nLoading language detection function..") #Takes 11:30 mins to execute

```

```

216 start=datetime.now()
217
218 def lang_det(st):
219     try:
220         lang=detect(st)
221         return lang
222     except:
223         lang="error"
224         return lang
225
226 view3=[ 'det_lang']
227 view3= view2+view3
228
229 print("Running language detection...")
230 df_updated[ 'det_lang']= df_updated[ 'caption_processed_2'].apply(lambda x: lang_det(x))
231 print("Language detection complete")
232 t=datetime.now() - start #datetime object
233 s=str(t) #string object
234 print("Execution time ", s[:-5])
235
236 #Export to app referenced df after this last transformation:
237 df_updated.to_csv("App_dataframe.csv")
238
239 # In[9]: Language countplot
240 ##visualise the new transformations
241 df_updated = pd.read_csv('App_dataframe.csv')
242
243 print("Making countplot..")
244 plt.figure(figsize=(16,6))
245 ax= sns.countplot(x= 'det_lang', data=df_updated, order = df_updated[ 'det_lang'].value_counts(
246     ascending=False).index)
247 ax.set_title('Language distribution')
248 ax.set_xlabel('Languages')
249
250 for tick in ax.get_xticklabels():
251     tick.set_rotation(90)
252
253 print("Language countplot graph loaded")
254
255 # In[10]: #Selective viewing of dataframe with "views" and "exemplars"
256 view_extracts=[ 'new_desc', 'det_lang', 'clean_captions', 'caption_processed_4', 'hashtags',
257 'cap_mentions', 'web_links' ]
258
259 df_updated[view_extracts].loc[ex_row_list]
260 df_updated[view_extracts].head(100)
261
262 #add new elements to exemplar list
263 ex_row_list.append(3) #add exemplar accents in differnt language
264 df_updated[view3][df_updated[ 'det_lang']=='error'].head(10)
265 #Errors for descriptions which do not have readable text. Either blank or emojis.
266 #Fonts normalised but not much improvement.
267
268 # In[11]: ##Wordcloud with English and French stopwords
269
270 print("\nLoading test Word Cloud..")
271 start=datetime.now()
272
273 # df_updated['clean_captions'][10] #test individual rows
274 word_string = " ".join(df_updated[ 'clean_captions'].dropna().astype(str)) #select column series
275 print(word_string[10:1000])
276
277
278 #Define stopwords
279 # type(STOPWORDS) #set
280 # len(STOPWORDS) #192
281 en_stopwords = stopwords.words('english') #179
282 fr_stopwords = stopwords.words('french') #157
283 web_links_sw = [ 'www', 'http', 'https', 'com' ]
284
285 combined_stopwords = en_stopwords + fr_stopwords
286
287 wordcloud = WordCloud(max_words=100,
288                         stopwords= combined_stopwords,
289                         collocations=False,
290                         color_func=lambda *args, **kwargs: "orange",
291                         background_color='white',
292                         width=1200,
293                         height=1000).generate(word_string)
294
295 fig_dims = (8, 8) #size of plot
296 fig, ax = plt.subplots(figsize=fig_dims)
297 plt.title("#Test word cloud")
298 plt.imshow(wordcloud, interpolation='bilinear')
299 plt.axis('off')
300 plt.show()
301 #wordcloud.to_file("static/images/en_fr_sw_wc.png")

```

```

302 print("Test Word Cloud loaded")
303 t=datetime.now() - start
304 s=str(t)
305 print("Execution time ", s[:-5])
306
307 #With only EN stopwords: "u", "de", "la", "e", still present
308 #With EN and FR : 'u', 'e', 'o' still present
309
310 # In[12] Filter; read_csv("App_dataframe.csv") -> to_csv("App_dataframe_2.csv")
311 #Filter the dataframe to English, #organic, and remove duplicate posts
312 ##Export to a new dataframe
313 print("Reading in dataframe for app with filters and duplicate removal..")
314 ad_1=pd.read_csv("App_dataframe.csv")
315 ad_1.shape
316 ad_2=ad_1[ad_1['det_lang']=='en'] #filter only english, removes appr 33%
317 print("Only English posts",len(ad_2))
318
319 ad_3=ad_2[ad_2['query']=="#organic"] #filter only #organic scrape query, removes 33%
320 print("Only organic posts: ", len(ad_3))
321
322 ad_4=ad_3.drop_duplicates(subset=['postId']) #remove dupl posts, removes 12.5%
323 print("After duplicate removal: ", len(ad_4))
324
325 ad_4.drop(['Unnamed: 0','ix','Unnamed: 0.1'],axis =1,inplace=True)
326
327 #Create new index for proper selection of the filtered table in appl
328 ad_4.reset_index(inplace= True,drop=True)
329 ad_4.reset_index(inplace= True)
330 ad_4.rename(columns={"Index": "ix"},inplace=True) #replace index to use as row ref
331 ad_4['ix']=ad_4['ix'].astype(str) #convert to string type
332
333 ad_4.to_csv("App_dataframe_2.csv", index_label = False)
334
335
336 # In[13] More preprocessing; read_csv("App_dataframe_2.csv") --> to_pickle("App_dataframe_3.pkl")
337 # Implement Vecotrization.
338 # Tokenize the processed description
339
340 import ast
341 import preprocessor as p
342
343 ad_4= pd.read_csv("App_dataframe_2.csv")
344
345 #define function to remove newly identified punctuations/encodings
346 def addi_treatment(s): #Move this to the main preprocessing block if required
347     s = np.char.replace(s, "", "")
348     s = np.char.replace(s, "", "")
349     s = np.char.replace(s, "", "")
350     s = np.char.replace(s, "-", "") #from viewing bag of wrods in first 100 of cp4
351     return s
352
353
354 def preprocess_tweet(row):
355     text = p.clean(row) # removes hashtags,emojis,urls, mentions, smileys
356     return text
357
358 def remove_stop_words(data): #takes 5:25 mins for the full dataset
359     """
360         Split/tokenize the words first. Create a list of tokens for each.
361         Iterate each token removing stop words from text.
362     """
363     row_tokens = data.split()
364     return [word for word in row_tokens if not word in stopwords.words("english")] # 25 s for 100r
365
366 start= datetime.now()
367
368 #Create new columns
369 ad_4['caption_processed_4'] = ad_4['caption_processed_4'].apply(lambda x : addi_treatment(x))
370 ad_4['cp4_no_emojis']=ad_4['caption_processed_4'].apply(lambda x: preprocess_tweet(x))
371 ad_4['cp4_rm_sw']= ad_4['cp4_no_emojis'].apply(remove_stop_words)
372 ad_4['cp4_rm_sw']= ad_4['cp4_rm_sw'].astype(str) #convert to string for literal eval
373 ad_4['cp4_rm_sw']= ad_4['cp4_rm_sw'].apply(lambda x : ast.literal_eval(x)) #preserve lists.
374 al=datetime.now() - start #start time logged at start
375 sal=str(al) #string object
376 print("Execution time for all rows : ", sal[:-5])
377
378 print("Type of rows in cp4_rm_sw are ", type(ad_4['cp4_rm_sw'][32000]))
379 ch_list=['caption_processed_4', 'cp4_rm_sw']
380 view=ad_4[['caption_processed_4', 'cp4_rm_sw']] #index 35129
381 # print(view)
382
383 ad_4['cp4_rm_sw'][32000]
384
385 ad_4.to_pickle("App_dataframe_3.pkl", protocol=0) #preserves the list datatypes. prot=0 for colab
386
387 ##App_dataframe_3 now has

```

```

388 #words minus stopwords
389 #treated for "addi_treatment"
390 #tweet preprocessor removed emojis
391
392
393 # In[14]# Class for cooccurrence embeddings initial- from prof Redha Moulla notes
394 # Finally used different co-occ method with window selection.
395 from sklearn.decomposition import TruncatedSVD
396
397 class CooccEmbedding:
398     def __init__(self, corpus):
399         """
400             Takes in the corpus in the form of a list of lists (list of tweets,
401             each tweet being a list of tokens)
402         """
403         self.corpus = corpus
404
405     def vocabulary(self):
406         """
407             Returns the vocabulary associated with the corpus (a list of unique tokens from the corpus
408         """
409         self.vocab = []
410         [self.vocab.append(word) for tweet in self.corpus for word in tweet if word not in self.vocab]
411         self.len_vocab = len(self.vocab) #Initialise this return because coocc module not used.
412         return self.vocab #outputs single list of unique words
413
414     def coocc_matrix(self): #Note that this is not looking at window size for cooccurrence. Consider
415         """
416             Returns the co-occurrence matrix associated with the corpus.
417             The co-occurrence matrix is first calculated as a list using comprehensive lists to speed
418             the iteration process over the vocabulary and the corpus.
419         """
420         self.len_vocab = len(self.vocab)
421         coocc_list = [sum([1 if (self.vocab[i] in tweet and self.vocab[j] in tweet and i != j) else
422                           0 for i in range(self.len_vocab) for j in range(self.len_vocab)]) #1st word at
423         coocc_list = np.array(coocc_list)
424         # transform the coocc_list into a matrix
425         self.coocc_mat = coocc_list.reshape([self.len_vocab, self.len_vocab])
426         return self.coocc_mat
427
428     def svd_reduction(self, n_components = 2, n_iter = 10):
429         """
430             Performs a singular value decomposition on the co-occurrence matrix M and returns the truncated
431             matrix U, where M = UDV^T.
432         """
433         svd = TruncatedSVD(n_components = n_components, n_iter = n_iter)
434         self.reduced_matrix = svd.fit_transform(self.coocc_mat)
435         return self.reduced_matrix
436
437     def vocab_ind_to_plot(self, vocab_to_plot):
438         """
439             Takes in the a list of vocabulary to plot and returns a dictionary in the form of
440             {word:index} which is a subset of the vocabulary.
441         """
442         vocab_dict = dict(zip(self.vocab, range(self.len_vocab)))
443         self.dic_vocab_to_plot = {key:value for key, value in vocab_dict.items() if key in vocab_to_plot}
444         return self.dic_vocab_to_plot
445
446
447 # In[15]: #Test bag of words sizes and instantiate Coocc_EMBEDDING class
448 # Trim the vocab so that cooccurrence matrix can be computed.
449
450 ad_5=pd.read_pickle('App_dataframe_3.pkl')
451
452 def list_of_words(data): #takes 5:25 mins for the full dataset
453     row_tokens = data.split()
454     return [word for word in row_tokens] # 25 s for 100return the word only if not in the stop words
455
456 #Prepare the list of lists for the bag of words input in Coocc class:
457 ad_5['caption_processed_4'] = ad_5['caption_processed_4'].apply(list_of_words)
458 words_cp4=[row_list for row_list in ad_5['caption_processed_4']]
459 words_cp4[0:3]
460
461 ad_5['cp4_no_emojis'] = ad_5['cp4_no_emojis'].apply(list_of_words)
462 words_cp4_no_emojis=[row_list for row_list in ad_5['cp4_no_emojis']]
463
464 words_cp4_rm_sw=[row_list for row_list in ad_5['cp4_rm_sw']] # List of lists
465
466 ##Tests on the first 100 rows
467 inst = CooccEmbedding(words_cp4[:100])
468 print("Number of unique words: ",len(inst.vocabulary())) #1587 words in cp4
469
470 inst = CooccEmbedding(words_cp4_no_emojis[:100])
471 print("Number of unique words: ",len(inst.vocabulary())) #1418 words in words_cp4_no_emojis
472
473 inst = CooccEmbedding(words_cp4_rm_sw[:100])

```

```

474 print("Number of unique words: ",len(inst.vocabulary())) #1309 words in words_cp4_rm_sw
475 |
476 #For full dataframe:
477 inst_1 = CooccEmbedding(words_cp4)
478 inst_2 = CooccEmbedding(words_cp4_no_emojis)
479 inst_3 = CooccEmbedding(words_cp4_rm_sw)
480
481 #Takes time so time it
482 time_vocab=datetime.now()
483
484 print("cp4 unique words: ",len(inst_1.vocabulary())) #71616 words,#46066 words,#45925 words
485 t=datetime.now() - time_vocab
486 s=str(t)
487 print("cp4 counting time: ", s[:-5])
488
489 print("cp4_no_emojis unique words: ",len(inst_2.vocabulary())) #46066 words
490 t=datetime.now() - time_vocab
491 s=str(t)
492 print("cp4_no_emojis counting time: ", s[:-5])
493
494 print("cp4_rm_sw unique words: ",len(inst_3.vocabulary())) #45925 words
495 t=datetime.now() - time_vocab
496 s=str(t)
497 print("cp4_rm_sw counting time: ", s[:-5])
498
499
500 # In[16]: #Lemmatization and Stemming ->create new pkl App_dataframe_4
501 ad_5=pd.read_pickle('App_dataframe_3.pkl')
502
503 #Porter Stemming
504 from nltk.stem import PorterStemmer
505 porter = PorterStemmer()
506
507 def porter_stemming(row_list):
508     new_word_list=[porter.stem(word) for word in row_list]
509     return new_word_list
510
511 start = datetime.now()
512 ad_5['rm_sw_stem']=ad_5['cp4_rm_sw'].apply(porter_stemming)
513 t=str(datetime.now()-start)
514 print("Time for stemming full dataframe: ", t[:-5]) #27.2 seconds for full dataframe
515
516 words_rm_sw_stem=[row_list for row_list in ad_5['rm_sw_stem']]
517
518 print(words_cp4_rm_sw[:2])
519 print("\n",words_rm_sw_stem[:2])
520
521 inst_4=CooccEmbedding(words_rm_sw_stem)
522
523 start=datetime.now()
524 print("rm_sw_stem unique words: ",len(inst_4.vocabulary())) #34051 unique words
525 t=str(datetime.now() - start)
526 print("rm_sw_stem counting time: ", t[:-5]) #1:26 mins
527
528
529 #Lemmatization
530 from nltk.stem import WordNetLemmatizer
531 wordnet_lemmatizer = WordNetLemmatizer()
532
533 def wordnet_lemmatizer_func(row_list):
534     new_word_list=[wordnet_lemmatizer.lemmatize(word) for word in row_list]
535     return new_word_list
536
537
538 ad_5['cp4_rm_sw'][3:5]
539
540 #For testing the functions
541 test_series= ad_5['cp4_rm_sw'][3:5].apply(wordnet_lemmatizer_func)
542 print(ad_5['cp4_rm_sw'][3:5])
543 print("\n Stemmed",ad_5['rm_sw_stem'][3:5])
544 print("\nLemmatized",test_series)
545
546 ##Repeat as in previous step for stemming
547
548 start = datetime.now()
549 ad_5['rm_sw_lemt']=ad_5['cp4_rm_sw'].apply(wordnet_lemmatizer_func)
550 t=str(datetime.now()-start)
551 print("Time for lemmatizing full dataframe: ", t[:-5]) #5.5 seconds for full dataframe
552
553 words_rm_sw_lemt=[row_list for row_list in ad_5['rm_sw_lemt']]
554
555 inst_5=CooccEmbedding(words_rm_sw_lemt)
556
557 start=datetime.now()
558 print("rm_sw_lemt unique words: ",len(inst_5.vocabulary())) #41526 unique words
559 t=str(datetime.now() - start)

```

```

560 print("rm_sw_lemt counting time: ", t[:-5]) #1:48 mins
561
562 ad_5.to_pickle("App_dataframe_4.pkl", protocol=0) #with 2 new columns stemming and lemmatization
563
564 # In[17]: #More filtering of the dataframe App_dataframe_4->App_dataframe_5
565 ##Filter out duplicate posts and create another pickle App_dataframe_5
566
567 ad_6=pd.read_pickle('App_dataframe_4.pkl')
568 len(ad_6) # 35131
569
570 ad_7 = ad_6.drop_duplicates(subset=['caption_processed_4'])
571 len(ad_7) # from 35131 of ad_6, ad_7 is reduced to 30240; Another 14.2% drop in the rows.
572
573 ad_7.drop(['ix'], axis=1, inplace = True) #drop any older 'ix' columns
574 ad_7.reset_index(inplace= True,drop=True) #reset and drop old index
575 ad_7.reset_index(inplace= True) #create new ix column
576 ad_7.rename(columns={"index": "ix"},inplace=True)
577
578 ad_7.to_pickle("App_dataframe_5.pkl", protocol=0)
579
580 len(ad_7) #29681
581
582 # In[18]: #Window Cooccurrence array; save('co_occ_arr.npy'); choose 'rm_sw_lemt'
583 # Takes a long time. Run this only for new datasets.
584
585 from collections import defaultdict
586
587 ad_6=pd.read_pickle('App_dataframe_5.pkl') #duplicate posts removed.
588
589 words_rm_sw_lemt=[row_list for row_list in ad_6['rm_sw_lemt']] #row_list represents each words in
590 print("Bag of words created..")
591 len(words_rm_sw_lemt)
592
593 text=words_rm_sw_lemt #full dataset
594
595 def co_occ_windows(sentences, window_size):
596     d = defaultdict(int) #gives default value for non-existing keys
597     vocab = set()
598     for text in sentences: #text is each post within full bag
599         print(text)
600         for i in range(len(text)):
601             token = text[i]
602             vocab.add(token) #add to vocab
603
604             # coocc_window = text[i+1 : i+1+window_size] #Rolling ahead window scope; sparse
605             coocc_window = text[i-window_size : i] + text[i+1 : i+1+window_size] #Rolling bidirec
606
607             for t in coocc_window:
608                 key = tuple( sorted([t, token]) )
609                 # print("key is: ",key)
610                 d[key] += 1 #at the tuple key, increase the value by 1
611                 # print("default dict value at key is: ",d[key]) #Each key is a tuple and is uniqu
612
613             # formulate the dictionary into dataframe
614             vocab = sorted(vocab) # sort vocab
615             global custom_vocab
616             custom_vocab = vocab #test object
617
618             # print("Sorted vocab at index is", vocab[0])
619             df = pd.DataFrame(data=np.zeros((len(vocab), len(vocab))), dtype=np.int16),
620                               index=vocab,
621                               columns=vocab) #Initialise a dataframe of zeroes
622             for key, value in d.items():
623                 df.at[key[0], key[1]] = value
624                 df.at[key[1], key[0]] = value
625             return df
626
627 #Test out the function
628 start=datetime.now()
629 # co_occ_df = co_occ_windows(text,5) #full array w4 takes 24:25 mins rolling ahead
630 co_occ_df = co_occ_windows(text,2) #bidirectional. Effective window size is 5 takes mins
631 print("shape of co_occ matrix is: ",co_occ_df.shape)
632 t=str(datetime.now() - start)
633 print("Time for execution: ", t[:-5])
634
635 print("Custom vocab shape: ", custom_vocab)
636 print("First 10 vocab shape: ", custom_vocab[:10])
637
638 co_occ_arr =co_occ_df.to_numpy() #convert to an array
639 co_occ_arr
640 from numpy import save #Have to import explicitly to save array as binary
641 save('ad5_co_occ_arr_w5.bi.npy', co_occ_arr)
642
643 # In[19]Perform Singular Value Decomposition on the array. SVD_matrix.npy
644
645 from numpy import load
646 # co_occ_arr = load('co_occ_arr_w2.npy') #41526

```

```

646 co_occ_arr = load('ad5_co_occ_arr_w5_bi.npy') #shape 41526
647
648 print("shape of co-occ matrix is: ", co_occ_arr.shape)
649
650 ad_6=pd.read_pickle('App_dataframe_5.pkl') #duplicate posts removed.
651
652 from sklearn.decomposition import TruncatedSVD
653
654 # co_occ_arr_sl=co_occ_arr[0:10000,0:10000] #sliced array
655 co_occ_arr_sl=co_occ_arr #Full cooccurrence matrix
656
657 start = datetime.now()
658
659 from scipy.sparse import coo_matrix
660 co_occ_arr_sl_coo=coo_matrix(co_occ_arr_sl)
661 co_occ_arr_sl_coo = co_occ_arr_sl_coo.asfptype() #convert to COO format before arpack
662
663 svd = TruncatedSVD(n_components = 2, n_iter = 10, algorithm = "arpack") #28.4 secs
664 # svd = TruncatedSVD(n_components = 2, n_iter = 10, algorithm = "randomized") #18.3 secs
665 Coocc_svd_matrix = svd.fit_transform(co_occ_arr_sl_coo)
666 Coocc_svd_matrix.shape
667
668 t = str(datetime.now()-start)
669 print("Time taken for svd: ",t[:-5])
670
671 print(type(Coocc_svd_matrix))
672
673 from numpy import save
674 save('ad5_svd_arpack_w5_bi.npy', Coocc_svd_matrix) #arpack algo used
675
676 # In[20]: svd array->plotting dataframe
677
678 coocc_svd_matrix = load('ad5_svd_arpack_w5_bi.npy') #Load arpack mode
679
680 ad_6=pd.read_pickle('App_dataframe_5.pkl') #df_8hrs is converted to string for col combining in ap
681
682 print("Custom vocab head 10: ",custom_vocab[:10])
683 custom_vocab[-10:-1]
684 cust_vocab = custom_vocab #entire word vocab. Common for all window sizes
685 # len(vocab_full) #41526
686 certs = ['gots','oeko','cosmos','ecocert','fsc'] #certifications
687 fruits = ['apple', 'orange']
688 rel_concs = ['fruits', 'certification','textile']
689
690 # vocab_to_plot = certs + fruits + rel_concs #specific plot list
691 vocab_to_plot = cust_vocab #full vocab plot
692
693 print("Making dict_to_plot..")
694 start=datetime.now()
695
696 #custom_vocab and co_occ vocab are indexed differently. coocc mat is based on custom(sorted) vocab
697 # dict_to_plot = inst.vocab_ind_to_plot(vocab_to_plot) #with above finding instantiating coocc cla
698
699 #dict_to_plot for the custom_vocab
700 def vocab_ind_to_plot(cust_vocab, vocab_to_plot):
701     """
702         Takes in a list of full vocabulary and list to plot and returns a dictionary in the fo
703         {word:index} which is a subset of the vocabulary.
704     """
705     vocab_dict = dict(zip(cust_vocab, range(len(cust_vocab)))) #custom vocab is sorted
706     dic_to_plot = {key:value for key, value in vocab_dict.items() if key in vocab_to_plot} #Re
707     return dic_to_plot
708
709 dict_to_plot = vocab_ind_to_plot(cust_vocab, cust_vocab)
710
711 print("Dict_to_plot is ",dict_to_plot)
712
713 wl1=str(datetime.now()-start)
714 print("Time to make dict_to_plot: ", wl1[:-5])
715
716 data_list=[]
717 for word, ind in dict_to_plot.items():
718     row_list=[word, coocc_svd_matrix[ind, 0], coocc_svd_matrix[ind, 1]]
719     data_list.append(row_list)
720
721 # print(data_list)
722 type(coocc_svd_matrix)
723 print(coocc_svd_matrix.shape)
724
725 vocab_words_df= pd.DataFrame.from_records(data_list, columns=['word', 'x', 'y'])
726 wl2=str(datetime.now()-start)
727 print("Time till vocab_words_df ready: ", wl2[:-5]) #55secs for full vocab 41526 words
728 vocab_words_df
729 vocab_words_df.dtypes
730
731 vocab_words_df.to_pickle("ad5_vocab_words_svd_w5_bi.pkl", protocol=0)

```

```

732 # print(vocab_words_df)
733
734 # In[21]: Visualise plotly scatter - Coocc SVD vectors
735 import plotly.express as px
736 import plotly.io as pio
737 pio.renderers.default='browser' #To plot in browser
738
739 vocab_words_df_full=pd.read_pickle('ad5_vocab_words_svd_w5_bi.pkl') #corrected word index
740
741 certs = ['gots', 'oeko', 'cosmos', 'ecocert', 'fsc'] #certifications
742 fruits = ['apple', 'orange', 'vanilla', 'cacao']
743 concs = ['sustainably', 'ethically']
744 rel_concs = ['fruits', 'certification', 'textile', 'cosmetics']
745 vocab_to_plot = certs + fruits + rel_concs
746
747 #filtered plot df
748 vocab_words_df = vocab_words_df_full[vocab_words_df_full['word'].isin(vocab_to_plot)]
749
750 fig = px.scatter(vocab_words_df_full, x="x", y="y", text="word", log_x=False, size_max=60,
751                   # trendline = "rolling",
752                   hover_name = 'word'
753                   )
754 fig.update_traces(textposition='top center')
755 fig.show()
756
757
758 # In[22]: ## Word2vec model and saving .txt files
759 import os
760 import re
761 import time
762 from gensim.models import Word2Vec
763 from gensim.models import KeyedVectors #For loading the model
764 from sklearn.manifold import TSNE #for reducing the dimensions
765 from tqdm import tqdm
766 tqdm.pandas()
767
768 ad_7 = pd.read_pickle('App_dataframe_5.pkl')
769 words_rm_sw_lemt = [row_list for row_list in ad_7['rm_sw_lemt']]
770 print("number of sentences is:",len(words_rm_sw_lemt)) #41526
771 train_sentences = words_rm_sw_lemt #list of lists sentences and words
772
773 start=datetime.now()
774 model = Word2Vec(sentences=train_sentences, sg=1,vector_size=300, window=5, min_count=1, workers=4
775 model.save("thesis_word2vec.model") #save a part of the training. avoid saving, is large
776 model = Word2Vec.load("thesis_word2vec.model") #load and continue training the sample
777 model.train([["hello", "world"]], total_examples=1, epochs=1) #train new samples if required
778 model.wv.most_similar('organic', topn = 20)
779 model.wv.save_word2vec_format('organic_glove_300d_w5.txt') #save all embeddings in this format
780
781 print("Time taken for creating model: ", str(datetime.now()-start)[:5])#44 secs for 300d_w4
782
783 # In[23]: Create and Export Tsne dataframe by transforming w2c text files
784 # Note how before the saving the model, called as model.wv
785 # Takes 45-60 mins for full dataset.
786
787 w2v = KeyedVectors.load_word2vec_format('organic_glove_300d_w5.txt') # To load text file
788 # w2v = model.wv #use this if not loading from file and model is instantiated
789
790 def tsne_df(w2v): # dataframe conversion function
791     """Creates and TSNE model and plots it"""
792     labels = []
793     tokens = []
794
795     words = list(w2v.index_to_key) #instead of model.wv.vocab
796     len(words)
797     words_cut=words[:1000] #To test time on a smaller df
798
799     for word in words: #Test with 'words_cut' first
800         tokens.append(w2v[word]) #Fetches the vector for the word. same as model.wv
801         # print(model.wv[word]) |
802         labels.append(word) #Word
803         # print(word)
804
805     tsne_model = TSNE(perplexity=40, n_components=2, init='pca', n_iter=2500, random_state=23)
806     new_values = tsne_model.fit_transform(tokens) #fit_transform on vectors
807
808     x = []
809     y = []
810     for value in new_values:
811         x.append(value[0])
812         y.append(value[1])
813
814     tsne_df = pd.DataFrame(
815         {'word': labels,
816          'x': x,
817          'y': y}
818

```

```

818     })
819
820     tsne_df.to_pickle("tsne_300d_w2_df.pkl", protocol=0)
821     print("Time taken for TSNE: ", str(datetime.now()-start)[-5])
822
823     print("Creating new TSNE dataframe..")
824     start=datetime.now()
825     tsne_df(w2v) #call the function
826
827 # In[24]: Addition and subtraction logic word2vec
828 w2v = KeyedVectors.load_word2vec_format('organic_glove_300d_w5.txt') # To load text file
829 w2v.most_similar(positive=['king', 'woman'], negative= ['man'])
830 print(w2v.most_similar(positive=['oil', 'hempseed'])) #addition
831 print(w2v.most_similar(['oil', 'hempseed', 'free'])) #same as above
832
833 # In[25]: App simulation - w2v Scatter plot logic including hue.
834 # Toggle text display on and off.
835 # Bold word name in hover text
836 import plotly.express as px
837 import plotly.io as pio #To plot in browser
838 pio.renderers.default='browser'
839
840 vocab_to_plot = ['vanilla', 'organic', 'sustainable', 'cacao']
841 tsne_df_full=pd.read_pickle('tsne_100d_w5_df.pkl')
842 w2v = KeyedVectors.load_word2vec_format('organic_glove_300d.txt')
843 nearest_size = 5 #make input for nearest neighbor size
844
845 #index of input words
846 tsne_df_full['type'] = 'Sel'
847 ip_vocab_index=list(tsne_df_full[tsne_df_full['word'].isin(vocab_to_plot)].index.values) #
848
849 #index of 10 words closest to inputs
850 clos_ten_out = []
851 for word in vocab_to_plot:
852     a=w2v.most_similar(word, topn = nearest_size)
853     clos_ten_in = [i[0] for i in a]
854     clos_ten_out.append(clos_ten_in)
855
856 clos_ten_flat = [item for sublist in clos_ten_out for item in sublist]
857 clos_ten_index = list(tsne_df_full[tsne_df_full['word'].isin(clos_ten_flat)].index.values) #
858
859 plot_index = ip_vocab_index + clos_ten_index
860 indices = [0,1,3,6,10,15]
861
862 indices = [0,1,3,6,10,15]
863 tsne_df_full.loc[ip_vocab_index, 'type'] = 'ip'
864 tsne_df_full.loc[clos_ten_index, 'type'] = 'near'
865
866 tsne_plot_df = tsne_df_full.loc[plot_index] #filter by index
867 tsne_plot_df
868 # tsne_plot_df = tsne_df_full[tsne_df_full['word'].isin(vocab_to_plot)] #filter by word
869
870 fig3 = px.scatter(tsne_plot_df, x="x", y="y", text="word", color='type', log_x=False, size_max=60,
871                     # trendline = "rolling",
872                     # hover_name = 'word'
873                     )
874 fig3.update_traces(textposition='top center')
875 fig3.show()
876
877 # In[26]: co-occurrence visualisations
878 #instantiate the sorted vocab in the co-occ code block
879 from numpy import load
880 co_occ_arr = load('ad5_co_occ_arr_w5.bi.npy') #shape 41526 * 41526
881 clf_table = pd.read_pickle("App_dataframe_5.pkl") #list col rm_sw_lemt
882
883 flat_list = [item for row_list in clf_table['rm_sw_lemt'] for item in row_list] #list col
884 len(flat_list) #total words in the corpus 799068 (df_8hr), 813922(App_dataframe_5)
885 len(set(flat_list)) #41526 (AD_5)
886
887 vocab = sorted(set(flat_list))
888
889 co_occ_df = pd.DataFrame(data = co_occ_arr,
890                           index = vocab,
891                           columns = vocab)
892 # co_occ_df.to_pickle("co_occ_df.pkl", protocol=0) #hangs the kernel. do not use!
893
894 #datastructure for bar chart
895 type(co_occ_df['organic'])
896 word = 'textile' #test
897 co_occ_plot = co_occ_df[word].sort_values(ascending=False).head(30)
898 co_occ_plot = co_occ_plot.reset_index() #make 'words' a column
899 mapping = {co_occ_plot.columns[0]: 'word', co_occ_plot.columns[1]: 'counts'}
900 co_occ_plot = co_occ_plot.rename(columns=mapping)
901 co_occ_plot

```

## 7.4 Python Code – Plotly dash web app

```

901 # In[27]: #App code block. To be run as a script. Can host on heroku app server.
902 print("\nStarting dashboard app...")
903 #####Import libraries and dataset
904 import pandas as pd      #(version 1.0.0)
905 import plotly            #(version 4.5.4) pip install plotly==4.5.4
906 import plotly.express as px
907 import dash               #(version 1.9.1) pip install dash==1.9.1
908 import dash_table          #interactive datatable
909 import dash_core_components as dcc #display elements
910 import dash_html_components as html #display elements
911 from dash.dependencies import Input, Output, State #for app callbacks
912 import dash_bootstrap_components as dbc #display elements
913 import pandas as pd
914 import matplotlib.pyplot as plt
915 import seaborn as sns
916 from wordcloud import WordCloud, STOPWORDS
917 from nltk.corpus import stopwords
918 from datetime import datetime
919 from numpy import load
920 import os
921 import re
922 import time
923 from gensim.models import Word2Vec #for word2vec model
924 from gensim.models import KeyedVectors #For loading a saved model
925 from scipy import stats # for zscore
926 import numpy as np # for filtering zscore
927 from collections import Counter #co-occurrence counts
928 import ast #literal evals for list types
929
930 #####
931 app_launch_start=datetime.now() #Set start time for program start
932 print("Loading the dataframe...")
933
934 #Supporting files to load
935 tsne_300d_w5_df = pd.read_pickle('tsne_300d_w5_df.pkl')
936
937 try: #try and except for quick reloading of the W2v object.
938     w2v
939 except:
940     w2v = KeyedVectors.load_word2vec_format('organic_glove_300d_w5.txt') #initialise only once per
941
942 #Dashtables
943 df_updated = pd.read_pickle("df_8hrs.pkl") #filtered 8hr dataset
944 disp1=['ix','caption_processed_4','hashtags','cap_mentions','web_links'] #select columns to display
945 df_disp_1 = df_updated[disp1]
946 df_disp_1.rename(columns={"caption_processed_4": "description"},inplace=True)
947
948 clf_table=pd.read_pickle("df_8hrs.pkl") #z-score table
949 clf_table['Inf_Non']='' #new entry columns
950 clf_table['Domain']='' #new entry columns
951
952 clf_view_6=['ix','time_delta_hours','likeCount','commentCount','likes_pr_min','caption_processed_4']
953 clf_dashtable=clf_table[clf_view_6]
954 clf_dashtable.rename(columns={"caption_processed_4": "description"},inplace=True)
955
956 #-----get vocab & initialise co-occurrence dataframe
957 co_occ_table = pd.read_pickle("App_dataframe_5.pkl") #list col rm_sw_lemt
958 flat_list_co_occ = [item for row_list in co_occ_table['rm_sw_lemt'] for item in row_list] #list co
959 vocab = sorted(set(flat_list_co_occ))
960
961 co_occ_arr = load('ad5_co_occ_arr_w5_b1.npy') #shape 41526
962 co_occ_df = pd.DataFrame(data = co_occ_arr,
963                           index = vocab,
964                           columns = vocab)
965
966 #-----
967
968 print("Dataframes loaded")
969 t=datetime.now() - app_launch_start
970 s=str(t)
971 print("Dataframes loaded in: ", s[:-5], "\n\n")
972
973 #####
974
975 app = dash.Dash(__name__) #loads any CSS directly added in the "assets" folder
976
977 al=str(datetime.now() - app_launch_start) #start time logged at start of program code
978 print("Total time to load files and launch app", al[:-5])
979
980 #-----
981 #Initialising variables outside the callbacks
982
983 #for dropdowns
984 col_sels = ['description','hashtags','cap_mentions','web_links'] #values for dropdown_1
985 col_sels_2 = ['description','hashtags','web_links','cap_mentions','username','fulUserName'] #values
986

```

```

987
988     input_boxes = ['text', 'text']
989     input_boxes1 = ('text', 'text', 'text')
990
991     #Default vocab to plot
992     certs = ['gots', 'oeko', 'cosmos', 'ecocert', 'fsc'] #certifications
993     fruits = ['apple', 'orange', 'vanilla', 'cacao']
994     concs = ['sustainably', 'ethically']
995     rel_concs = ['fruits', 'certification', 'textile', 'cosmetics']
996     vocab_plot_list = certs + fruits + rel_concs
997
998     #other arbitrary initialisations
999     z_init = 5
1000    z_init2 = 5
1001    wc_init = 0
1002
1003    #
1004
1005    app.layout = dbc.Container([
1006        ##### ROW1-Instagram Data, Sel-Desel button_1,Wordcloud #####
1007        dbc.Row([
1008            dbc.Col(html.H2("Instagram DataTable"), width={'size':3}),
1009            dbc.Col(
1010                dbc.Button(id='sel-button', n_clicks=0, children="Sel_all", className="mt-5 mr-2"),
1011                width={'size': 0.5}, style={'textAlign': "left"}), #another way to align
1012            dbc.Col(
1013                dbc.Button(id='desel-button', n_clicks=0, children="Des_all", className="mt-5"),
1014                width={'size': 0.5}, style={'textAlign': "left"}),
1015            dbc.Col(html.H2("Wordcloud - Word frequency"), width={'size':5, 'offset':2}, style={'textA'},
1016                no_gutters=False),
1017        ],
1018        ##### ROW2-Dropdown #####
1019        dbc.Row([
1020            dbc.Col([
1021                dcc.Dropdown(id='my-dropdown', multi=True,
1022                    options=[{'label': x, 'value': x} for x in col_sels],
1023                    value=["description"], #initial values to pass
1024                    style={'width': '690px'}
1025                )])
1026        ]),
1027
1028        ##### ROW3-Dashtable & Wordcloud #####
1029        dbc.Row([
1030            dbc.Col(
1031                dash_table.DataTable(
1032                    id='datatable_id',
1033                    data=df_disp_1.to_dict('records'),
1034                    columns=[
1035                        {"name": i, "id": i, "deletable": False, "selectable": True} for i in df_disp_
1036                    ],
1037                    editable=False,
1038                    filter_action="native",
1039                    sort_action="native",
1040                    sort_mode="multi",
1041                    row_selectable="multi",
1042                    row_deletable=False,
1043                    selected_rows=[], #this parameter gets updated by interactive selection
1044                    page_action="native",
1045                    page_current= 0,
1046                    page_size= 10,
1047                    fixed_rows={ 'headers': True, 'data': 0 }, #if the header style is not defined and
1048                    style_cell=[ #creates a wrapping of the data to constrain column widths. applies to
1049                        'whiteSpace': 'normal',
1050                        'height': 'auto',
1051                        'minWidth': '50px', 'width': '50px', 'maxwidth': '50px', # minwidth of col.
1052                    ],
1053                    style_table= { #For parameters of the table container
1054                        'height': '300px',
1055                        'width': '690px',
1056                        'overflowY': 'auto'
1057                    },
1058                    style_cell_conditional=[
1059                        {'if': {'column_id': 'ix'},
1060                            'width': '50px', 'textAlign': 'left'}, #Setting the px values 1400px seems to
1061                        {'if': {'column_id': 'description'},
1062                            'width': '200px', 'textAlign': 'left'},
1063                        {'if': {'column_id': 'hashtags'},
1064                            'width': '200px', 'textAlign': 'left'},
1065                        {'if': {'column_id': 'cap_mentions'},
1066                            'width': '125px', 'textAlign': 'left'},
1067                        {'if': {'column_id': 'web_links'},
1068                            'width': '125px', 'textAlign': 'left'},
1069                    ],
1070                    width={'size': 6}),
1071            dbc.Col(
1072

```

```

1073     dcc.Graph(id='wordcloud', figure={}, config={'displayModeBar': True},
1074                 style={'width': '600px', 'height': '350px'}),
1075             )
1076         ], no_gutters=False),
1077
1078 ##### ROW4-Headers 2 #####
1079 dbc.Row([
1080     dbc.Col(html.H2("Vector exploration"), width={'size':2}),
1081     dbc.Col(
1082         dbc.Button(id='add1_button', n_clicks=0, children="Add", className="mt-5 mr-2",
1083                     width={'size': 0.5}, style={'textAlign': "left"}),
1084     dbc.Col(
1085         dbc.Button(id='rem1_button', n_clicks=0, children="Rem", className="mt-5 mr-2",
1086                     width={'size': 0.5}, style={'textAlign': "left"}),
1087     dbc.Col(
1088         dbc.Button(id='clr1_button', n_clicks=0, children="Clr", className="mt-5 mr-2",
1089                     width={'size': 0.5}, style={'textAlign': "left"}),
1090     dbc.Col(
1091         dbc.Button(id='plt1_button', n_clicks=0, children="Plot", className="mt-5 mr-2",
1092                     width={'size': 0.5}, style={'textAlign': "left"}),
1093     dbc.Col(
1094         dbc.Button(id='words_tog', n_clicks=0, children="words", className="mt-5 mr-2",
1095                     width={'size': 1, 'offset':1}, style={'textAlign': "left"}), #size 0.5
1096         dbc.Col(html.H2("Vector math"), width={'size':5}, style={'textAlign': "center"}), #If
1097     ], no_gutters=False),
1098
1099 ##### ROW5-Input_boxes #####
1100 dbc.Row([
1101     dbc.Col(
1102         dcc.Input(
1103             id='input_1',
1104             type='text',
1105             placeholder="", #pass list vocab to display here
1106             debounce=False, # changes to input are sent to Dash server only on enter or
1107             pattern=r"^[A-Za-z].*", # Regex: string must start with letters only
1108             spellCheck=True,
1109             inputMode='latin', # provides a hint to browser on type of data that might be en
1110             name='text', # the name of the control, which is submitted with the form d
1111             list='browser', # identifies a list of pre-defined options to suggest to the
1112             n_submit=0, # number of times the Enter key was pressed while the input h
1113             n_submit_timestamp=-1, # last time that Enter was pressed
1114             autoFocus=False, # the element should be automatically focused after the page
1115             n.blur=0, # number of times the input lost focus
1116             n.blur_timestamp=-1, # last time the input lost focus.
1117             size="30" # Width of box which can display placeholders
1118             ), width={'size': 4}, style={'textAlign': "left"} #size was 6
1119     ),
1120     #-----Word2vec_2 math buttons-----
1121     dbc.Col(
1122         dcc.Input(
1123             id='input_2',
1124             type='text',
1125             placeholder='', # A hint to the user of what can be entered in the control
1126             debounce=True, # Changes to input are sent to Dash server onl
1127             min=2015, max=2019, step=1, # Ranges of numeric value. Step refers to incr
1128             minLength=0, maxLength=50, # Ranges for character length inside input box
1129             autoComplete='on',
1130             disabled=False, # Disable input box
1131             readOnly=False, # Make input box read only
1132             required=False, # Require user to insert something into input
1133             size="7", # Sets width of box. If exceeds col width then
1134             ), width={'size': 1, 'offset' : 2}, style={'textAlign': "center"} #text size 5 co
1135     ),
1136     dbc.Col(html.H3("-"), width={'size':0.5}, style={'textAlign' : "center"}), #total width to
1137     dbc.Col(
1138         dcc.Input(
1139             id='input_3',
1140             type='text',
1141             placeholder='', # A hint to the user of what can be entered in the control
1142             debounce=True, # Changes to input are sent to Dash server onl
1143             min=2015, max=2019, step=1, # Ranges of numeric value. Step refers to incr
1144             minLength=0, maxLength=50, # Ranges for character length inside input box
1145             autoComplete='on',
1146             disabled=False, # Disable input box
1147             readOnly=False, # Make input box read only
1148             required=False, # Require user to insert something into input
1149             size="7", # Number of characters that will be visible ins
1150             ), width={'size': 1}, style={'textAlign': "center"})
1151     ),
1152     dbc.Col(html.H3("+"),width={'size': 0.5}, style={'textAlign' : "center"}),
1153     dbc.Col(
1154         dcc.Input(
1155             id='input_4',
1156             type='text',
1157             placeholder='', # A hint to the user of what can be entered in the control
1158             debounce=True, # Changes to input are sent to Dash server onl
1159             ), width={'size': 1}, style={'textAlign': "center"})
1160

```

```

1159         min=2015, max=2019, step=1,          # Ranges of numeric value. Step refers to incr
1160         minLength=0, maxLength=50,          # Ranges for character length inside input box
1161         autoComplete='on',
1162         disabled=False,
1163         readOnly=False,
1164         required=False,
1165         size="7",
1166         ), width={'size': 1}, style={'textAlign': "center"}),
1167     ),
1168     dbc.Col(html.H3(":"), width={'size':0.5}, style={'textAlign' : "center"}),
1169     dbc.Col(
1170       dcc.Input(
1171         id='output_1',
1172         type='text',
1173         placeholder='', # A hint to the user of what can be entered in the control
1174         debounce=True, # Changes to input are sent to Dash server only when debounce is True
1175         min=2015, max=2019, step=1,          # Ranges of numeric value. Step refers to incr
1176         minLength=0, maxLength=50,          # Ranges for character length inside input box
1177         autoComplete='on',
1178         disabled=True,
1179         readOnly=False,
1180         required=False,
1181         size="7",
1182         ), width={'size': 1}, style={'textAlign': "center"})
1183     ),
1184   ], no_gutters=False),
1185
1186 ##### ROW6-word2vec_1 #####
1187 dbc.Row([
1188   dbc.Col(
1189     dcc.Graph(id='word2vec_1'),
1190     width={'size': 6}
1191   ),
1192   dbc.Col(
1193     dcc.Graph(id='word2vec_2'),
1194     width={'size': 6}
1195   )
1196 ], no_gutters=False),
1197
1198 ##### ROW7-range slider #####
1199 dbc.Row([
1200   dbc.Col(
1201     dcc.Slider(
1202       id='my_slider',
1203       min=0,
1204       max=20,
1205       step=1,
1206       value=0,
1207       marks={
1208         0: {'label': '0', 'style': {'color': '#77b0b1'}},
1209         5: {'label': '5'},
1210         10: {'label': '10'},
1211         15: {'label': '15', 'style': {'color': '#f50'}},
1212         20: {'label': '20'}
1213       },
1214       width={'size': 6}),
1215     dbc.Col(
1216       dcc.Slider(
1217         id='my_slider_2',
1218         min=0,
1219         max=20,
1220         step=1,
1221         value=0, #default value
1222         marks={
1223           0: {'label': '0', 'style': {'color': '#77b0b1'}},
1224           5: {'label': '5'},
1225           10: {'label': '10'},
1226           15: {'label': '15', 'style': {'color': '#f50'}},
1227           20: {'label': '20'}
1228         },
1229         width={'size': 6})
1230   ], no_gutters=False),
1231
1232 ##### ROW8—"Classification Dashtable", Sel-Desel #####
1233 dbc.Row([
1234   dbc.Col(html.H2("Z-Score adjustable DataTable"), width={'size':2}),
1235   dbc.Col(
1236     dbc.Button(id='sel-button_2', n_clicks=0, children="Sel_all", className="mt-5 mr-2"),
1237     width={'size': 0.5}, style={'textAlign': "left"}),
1238   dbc.Col(
1239     dbc.Button(id='desel-button_2', n_clicks=0, children="Des_all", className="mt-5"),
1240     width={'size': 0.5}, style={'textAlign': "left"}),
1241   dbc.Col(html.H2("Wordcloud - Post frequency"), width={'size':5, 'offset':3}, style={'textAlign': "center"}, no_gutters=False),
1242
1243 ##### ROW9-Dropdown_2 #####
1244

```

```

1245 |     dbc.Row([
1246 |         dbc.Col(
1247 |             dcc.Dropdown(id='my_dropdown_2', multi=True,
1248 |                         options=[{'label': x, 'value': x} for x in col_sels_2],
1249 |                         value=["description"], #initial values to pass
1250 |                         style={'width':'490px'},
1251 |                         ), width={ 'offset':0}),
1252 |         dbc.Col(
1253 |             dcc.Slider(
1254 |                 id='my_slider_zscore',
1255 |                 min=0,
1256 |                 max=4,
1257 |                 step=0.25,
1258 |                 value=0,
1259 |                 marks={
1260 |                     0: {'label': '0', 'style': {'color': '#77b0b1'}},
1261 |                     1: {'label': '1'},
1262 |                     2: {'label': '2'},
1263 |                     3: {'label': '3', 'style': {'color': '#f50'}},
1264 |                     4: {'label': '4'}
1265 |                 },
1266 |                 width={'size': 2}),
1267 |             ], no_gutters = False),
1268 |
1269 ##### ROW10-Clasification dashtable, WordCloud_2 #####
1270 |     dbc.Row([
1271 |         dbc.Col(
1272 |             dash_table.DataTable(
1273 |                 id='datatable_2',
1274 |                 data=clf_dashtable.to_dict('records'),
1275 |                 columns=[
1276 |                     {"name": i, "id": i, "deletable": False, "selectable": True} for i in clf_
1277 |                 ],
1278 |                 editable=True,
1279 |                 filter_action="native",
1280 |                 sort_action="native",
1281 |                 sort_mode="multi",
1282 |                 row_selectable="multi",
1283 |                 row_deletable=False,
1284 |                 selected_rows=[], #this parameter gets updated by interactive selection
1285 |                 page_action="native",
1286 |                 page_current= 0,
1287 |                 page_size= 10,
1288 |                 fixed_rows={ 'headers': True, 'data': 0 }, #if the header style is not defined
1289 |                 style_cell={ #creates a wrapping of the data to constrain column widths. appli
1290 |                     'whiteSpace': 'normal',
1291 |                     'height': 'auto',
1292 |                     'minWidth': '50px', 'width': '50px', 'maxWidth': '50px', # minwidth of co
1293 |                 },
1294 |                 style_table={ #For parameters of the table container
1295 |                     'height': '600px', #was 300px
1296 |                     'width': '690px',
1297 |                     'overflowY': 'auto'
1298 |                 },
1299 |                 style_cell_conditional=[
1300 |                     {'if': {'column_id': 'description'},
1301 |                      'width': '200px', 'textAlign': 'left'},
1302 |                 ],
1303 |                 export_format='xlsx',
1304 |             ),
1305 |                 width={'size': 6}),
1306 |             dbc.Col(
1307 |                 dcc.Graph(id='wordcloud_2', figure={}, config={'displayModeBar': True},
1308 |                           style={'width': '600px', 'height': '600px'}), #ht was 350px
1309 |             )
1310 |         ],no_gutters=False),
1311 |
1312 ##### ROW11-Header "Posts count of unique words" #####
1313 |     dbc.Row([
1314 |         dbc.Col(html.H2("Posts count of distinct words"), width={'size':2}),
1315 |         dbc.Col(html.H2("Co-occurrence counts"), width={'size':2, 'offset':4}),
1316 |     ],no_gutters=False),
1317 |
1318 ##### ROW12-Slider for most common #####
1319 |     dbc.Row([
1320 |         dbc.Col(
1321 |             dcc.Slider(
1322 |                 id='my_slider_mc',
1323 |                 min=10,
1324 |                 max=30,
1325 |                 step=1,
1326 |                 value=10,
1327 |                 marks={
1328 |                     10: {'label': '10', 'style': {'color': '#77b0b1'}},
1329 |                     20: {'label': '20'},
1330 |                     30: {'label': '30'},
1331 |                 }
1332 |             )
1333 |         )
1334 |     ]
1335 | )

```

```

1331                               40: {'label': '40', 'style': {'color': '#f50'}},
1332                               }),
1333                               width={'size': 2}),
1334                         dbc.Col(
1335                           dbc.Input(
1336                             id='co_occ_ip',
1337                             type='text',
1338                             placeholder='', # A hint to the user of what can be entered in the control
1339                             debounce=True, # Changes to input are sent to Dash server only
1340                             minLength=0, maxLength=50, # Ranges for character length inside input box
1341                             autoComplete='on',
1342                             disabled=False, # Disable input box
1343                             readOnly=False, # Make input box read only
1344                             required=False, # Require user to insert something into input
1345                             size="7",
1346                             ), width={'size': 2, 'offset' : 4}, style={'textAlign': "center"}), #text size 5
1347                         dbc.Col(
1348                           dcc.Slider(
1349                             id='slider_co_occ',
1350                             min=10,
1351                             max=40,
1352                             step=1,
1353                             value=10,
1354                             marks={
1355                               10: {'label': '10', 'style': {'color': '#77b0b1'}},
1356                               20: {'label': '20'},
1357                               30: {'label': '30'},
1358                               40: {'label': '40', 'style': {'color': '#f50'}}},
1359                             ),
1360                             width={'size': 2})
1361                           ], no_gutters=False),
1362
1363 ##### ROW13-Bar_chart for distinct word count #####
1364   dbc.Row([
1365     dbc.Col(
1366       dcc.Graph(id='bar_chart_1'),
1367       width={'size': 6}
1368     ),
1369     dbc.Col(
1370       dcc.Graph(id='bar_chart_2'),
1371       width={'size': 6}
1372     ),
1373   ], no_gutters=False),
1374 ], fluid=True) #Closes initial dbc container
1375
1376 ##### App Callbacks #####
1377 ##### CALLBACKS #####
1378 ##### App Callbacks #####
1379
1380 ##### Barcharts co-occ counts#####
1381 @app.callback(
1382   Output(component_id='bar_chart_2', component_property='figure'),
1383   [Input(component_id='co_occ_ip', component_property='value'),
1384   Input(component_id='slider_co_occ', component_property='value')],
1385   prevent_initial_call=True,
1386 )
1387
1388 def barchart_2(co_occ_ip,slider):
1389   word = co_occ_ip
1390   head_values = slider
1391   try:
1392     co_occ_plot = co_occ_df[word].sort_values(ascending=False).head(head_values)
1393     co_occ_plot = co_occ_plot.reset_index()
1394   except KeyError as e:
1395     print("\n word not in vocab: ", e)
1396   try:
1397     mapping = {co_occ_plot.columns[0]: 'word', co_occ_plot.columns[1]: 'counts'}
1398     co_occ_plot = co_occ_plot.rename(columns=mapping)
1399     fig_br2 = px.bar(co_occ_plot, x="counts", y="word", title="Count of co-occurrences", orientation='v')
1400     fig_br2.update_layout(yaxis=dict(autorange="reversed"))
1401     return fig_br2
1402   except UnboundLocalError as e:
1403     print(e)
1404     raise dash.exceptions.PreventUpdate
1405
1406 ##### Barcharts distinct words per post #####
1407 @app.callback(
1408   Output(component_id='bar_chart_1', component_property='figure'),
1409   [Input(component_id='datatable_2', component_property='selected_rows'),
1410   Input(component_id='my_dropdown_2', component_property='value'),
1411   Input(component_id='my_slider_zscore', component_property='value'),
1412   Input(component_id='my_slider_mc', component_property='value')],
1413   prevent_initial_call=False,
1414 )
1415
1416 def barchart_1(chosen_rows, chosen_cols, zscore, mostcomm):

```

```

1417     global z_init2, flat_list, fig_br1
1418     ctx = dash.callback_context
1419     trigger = (ctx.triggered[0]['prop_id'].split('.')[0])
1420
1421     z = np.abs(stats.zscore(clf_dashtable[['time_delta_hours', 'likes_pr_min']])) #assign z-score to
1422     view_rm_lm= ['rm_sw_lemt']
1423     view_rm_lm = clf_view_6 + view_rm_lm
1424     clf_dashtable_z=clf_table[view_rm_lm].iloc[np.unique(np.where(z > zscore)[0])] #table should have
1425     #clean up index for "select_all" functionality
1426     clf_dashtable_z.drop(['ix'], axis=1, inplace = True) #drop any older 'ix' columns
1427     clf_dashtable_z.reset_index(inplace= True,drop=True) #reset and drop old index
1428     clf_dashtable_z.reset_index(inplace= True) #create new ix column
1429     clf_dashtable_z.rename(columns={"index": "ix"},inplace=True)
1430
1431     if len(chosen_cols) > 0: #consider rm_sw_lemt instead
1432         if 'description' in chosen_cols:
1433             chosen_cols.remove('description')
1434             chosen_cols.append('rm_sw_lemt')
1435             print("Chosen cols now contain",chosen_cols)#atleast 1 col to be selected
1436     if len(chosen_rows)==0: #if no rows selected consider all rows of z score filtered table-d
1437         #----fast render optimization
1438         if ((z_init2 == zscore) and (trigger not in ['my_slider_mc','my-dropdown_2'])): #if no
1439             print("\nz_init2 is equal to zscore, z_init2 value is: ", z_init2)
1440             return (fig_br1)
1441         if ((z_init2 == zscore) and (trigger in ['my_slider_mc'])): #uses flat_list global for
1442             counts = dict(Counter(flat_list).most_common(mostcomm))
1443             labels, values = zip(*counts.items())
1444             indSort = np.argsort(values)[::-1] # sort your values in descending order
1445             labels = np.array(labels)[indSort]
1446             values = np.array(values)[indSort]
1447             indexes = np.arange(len(labels))
1448             word_count_df = pd.DataFrame({'word': labels, 'posts_count': values}, columns=['wo
1449             fig_br_12 = px.bar(word_count_df, x="posts_count", y="word", title="Lemmatised wor
1450             fig_br_12.update_layout(yaxis=dict(autorange="reversed"))
1451             return (fig_br_12)
1452         else: #initialise object and set z_init value
1453             print("\nElse block reached, z_init2 value is; ", z_init2)
1454             z_init2 = zscore
1455             print("z_init2 value updated to: ", z_init2)
1456
1457         df_filtered = clf_dashtable_z[chosen_cols]
1458         df_filtered['comb_cols'] = df_filtered[df_filtered.columns[0:]].apply(
1459             lambda x: ' '.join(x.dropna().astype(str)),
1460             axis=1)#Combine multiple columns into a single series
1461         df_filtered['comb_cols'] = df_filtered['comb_cols'].apply(
1462             lambda x: ' '.join(set(x.split()))))#Use only unique words per post in wordcloud
1463         #Create word counts dataframe
1464         flat_list = [item for row_list in df_filtered['comb_cols'] for item in row_list.sp
1465         counts = dict(Counter(flat_list).most_common(mostcomm))
1466         labels, values = zip(*counts.items())
1467         indSort = np.argsort(values)[::-1] # sort your values in descending order
1468         labels = np.array(labels)[indSort]
1469         values = np.array(values)[indSort]
1470         indexes = np.arange(len(labels))
1471         word_count_df = pd.DataFrame({'word': labels, 'posts_count': values}, columns=['wo
1472         fig_br1 = px.bar(word_count_df, x="posts_count", y="word", title="Lemmatised words
1473         fig_br1.update_layout(yaxis=dict(autorange="reversed"))
1474         return (fig_br1)
1475     elif len(chosen_rows) > 0 : #filter to selected rows
1476         df_filtered = clf_dashtable_z[chosen_cols]
1477         df_filtered = df_filtered[df_filtered.index.isin(chosen_rows)]
1478     elif len(chosen_cols) == 0:
1479         raise dash.exceptions.PreventUpdate #no update if no cols selected
1480
1481     ##figure logic repeated for row selections
1482     df_filtered['comb_cols'] = df_filtered[df_filtered.columns[0:]].apply(
1483         lambda x: ' '.join(x.dropna().astype(str)),
1484         axis=1)
1485     df_filtered['comb_cols'] = df_filtered['comb_cols'].apply(
1486         lambda x: ' '.join(set(x.split())))
1487     flat_list = [item for row_list in df_filtered['comb_cols'] for item in row_list.split()]
1488     counts = dict(Counter(flat_list).most_common(mostcomm))
1489     labels, values = zip(*counts.items())
1490     indSort = np.argsort(values)[::-1] # sort your values in descending order
1491     labels = np.array(labels)[indSort]
1492     values = np.array(values)[indSort]
1493     indexes = np.arange(len(labels))
1494     word_count_df = pd.DataFrame({'word': labels, 'posts_count': values}, columns=['word', 'posts_
1495     fig_br_1 = px.bar(word_count_df, x="posts_count", y="word", title="Lemmatised words by post co
1496     fig_br_1.update_layout(yaxis=dict(autorange="reversed"))
1497     return (fig_br_1) #note change in return object name
1498
1499 ##### Column highlighting Dashtable2 #####
1500 @app.callback(
1501     Output('datatable_2', 'style_data_conditional'),
1502     Input(component_id='my-dropdown_2', component_property='value')

```

```

1503 )
1504
1505 def update_styles(selected_columns):
1506     return [
1507         'if': { 'column_id': i },
1508         'background_color': '#D2F3FF'
1509     } for i in selected_columns]
1510
1511 ##### Dashtable_2 Zscore virtual table#####
1512 @app.callback(
1513     Output('datatable_2', 'data'), #virtual_data is displayed table. Even after filtering. #referenc
1514     Input('my_slider_zscore','value'),
1515     prevent_initial_call=False
1516 )
1517
1518 def zscore_outliers(zscore):
1519     z = np.abs(stats.zscore(clf_dashtable[['time_delta_hours', 'likes_pr_min']])) #assign z-score to
1520     clf_dashtable_z=clf_dashtable.iloc[np.unique(np.where(z > zscore)[0])] #create outlier table
1521     #clean up index for "select_all" functionality
1522     clf_dashtable_z.drop(['ix'], axis=1, inplace = True) #drop any older 'ix' columns
1523     clf_dashtable_z.reset_index(inplace= True,drop=True) #reset and drop old index
1524     clf_dashtable_z.reset_index(inplace= True) #create new ix column
1525     clf_dashtable_z.rename(columns={"index": "ix"},inplace=True)
1526     print("\nZscore slider value is: ", zscore)
1527     print("No of rows are: ", len(clf_dashtable_z))
1528     return clf_dashtable_z.to_dict('records')
1529
1530 ##### Dashtable_2 Sel_Desel button_2 #####
1531 @app.callback(
1532     [Output('datatable_2', 'selected_rows')], #references ordered 0 to n index of larger table irr
1533     [Input('sel-button_2', 'n_clicks'),
1534     Input('desel-button_2', 'n_clicks'),
1535     Input('my_slider_zscore','value')],
1536     [State('datatable_2', 'derived_virtual_selected_rows'), #virtual selected row is what is selec
1537     State('datatable_2', 'derived_virtual_data')], #virtual_data is displayed table. Even after f
1538     prevent_initial_call=True
1539 )
1540
1541 def select_deselect(selbtn, deselbtn, z_score_trigger,selected_rows,filtered_table):
1542     ctx = dash.callback_context
1543     if ctx.triggered:
1544         print(ctx.triggered)
1545         trigger = (ctx.triggered[0]['prop_id'].split('.')[0])
1546         if trigger == 'sel-button_2':
1547             print("\n\nSelect button clicked")
1548             print("Length of Selected_rows is: ", len(selected_rows))
1549             wc_list= []
1550             wc_list=[[row['ix'] for row in filtered_table]]
1551             wc_list = [[int(i) for i in wc_list[0]]] #convert to int for index reference
1552             print("Displayed table has:", len(filtered_table), 'rows')
1553             print("Wordcloud list contains:", len(wc_list[0]), "elements")
1554             return wc_list
1555         else:
1556             print("\n\nDeselect button clicked")
1557             print("Length of Selected_rows is: ", len(selected_rows))
1558             return []
1559
1560 ##### Wordcloud_2 #####
1561 @app.callback(
1562     Output('wordcloud_2', 'figure'),
1563     [Input(component_id='datatable_2',component_property='selected_rows'),
1564     Input(component_id='my-dropdown_2', component_property='value'),
1565     Input('my_slider_zscore','value')],
1566     prevent_initial_call=False
1567 )
1568
1569 def ren_wordcloud(chosen_rows, chosen_cols, zscore):
1570     global z_init, fig_wordcloud2 #for try and except. otherwise destroyed at end of callback.
1571     ctx = dash.callback_context
1572     trigger = (ctx.triggered[0]['prop_id'].split('.')[0])
1573
1574     z = np.abs(stats.zscore(clf_dashtable[['time_delta_hours', 'likes_pr_min']])) #assign z-score to
1575     clf_dashtable_z=clf_dashtable.iloc[np.unique(np.where(z > zscore)[0])] #create outlier table
1576     #clean up index for "select_all" functionality
1577     clf_dashtable_z.drop(['ix'], axis=1, inplace = True) #drop any older 'ix' columns
1578     clf_dashtable_z.reset_index(inplace= True,drop=True) #reset and drop old index
1579     clf_dashtable_z.reset_index(inplace= True) #create new ix column
1580     clf_dashtable_z.rename(columns={"index": "ix"},inplace=True)
1581
1582     if len(chosen_cols) > 0: #atleast 1 col to be selected
1583         if len(chosen_rows)==0:
1584             #----fast render optimization
1585             if ((z_init == zscore) and (trigger != 'my-dropdown_2')): #if no change in the score r
1586                 print("\nz_init is equal to zscore, z_init value is: ", z_init)
1587                 return fig_wordcloud2
1588             else: #initialise object and set z_init value

```

```

1589         print("\nElse block reached, z_init value is: ", z_init)
1590         z_init=zscore
1591         print("z_init value updated to: ", z_init)
1592     #-----
1593     #if no rows selected consider all rows of z score filtered table-defined again
1594     df_filtered = clf_dashtable_z[chosen_cols]
1595     print("Initialising full word cloud..")
1596     df_filtered['comb_cols'] = df_filtered[df_filtered.columns[0:]].apply( #Combine mu
1597         lambda x: ' '.join(x.dropna().astype(str)),
1598         axis=1)
1599     df_filtered['comb_cols'] = df_filtered['comb_cols'].apply(
1600         lambda x: ' '.join(set(x.split())))#Use only unique words per post in wordcloud
1601
1602     en_stopwords = stopwords.words('english')
1603     fr_stopwords = stopwords.words('french')
1604     web_links_sw = ['www', 'http', 'https', 'com']
1605     combined_stopwords = en_stopwords + fr_stopwords + web_links_sw
1606     print("\nNo of rows in WC are: ",len(df_filtered))
1607     wordcloud = WordCloud(max_words=100,
1608                           # stopwords= combined_stopwords,
1609                           colormap='tab20b',
1610                           background_color='white',
1611                           width=1200, #1200,1700
1612                           height=1000, #1000
1613                           random_state=1).generate_from_text(' '.join(df_filtered['comb_cols']))
1614
1615     fig_wordcloud2 = px.imshow(wordcloud, template='ggplot2')#title="Post frequency"
1616     fig_wordcloud2.update_layout(margin=dict(l=0, r=0,b=0,t=0))
1617     fig_wordcloud2.update_xaxes(visible=False)
1618     fig_wordcloud2.update_yaxes(visible=False)
1619     return fig_wordcloud2
1620
1621     elif len(chosen_rows) > 0 :
1622         df_filtered = clf_dashtable_z[chosen_cols]
1623         df_filtered = df_filtered[df_filtered.index.isin(chosen_rows)]
1624     elif len(chosen_cols) == 0:
1625         raise dash.exceptions.PreventUpdate
1626 ## load wordcloud only once
1627     df_filtered['comb_cols'] = df_filtered[df_filtered.columns[0:]].apply( #Combine multiple columns
1628         lambda x: ' '.join(x.dropna().astype(str)),
1629         axis=1)
1630     df_filtered['comb_cols'] = df_filtered['comb_cols'].apply(
1631         lambda x: ' '.join(set(x.split())) ) #because no order, bigrams become irrelevant
1632
1633     en_stopwords = stopwords.words('english')
1634     fr_stopwords = stopwords.words('french')
1635     web_links_sw = ['www', 'http', 'https', 'com']
1636     combined_stopwords = en_stopwords + fr_stopwords + web_links_sw
1637     print("\nNo of rows in WC are: ",len(df_filtered))
1638     wordcloud = WordCloud(max_words=100,
1639                           # stopwords= combined_stopwords,
1640                           colormap='tab20b',
1641                           background_color='white',
1642                           width=1200, #1200,1700
1643                           height=1000, #1000
1644                           random_state=1).generate_from_text(' '.join(df_filtered['comb_cols'])) #
1645
1646     fig_wordcloud_2 = px.imshow(wordcloud, template='ggplot2')#title="Post frequency"
1647     fig_wordcloud_2.update_layout(margin=dict(l=0, r=0,b=0,t=0))
1648     fig_wordcloud_2.update_xaxes(visible=False)
1649     fig_wordcloud_2.update_yaxes(visible=False)
1650     return fig_wordcloud_2
1651 ###### word2vec_2 math graph #####
1652 @app.callback(
1653     [Output(component_id='word2vec_2', component_property='figure'),
1654      Output(component_id='output_1', component_property='placeholder')],
1655      [Input('input_2','value'),
1656       Input('input_3','value'),
1657       Input('input_4','value'),
1658       Input('my_slider_2','value'),
1659       Input('words_tog','n_clicks')], prevent_initial_call=False
1660 )
1661
1662
1663 def word2vec_math(ip1,ip2,ip3,slider2_val,word_tog):
1664     tsne_df_full = tsne_300d_w5_df
1665     vocab_plot_list_2 = [ip1,ip2,ip3]
1666     print("\nvocab_plot_list is: ", vocab_plot_list_2)
1667     vocab_to_plot = list(filter(None, vocab_plot_list_2))
1668     print("Filtered vocab_to_plot is: ", vocab_to_plot) #Only non-empty inputs
1669     print("Slider2 val: ",slider2_val)
1670     nearest_size = slider2_val #nearest size defined
1671     tsne_df_full['type'] = 'Sel'
1672     ip_vocab_index=list(tsne_df_full[tsne_df_full['word'].isin(vocab_to_plot)].index.values)
1673     clos_ten_out = [] #index of closest 10 words
1674     try:

```

```

1675     for i,word in enumerate(vocab_to_plot):
1676         ix_word = tsne_df_full[tsne_df_full['word'].isin([word])].index.values
1677         tsne_df_full.loc[ix_word, 'type'] = "ip_{}".format(i)
1678         if nearest_size != 0:
1679             a=w2v.most_similar(word, topn = nearest_size)
1680             clos_ten_in = [i[0] for i in a] #grab vectors and not similarity probability
1681             clos_ten_out.append(clos_ten_in) #gather words in a single list
1682         else:
1683             clos_ten_out = []
1684     except ValueError as e:#Empty word None input
1685         print(e)
1686         pass
1687     except KeyError as e: #Word not in vocab
1688         print(e)
1689         pass
1690
1691     ip_nearest=[]
1692     clos_ten_dict={}
1693     try:
1694         for i,sub_list in enumerate(clos_ten_out):
1695             clos_ten_dict[i] = list(tsne_df_full[tsne_df_full['word'].isin(sub_list)].index.values)
1696             tsne_df_full.loc[clos_ten_dict[i], 'type'] = "near_ip_{}".format(i)
1697             ip_nearest.append(clos_ten_dict[i]) #index values
1698     except Exception as e:
1699         print("clos_ten_out is empty: ", e)
1700         pass
1701     ip_nearest_ix = [item for sublist in ip_nearest for item in sublist]
1702     try: #Try catches no input case
1703         r=w2v.most_similar(positive=[ip1,ip3], negative= [ip2], topn = nearest_size+1)
1704         res_ten_in = [i[0] for i in r]
1705         print("Res words are: ", res_ten_in)
1706         res_ten_index = list(tsne_df_full[tsne_df_full['word'].isin(res_ten_in)].index.values)
1707     except Exception as e: #TypeError
1708         print("Trying addition ")
1709     try:
1710         r=w2v.most_similar(positive=[ip1,ip3], topn = nearest_size+1)
1711         res_ten_in = [i[0] for i in r]
1712         print("Res words of addition: ", res_ten_in)
1713         res_ten_index = list(tsne_df_full[tsne_df_full['word'].isin(res_ten_in)].index.values)
1714     except Exception as e:
1715         print("Trying subtraction: ")
1716     try:
1717         r=w2v.most_similar(positive=[ip1], negative= [ip2], topn = nearest_size+1)
1718         res_ten_in = [i[0] for i in r]
1719         print("Res words of subtraction: ", res_ten_in)
1720         res_ten_index = list(tsne_df_full[tsne_df_full['word'].isin(res_ten_in)].index.values)
1721     except Exception as e:
1722         print("Insufficient inputs",e)
1723         pass
1724     try:
1725         tsne_df_full.loc[res_ten_index[0], 'type'] = 'res'
1726         tsne_df_full.loc[res_ten_index[1:], 'type'] = 'near_res'
1727     except Exception as e:
1728         print("res nearest not defined: ", e)
1729         pass
1730     try:
1731         plot_index = ip_vocab_index + ip_nearest_ix + res_ten_index
1732     except Exception as e:
1733         print("plot index not fully defined: ",e)
1734         pass
1735     try:
1736         print("res_ten_index, word, type is: ", res_ten_index,tsne_df_full.loc[res_ten_index, 'word'])
1737     except UnboundLocalError as e:
1738         print(e)
1739         pass
1740     try:
1741         tsne_plot_df = tsne_df_full.loc[plot_index]
1742     except UnboundLocalError as e:
1743         print(e)
1744         pass
1745     ctx = dash.callback_context
1746     if ctx.triggered:
1747         # print(ctx.triggered)
1748         trigger = (ctx.triggered[0]['prop_id'].split('.')[0])
1749         print("trigger is: ", trigger)
1750         print("word_tog nclick value is: ", word_tog)
1751         if trigger == 'words_tog':
1752             if (word_tog % 2) == 0:
1753                 print("Words display toggled")
1754                 text_disp = 'word'
1755             else:
1756                 text_disp = None
1757         if nearest_size != 0 :
1758             order = ["ip_0", "ip_1", "ip_2", "res","near_ip_0", "near_ip_1", "near_ip_2", "near_res"]
1759         else:
1760             order = ["ip_0", "ip_1", "ip_2", "res"]
```

```

1761     #accommodate edge cases where input ex "near_ip_0" does not exist etc.
1762     pop_list_flag = 1
1763     while pop_list_flag == 1: #To iteratively remove undefined inputs
1764         try:
1765             df = tsne_plot_df.set_index('type') #set as index so that order can be specified
1766             df_ordered = df.T[order].T.reset_index()
1767         except KeyError as e:
1768             print("Exception in df_ordered dataframe: ",e)
1769             print("Entering while loop to remove items from order list")
1770             print("Initial order is: ", order)
1771             s=str(e)
1772             rem_val=s[s.find("[") : s.find("]")+1]
1773             rem_val = ast.literal_eval(rem_val)
1774             order.remove(rem_val[0])
1775             print("order is now: ", order)
1776         except UnboundLocalError as e:
1777             print(e)
1778             pop_list_flag =0 #Else it continues in an endless while loop
1779             break
1780         else: #if try is succesful, end the while loop.s
1781             pop_list_flag =0
1782     try: #if fig_1 is in the main program, "referenced before assignment" error.
1783         fig_2 = px.scatter(df_ordered, x="x", y="y", text= text_disp, color='type',
1784                             color_discrete_map={
1785                                 "ip_0": "#7bc043", #green- most natural word by essence "concept1"
1786                                 "ip_1": "#fdf498", #yellow- noticeable contrast from "concept1";
1787                                 "ip_2": "#f37736", #orange- "concept2b"lighter contrast from "con
1788                                 "res": "#ee4035", #red- attracts first attention; strange that it
1789                                 "near_ip_0": "#339900",
1790                                 "near_ip_1": "#ffcc00",
1791                                 "near_ip_2": "#ff9966",
1792                                 "near_res": "#cc3300"}, log_x=False, hover_name = 'word')
1793     except UnboundLocalError as e:
1794         print (e)
1795         try:
1796             fig_2 = px.scatter(df_ordered, x="x", y="y", text= 'word', color='type',
1797                             color_discrete_map={
1798                                 "ip_0": "#7bc043", #green- most natural word by essence "concept1"
1799                                 "ip_1": "#fdf498", #yellow- noticeable contrast from "concept1";
1800                                 "ip_2": "#f37736", #orange- "concept2b"lighter contrast from "con
1801                                 "res": "#ee4035", #red- attracts first attention; strange that it
1802                                 "near_ip_0": "#339900",
1803                                 "near_ip_1": "#ffcc00",
1804                                 "near_ip_2": "#ff9966",
1805                                 "near_res": "#cc3300"}, hover_name = 'word', log_x=False)
1806     except UnboundLocalError as e:
1807         print(e)
1808         pass
1809     try:
1810         fig_2.update_traces(textposition='top center')
1811         fig_2.update_layout(legend_traceorder="normal")
1812         fig_2.update_layout(margin=dict(l=0, r=0),uirevision = True)
1813         return (fig_2), res_ten_in[0]
1814     except Exception as e:
1815         print("No return values. Error: ",e)
1816         raise dash.exceptions.PreventUpdate
1817         pass
1818 ###### Arbitrary_graphing_inputbox #####
1819 @app.callback(
1820     [Output(component_id='input_1', component_property='placeholder'),
1821      Output(component_id='input_1', component_property='value')],
1822      [Input('add1_button', 'n_clicks'),
1823       Input('rem1_button', 'n_clicks'),
1824       Input('clr1_button', 'n_clicks')],
1825      [State(component_id='input_1', component_property='value')],
1826      prevent_initial_call=False
1827 )
1828 def vocab_list(add,rem,clr,inp_1):
1829     ctx = dash.callback_context
1830     if ctx.triggered:
1831         print(ctx.triggered)
1832         trigger = (ctx.triggered[0]['prop_id'].split('.')[0])
1833         if trigger == 'add1_button':
1834             if ',' in inp_1: #if typing a list directly split into its words
1835                 inp_1=inp_1.replace(" ", "") #remove whitespace
1836                 split_words=inp_1.split(',')
1837                 print("Input words added are: ",split_words)
1838                 vocab_plot_list.extend(split_words)
1839                 value=''
1840                 return vocab_plot_list, value
1841             else:
1842                 print("Input word added: ",inp_1)

```

```

1847             vocab_plot_list.append(inp_1)
1848             print("vocab_list is: ",vocab_plot_list )
1849             value=''
1850             return vocab_plot_list, value
1851         if trigger == 'rem1_button':
1852             print("Last word removed")
1853             vocab_plot_list.pop()
1854             print("vocab_list is now: ",vocab_plot_list)
1855             value=''
1856             return vocab_plot_list, value
1857         if trigger == 'clr1_button':
1858             print("Clr button clicked")
1859             vocab_plot_list.clear()
1860             print("vocab_list is now: ",vocab_plot_list )
1861             value=''
1862             return vocab_plot_list, value
1863
1864     else:
1865         vocab_plot_list.append(inp_1)
1866         value=''
1867         return vocab_plot_list, value
1868
1869 ##### Arb_Word2vec_1 graph #####
1870 @app.callback(
1871     Output(component_id='word2vec_1', component_property='figure'),
1872     [Input('add1_button','n_clicks'),
1873      Input('rem1_button','n_clicks'),
1874      Input('clr1_button','n_clicks'),
1875      Input('plt1_button','n_clicks'),
1876      Input('my_slider','value'),
1877      Input('words_tog','n_clicks'),
1878      Input('input_1','value')], #debounce of input box triggers callback of figure
1879      prevent_initial_call=False
1880 )
1881
1882 def svd_user_inputs(ad,rem,clr,plot_butts,slider_val,word_tog,ip_tog):
1883
1884     vocab_to_plot = vocab_plot_list #initialised at the start of program
1885     tsne_df_full = tsne_300d_w5_df
1886     print("Slider val: ",slider_val)
1887     nearest_size = slider_val
1888
1889     tsne_df_full['type'] = 'Sel'
1890     ip_vocab_index=list(tsne_df_full[tsne_df_full['word'].isin(vocab_to_plot)].index.values)
1891     print("vocab to plot of slider graph is: ",vocab_to_plot )
1892
1893     clos_ten_out = [] #index of closest 10 words
1894     for word in vocab_to_plot:
1895         try:
1896             a=w2v.most_similar(word, topn = nearest_size)
1897             clos_ten_in = [i[0] for i in a]
1898             clos_ten_out.append(clos_ten_in)
1899         except ValueError as e:#Empty word None input
1900             print(e)
1901             pass
1902         except KeyError as e: #Word not in vocab
1903             print(e)
1904             pass
1905     print("close_ten_words are: ",clos_ten_out)
1906     clos_ten_flat = [item for sublist in clos_ten_out for item in sublist]
1907     clos_ten_index = list(tsne_df_full[tsne_df_full['word'].isin(clos_ten_flat)].index.values)
1908
1909     plot_index = ip_vocab_index + clos_ten_index
1910     tsne_df_full.loc[ip_vocab_index, 'type'] = 'ip' # indices = [0,1,3,6,10,15]
1911     tsne_df_full.loc[clos_ten_index, 'type'] = 'near'
1912     tsne_plot_df = tsne_df_full.loc[plot_index]
1913
1914     ctx = dash.callback_context
1915     if ctx.triggered:
1916         trigger = (ctx.triggered[0]['prop_id'].split('.')[0])
1917         print("trigger is: ", trigger)
1918         print("word_tog nclick value is: ", word_tog)
1919         if trigger == 'words_tog':
1920             if (word_tog % 2) == 0:
1921                 print("Words display toggled")
1922                 text_disp = 'word'
1923             else:
1924                 text_disp = None
1925             try: #if fig_1 is in the main program, "referenced before assignment" error.
1926                 fig_1 = px.scatter(tsne_plot_df, x="x", y="y", text= text_disp, color='type',
1927                                     log_x=False, hover_name = 'word')
1928             except UnboundLocalError as e:
1929                 print (e)
1930                 fig_1 = px.scatter(tsne_plot_df, x="x", y="y", text= 'word', color='type',
1931                                     hover_name = 'word', log_x=False)
1932             pass

```

```

1933
1934     fig_1.update_traces(textposition='top center')
1935     fig_1.update_layout(margin=dict(l=0, r=0),uirevision = True)
1936     return (fig_1)
1937
1938 ##### Select all button #####
1939 @app.callback(
1940     [Output('datatable_id', 'selected_rows')], #references ordered 0 to n index of larger table in
1941     [Input('sel-button', 'n_clicks'),
1942      Input('desel-button', 'n_clicks')],
1943     [State('datatable_id', 'derived_virtual_selected_rows'), #virtual selected row is what is sele
1944      State('datatable_id', 'derived_virtual_data')], #virtual_data is displayed table. Even after
1945      prevent_initial_call=True
1946 )
1947
1948 def select_deselect(selbtn, deselbtn, selected_rows,filtered_table):
1949     ctx = dash.callback_context
1950     if ctx.triggered:
1951         print(ctx.triggered)
1952         trigger = (ctx.triggered[0]['prop_id'].split('.')[0])
1953         if trigger == 'sel-button':
1954             print("\n\nSelect button clicked")
1955             print("Length of Selected_rows is: ", len(selected_rows))
1956             wc_list= []
1957             wc_list=[[row['ix'] for row in filtered_table]]
1958             wc_list = [int(i) for i in wc_list[0]] #convert to int for index reference
1959             print("Displayed table has:", len(filtered_table), 'rows')
1960             print("Wordcloud list contains:", len(wc_list[0]), "elements")
1961             return wc_list
1962         else:
1963             print("\n\nDeselect button clicked")
1964             print("Length of Selected_rows is: ", len(selected_rows))
1965             return []
1966
1967 ##### Wordcloud callback #####
1968 @app.callback(
1969     Output('wordcloud', 'figure'),
1970     [Input(component_id='datatable_id', component_property='selected_rows'),
1971      Input(component_id='my-dropdown', component_property='value')],
1972      prevent_initial_call=False
1973 )
1974
1975 def ren_wordcloud(chosen_rows, chosen_cols):
1976     global wc_init, fig_wordcloud1 #global for pre-app saving of the object
1977     if len(chosen_cols) > 0: #atleast 1 col to be selected
1978         if len(chosen_rows)==0:
1979             df_filtered = df_disp_1[chosen_cols] #if no rows selected consider all rows
1980             if wc_init == 1: #load once.
1981                 return fig_wordcloud1
1982             else:
1983                 wc_init = 1 #object initialised flag set
1984
1985             print("Full Word cloud not initialised")
1986             print("Initialising full word cloud..")
1987             df_filtered['comb_cols'] = df_filtered[df_filtered.columns[0:]].apply( #Combine mu
1988                 lambda x: ' '.join(x.dropna().astype(str)),
1989                 axis=1)
1990
1991             en_stopwords = stopwords.words('english')
1992             fr_stopwords = stopwords.words('french')
1993             web_links_sw = [ 'www', 'http', 'https', 'com' ]
1994             combined_stopwords = en_stopwords + fr_stopwords + web_links_sw
1995             print("\nNo of rows in WC are: ",len(df_filtered))
1996             wordcloud = WordCloud(max_words=100,
1997                 # stopwords= combined_stopwords,
1998                 colormap='tab20c',
1999                 background_color='white',
2000                 width=1700, #1200,1700
2001                 height=1000, #1000
2002                 random_state=1).generate(' '.join(df_filtered['comb_cols']))
2003
2004             fig_wordcloud1 = px.imshow(wordcloud, template='ggplot2')#title="Word frequency"
2005             fig_wordcloud1.update_layout(margin=dict(l=0, r=0,b=0,t=0))
2006             fig_wordcloud1.update_xaxes(visible=False)
2007             fig_wordcloud1.update_yaxes(visible=False)
2008             return fig_wordcloud1
2009         elif len(chosen_rows) > 0 :
2010             df_filtered = df_disp_1[chosen_cols]
2011             df_filtered = df_filtered[df_filtered.index.isin(chosen_rows)]
2012         elif len(chosen_cols) == 0:
2013             raise dash.exceptions.PreventUpdate
2014 ## load wordcloud only once
2015         df_filtered['comb_cols'] = df_filtered[df_filtered.columns[0:]].apply( #Combine multiple column
2016                 lambda x: ' '.join(x.dropna().astype(str)),
2017                 axis=1)
2018

```

```

2019     en_stopwords = stopwords.words('english')
2020     fr_stopwords = stopwords.words('french')
2021     web_links_sw = ['www', 'http', 'https', 'com']
2022     combined_stopwords = en_stopwords + fr_stopwords + web_links_sw
2023     print("\nNo of rows in WC are: ", len(df_filtered))
2024     wordcloud = WordCloud(max_words=100,
2025                           # stopwords= combined_stopwords,
2026                           colormap='tab20c',
2027                           background_color='white',
2028                           width=1700, #1200,1700
2029                           height=1000, #1000
2030                           random_state=1).generate(' '.join(df_filtered['comb_cols'])) #df_filtered
2031     fig_wordcloud = px.imshow(wordcloud, template='ggplot2')#title="Word frequency"
2032     fig_wordcloud.update_layout(margin=dict(l=0, r=0,b=0,t=0))
2033     fig_wordcloud.update_xaxes(visible=False)
2034     fig_wordcloud.update_yaxes(visible=False)
2035     return fig_wordcloud
2036
2037 ##### Column highlighting #####
2038 @app.callback(
2039     Output('datatable_id', 'style_data_conditional'),
2040     Input(component_id='my-dropdown', component_property='value')
2041 )
2042
2043 def update_styles(selected_columns):
2044     return [
2045         {
2046             'if': { 'column_id': i },
2047             'background_color': '#D2F3FF'
2048         } for i in selected_columns]
2049
2050 if __name__ == '__main__':
2051     app.run_server(debug=False)

```

## Summary – English and French

Keywords: NLP, dashboard, wordembeddings, socialmedia, instagram, word2vec

Paying attention to language use is very useful in understanding what people think about a certain subject. The way a Particular word is used in different contexts, defines the meaning of the word. As such, the meaning of a word evolves over time and new words may often come into existence. In this project we build an interactive tool which uses natural language processing (NLP) techniques, to analyze text and infer the current meaning of the word ‘organic’ as applied to different vocations, from Instagram post descriptions. When a person or company uses a particular hashtag for a post in social media (ex. #organic), it is implied that the text descriptions they write for the posts, must be related to this hashtag. We use visualizations of word frequency counts, co-occurrences, and word embeddings from the word2vec skip gram method, to gather insights about the subject.

As for the insights, the premise is that the word ‘organic’ assumes different meanings in different contexts. For example, the meaning of ‘organic’ in the context of textiles, food, and cosmetics, at the highest level, carries connotations of ‘health’, ‘certifications’, and ‘unadulterated’ respectively. This interactive dashboard tool can be useful for companies and individuals trying to understand how the meaning of a particular word/concept is being perceived by their target audience and inform their marketing campaigns. An approach for using the tool is described in detail, using “organic food” as an example context of study.

Résumé: NLP, tableau de bord, wordembeddings, socialmedia, Instagram, word2vec

Prêter attention à l'utilisation de la langue est très utile pour comprendre ce que les gens pensent d'un certain sujet. La façon dont un mot particulier est utilisé dans différents contextes définit le sens du mot. Ainsi, le sens d'un mot évolue avec le temps et de nouveaux mots peuvent souvent voir le jour. Dans ce projet, nous construisons un outil interactif qui utilise des techniques de traitement du langage naturel (PNL) pour analyser le texte et déduire la signification actuelle du mot « biologique » tel qu'il est appliqué à différentes vocations, à partir des descriptions de publications Instagram. Lorsqu'une personne ou une entreprise utilise un hashtag particulier pour une publication sur les réseaux sociaux (par exemple, #organique), il est sous-entendu que les descriptions textuelles qu'elle rédige pour les publications doivent être liées à ce hashtag. Nous utilisons des visualisations du nombre de fréquences de mots, des cooccurrences et des incorporations de mots de la méthode word2vec skip gram, pour recueillir des informations sur le sujet. En ce qui concerne les idées, la prémissse est que le mot « biologique » prend différentes significations dans différents contextes. Par exemple, la signification de « biologique » dans le contexte des textiles, des aliments et des cosmétiques, au plus haut niveau, a des connotations de « santé », de « certifications » et de « pur » respectivement. Une approche d'utilisation de l'outil est décrite en détail, en utilisant « l'alimentation biologique » comme exemple de contexte d'étude.