

## **Advanced Data Mining Analytics Multilevel Marketing on Twitter**

### **Business Understanding and Background**

The principal user activities on social media platforms like Twitter, Facebook, Instagram, and LinkedIn have expanded beyond connecting socially to marketing and sales for small businesses. Some of these small money making ventures include legitimate small businesses like the sale of handmade items on platforms like Etsy and Minted, or secondhand sales of clothing and jewelry on platforms like Poshmark and RetailMeNot. A large number of these small businesses are actually rungs within the pyramid of a multilevel marketing company.

Multilevel marketing is problematic in that the companies are considered legitimate businesses. These companies sell products, like Lularoe sells patterned activewear for women, or services, like Primerica sells insurance policies; however, the majority of the individual seller's income is generated by recruiting new salespeople and receiving a cut of their recruits' sales. These businesses easily oversaturate the market and only early participants will receive significant profit, while people who join later rarely recuperate their sign-on fees.

One of the top strategies of multilevel marketing companies is participants posting overly enthusiastic commentary on social media platforms to hook new recruits. With all the noise within social media platforms, it can be difficult to sniff out a multilevel marketing company from a legitimate small business or side hustle. Since these companies are inherently predatory to the consumer, it is important to identify potential filters that can be put into place that might protect social media users from exposure to posts about these scams. Performing text analysis to determine what language is more frequently used in posts about multilevel marketing companies versus posts about legitimate small businesses is a crucial first step in the process.

### **Data Source and Attribute Definitions**

Between all the social media platforms that are widely used, Twitter was selected as the focus. Twitter has a limit of 280 characters per post, which means users need to get to the point quite quickly. This platform was the origin of the current use of the hashtag, which is a way of marking words or phrases to make post content easily searchable. Twitter also has an Application Programming Interface (API) that makes it easy to scrape data.

For this project, there were time constraints that made it challenging to effectively scrape Twitter through a self-built application. Instead, we used a bot to scrape Twitter for a series of hashtags, specifically, eleven known multilevel marketing companies and eleven known legitimate small businesses, or companies that offer a similar product or service as the multilevel marketing companies. The first few rows of data are pictured in the image below:

	A	B	C	D	E	F	G	H
1	Hashtag	Username	User handle	Date of posting	Text	Retweet count	Like count	MLM
2	#etsy	Pinup Artist • sfw	@__ryusart	Mon Mar 28 22:1	Someone suggest <a href="https://t.co/ayeY">https://t.co/ayeY</a>	0	3	0
3	#lularoe		@_aaashleey	Mon Mar 28 03:3	So good I had to	0	0	1
4					We got an agent 👉 @basedfishmafi #BFM #NFTs #N			
	#minted	Julz 🍷	@_AusJulz	Mon Mar 28 01:1	Can't wait to see	10	51	0
5	#poshmark	Candice Couture	@_candicecoutu	Tue Mar 29 03:2	So good I had to	2	2	0
6					Great news! #20  <a href="https://t.co/30mV">https://t.co/30mV</a>			
	#minted	DoodleFriends.s	@_DoodleFrienc	Tue Mar 22 03:1	#NFTCommunity	0	2	0
7					#herbalife #skinc			
	#herbalife	MaFifi	@_Fifinky	Mon Mar 21 19:3	I'm placing order	0	0	1
8	#avon	Keshia Owen	@_GoldenKe	Mon Mar 28 00:2	So good I had to	0	0	1
9	#gofundme	Heather Buckley	@_HeatherBuck	Tue Mar 29 03:0	This post alerted	1	2	0
10	#gofundme	Heather Buckley	@_HeatherBuck	Mon Mar 28 18:1	:D Asking for sto	3	4	0
11	#gofundme	Heather Buckley	@_HeatherBuck	Mon Mar 28 00:4	Asking for some	1	1	0
12	#gofundme	Heather Buckley	@_HeatherBuck	Sun Mar 27 17:2	Asking for some	7	7	0
13	#gofundme	Di Devi	@_idksophie_	Mon Mar 28 13:3	JUST GOT ACC	10	18	0

The attributes of this data set are as follows:

- Hashtag - this is how the tweets were found within Twitter
- Username and User handle - these simply identify the posters
- Date of posting - this is both the date and time of posting
- Text - herein lies the corpus of documents, this is the content of the tweets
- Retweet count - this is the number of other users that re-posted this tweet to their personal account
- Like count - this is the number of other users that liked the tweet
- MLM - this is a binary target variable that is used to categorize the content of the tweets as being about a multilevel marketing company or a legitimate business

## Data Cleaning

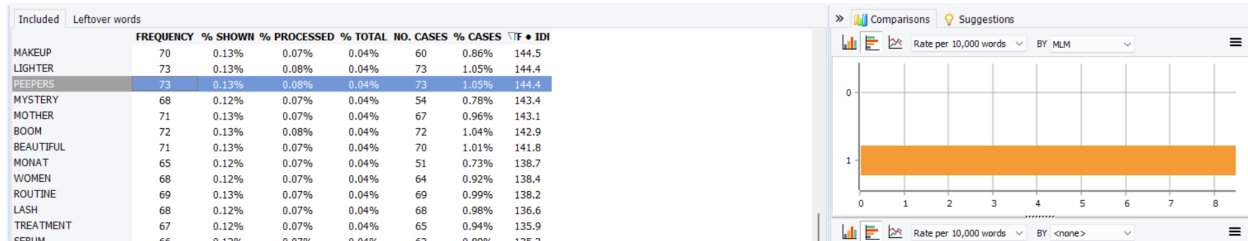
The first steps that were taken to clean this dataset were to remove the emojis, html entities, html tags, and non-printing characters. Twitter users tend to add a lot of extraneous symbols and spacing to call attention to their tweets, so removing this noise was a necessary step.

## Exploratory Analysis

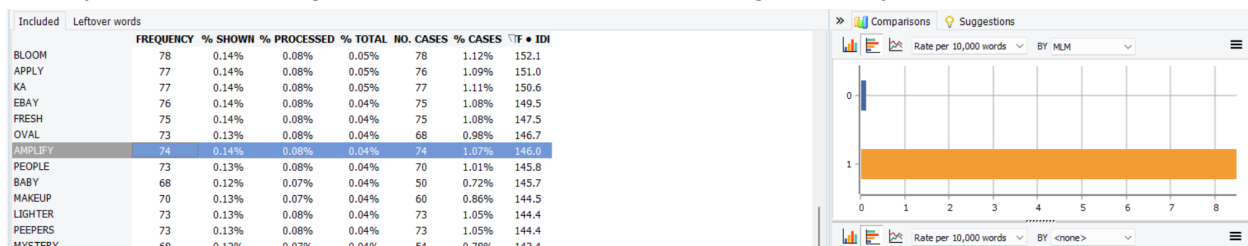
The initial exploration of the text was performed in Wordstat as a content categorization with the tweet being the content and the binary classes of multilevel marketing company or legitimate business being the categories. As one might imagine, there was an enormous amount of words and phrases selected through the auto-settings on the application. There were some obvious

classifiers, like “Etsy” and “Etsy shop” which clearly belong in one of the binary categorizations since “Etsy” was one of the hashtags used to scrape Twitter. Setting aside those n-grams, there were several strong classifiers.

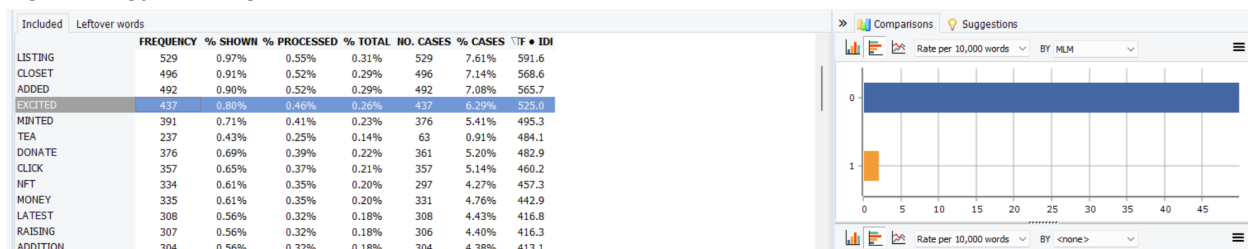
Peepers (synonym for eyes) is a strong indicator of a multilevel marketing company.



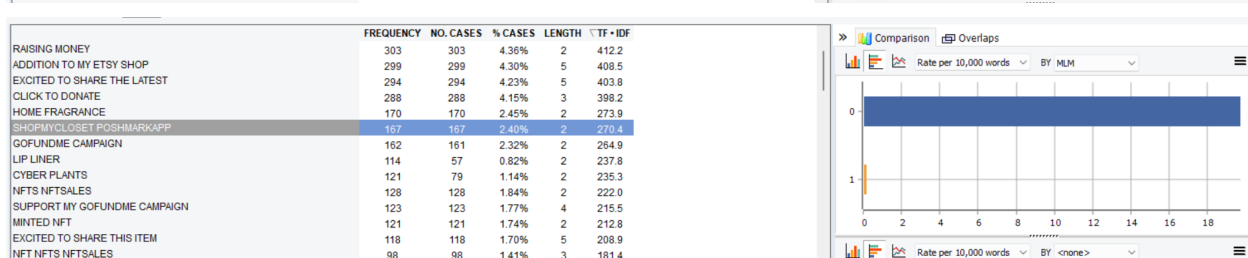
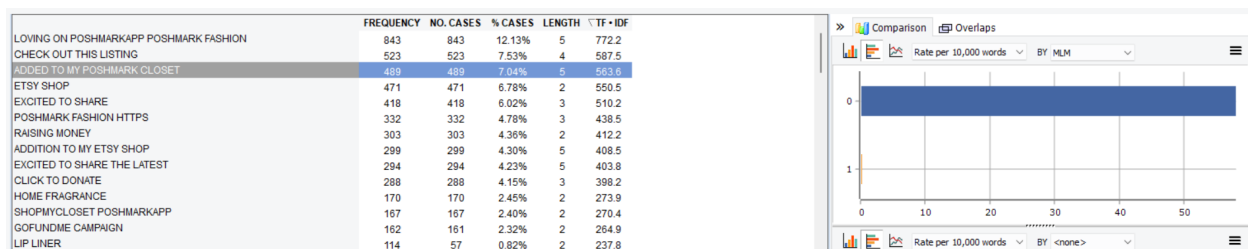
Amplify is another strong indicator of a multilevel marketing company.



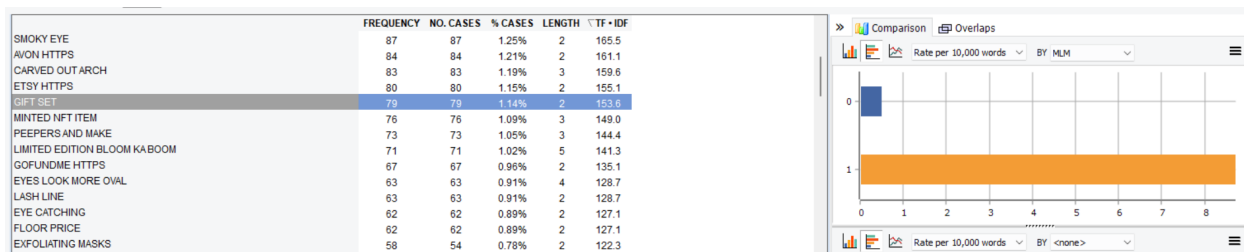
Excited is a strong indicator of a legitimate business, which is a surprise, since so much of the common language used by direct sellers of multilevel marketing companies is biased toward high energy and high enthusiasm.



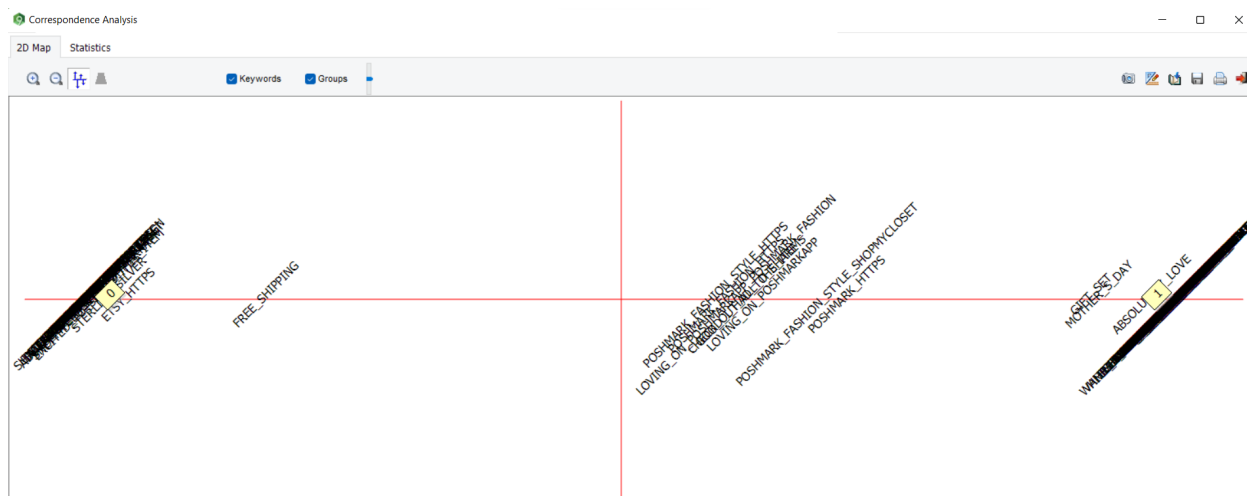
Pairing “Poshmark” with “Closet” is a strong indicator of a legitimate business. This is not initially a surprise since Poshmark was one of the hashtags for legitimate business.



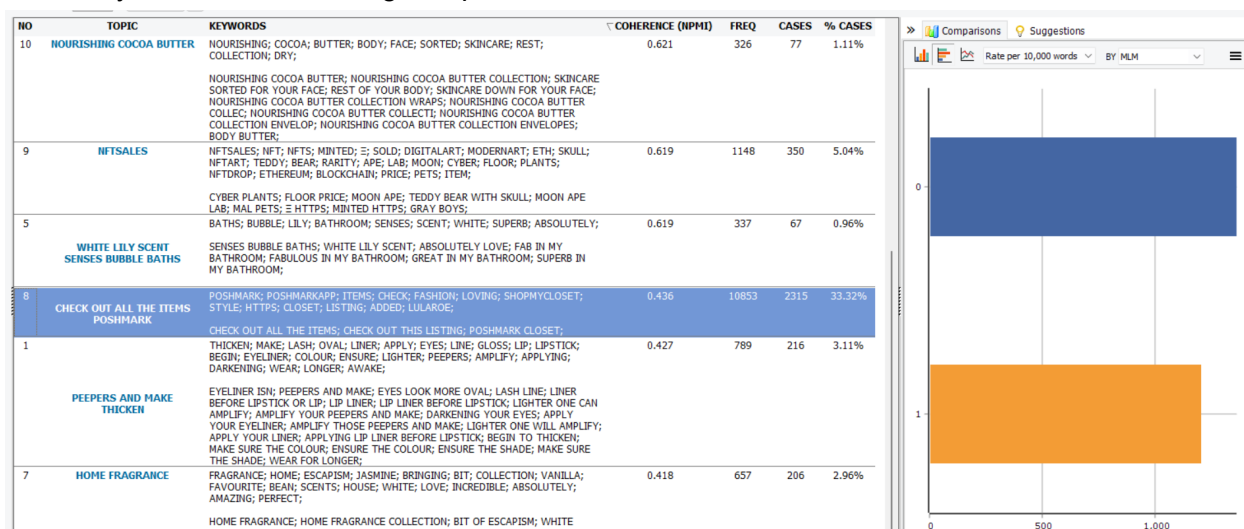
“Gift set” is a strong indicator of a multilevel marketing company.



Putting the phrases into a spectrum between the two categories shows interesting results. For example, “Free shipping” tends to appear in more content about legitimate businesses, while, again, “Gift set” is an indicator of a multilevel marketing company, as is the phrase, “Absolutely love”



An interesting result from this exploration is that Poshmark was originally listed as a legitimate business, which it is, but apparently content about Poshmark appears in tweets about both legitimate businesses and multilevel marketing companies. At this point, I looked back into Poshmark to find out that some of the direct sales on Poshmark are for products and services offered by multilevel marketing companies.



## Modeling

The model that was selected for this project was to use a Passive Aggressive Classifier with a TFIDF Vectorizer.

The TFIDF Vectorizer consists of two parts: term frequency and inverse document frequency. This is a measure of the terms that make a big impact when they appear in the documents within the corpus. To measure term frequency alone means that across the corpus, a term appears often, which may be significant, but may not be. To measure only inverse document frequency means to give weight to rarer words in documents. TFIDF measures a combination of the two, both the important terms, and the documents in which they appear.

The Passive Aggressive Classifier is an iterative algorithm that learns incrementally. The algorithm responds passively to correct classifications (it lets the correct classifications remain as they are), and aggressively to incorrect classifications (the adjustment to the classification method occurs upon incorrect classification results).

The training set is 20% of the total dataset and the test set is the remaining 80%.

```
✓ [12] x_train,x_test,y_train,y_test = train_test_split(t_df['Text'], labels, test_size=0.2, random_state=7)
0s

✓ [13] tfidf_vectorizer = TfidfVectorizer(stop_words='english', max_df=0.7)
0s

✓ [14] tfidf_train = tfidf_vectorizer.fit_transform(x_train)
0s      tfidf_test = tfidf_vectorizer.transform(x_test)

✓ [15] pac = PassiveAggressiveClassifier(max_iter = 50)
0s      pac.fit(tfidf_train,y_train)

      PassiveAggressiveClassifier(max_iter=50)

✓ [16] y_pred = pac.predict(tfidf_test)
0s

✓ [17] score = accuracy_score(y_test,y_pred)
0s      print(f'Accuracy: {round(score*100,2)}%')

      Accuracy: 91.8%
```

The accuracy of the model is 91.8%, which is a very high level of accuracy for such a small dataset and with minimal feature engineering.

## Model Evaluation

One of the model evaluation methods employed for this project was a confusion matrix:

```
✓ [19] confusion_matrix(y_test,y_pred, labels=[1,0])  
0s  
array([[657, 18],  
       [ 96, 619]])
```

We can see in the array that out of all the tweets about multilevel marketing companies in the test set, 657 of them were correctly categorized as multilevel marketing tweets. Of all the tweets about legitimate companies, 619 were correctly categorized as legitimate companies. Ninety-six of the legitimate companies were categorized as multilevel marketing companies, and 18 of the multilevel marketing companies were categorized as legitimate companies.

If our purpose is to prevent vulnerable individuals from being exposed to and potentially scammed by content about multilevel marketing, this skew towards over-predicting twitter posts as multilevel marketing content is a positive.

### Recommendations for Future Steps

If further analysis is performed, the first step would be to scrape other social media platforms like Facebook, Instagram, and LinkedIn. These platforms do not restrict users to the character limit that Twitter enforces, so the text analysis will likely uncover more patterns within the text that can be used to identify strategies commonly used by multilevel marketing companies.

For a more in-depth analysis and higher prediction accuracy, I recommend using a dictionary method, such as dictionary-based sentiment analysis. These analytical methods can identify brand personality which can be used to more immediately recognize these types of posts.

Finally, if a much larger dataset was collected, the posting behavior of the content creators could potentially be used in conjunction with the text analysis as a predictor of multilevel marketing content. From my own experience as a social media user, I recognize that multilevel marketing posts are published on a schedule and with a high frequency of post output. This pattern of behavior could be analyzed and incorporated into the model to create highly accurate predictions of multilevel marketing content.

### Works Cited

Botster Team. "Twitter Hashtag Scraper." *botster*, 2022, <https://botster.io/bots/twitter-hashtag-scraper>.