

APCV361 Spring 2022
 Dr. Li Xu
 Team Generic:
 Cortland Diehm
 Katelyn Griffith
 Jake Sherwood

Words of Engagement (APCV 361 Final Project)

Introduction

The popularity of social media websites such as Twitter, Instagram, and Facebook have become an irreplaceable mechanism for organizations and individuals alike to communicate ideas and thoughts to a wide range of people, both in scope and character. Social media, and Twitter in this case, can serve as a powerful apparatus for communicating news and events to users. Due to the potential to influence or liaise with a large audience, the factors that go into the popularity of a single posted tweet is invaluable information.

The nature of a tweet is to be succinct and concise, which makes the wording of these posts more important to maximize the amount of engagement with an audience. In this study, we analyzed the amount of engagement a tweet has from the University of Arizona's College of Agriculture and Life Sciences account "@UArizonaCALS" based on the wording found within the tweet itself. This was done by collecting tweets from the account and sorting them into categories based on the specific words found in the tweet. The number of favorites, known as "likes," was used as a measure to analyze the engagement with the audience. Those likes were associated with each category and analyzed through graphing and comparison tests in this study.

Project Design and Development

For the design of this project, the data was extracted and sorted via Python and then graphed/analyzed using tools provided by the Matplotlib package. For the extraction of twitter data, we used tools available through Python and the Twitter Developer API. Using Tweepy, an open-source Python package, the data from the "@UArizonaCALS" twitter account (around 1100 tweets) was extracted and saved to a separate CSV file. In this CSV file, each tweet was saved along with its corresponding 'like' count.

The categories in which the tweets were separated into are based on multiple text files (with each text file being their own category of words). Each file holds their own unique list of keywords that correspond to that category. For example, if the category of tweets we want to analyze falls under "Climate", then the tweets that contain keywords associated with climate found in the category's text file were separated into a separate group. Once all of the tweets were separated into their rightful categories, the data for each category was saved to their own CSV file, which served as the basis for the analysis and comparisons later.

Once the tweets were sorted into their appropriate categories, we implemented the Matplotlib package for the visualizations and models. The data collected was graphed onto various models, including bar plots and Wordcloud visualizations. Each category was graphed onto a bar plot, illustrating all of the words within that category and their corresponding like counts. Along with the graphs for each category, two more graphs comparing the total like counts of each category and the top ten most-liked words respectively were created to further illustrate which words are associated with increased popularity. With the results of these graphs, we were able to analyze which categories and which words were most likely to have an increased like count in a tweet.

Step One: Gathering Tweets to Use

The first process was gathering tweets to use for our project. We started by creating an API instance using `tweepy`. Keys were stored in a separate `keys.py` file to enhance protection and allow us to interchange our keys to access the API.

```
In [1]: 1 # -----
2 # Create API instance
3 # -----
4 # Author: Jake
5
6 import tweepy
7
8 def create_api(key, secret, authtoken, tokensecret):
9     auth = tweepy.OAuthHandler(key, secret)
10     auth.set_access_token(authtoken, tokensecret)
11
12     # Create API object
13     api = tweepy.API(auth, wait_on_rate_limit=True)
14
15     # Test authentication
16     try:
17         api.verify_credentials()
18         #print('You\'ll Succeed, Eventually') # Secret confirmation message
19         #api.update_status('You\'ll Succeed, Eventually.') # Debugging only
20     except:
21         print('There was an error during authentication\n')
22     return api # Returns API object
23
24 # Keys for authentication
25
26 from keys import consumer_key as key
27 from keys import consumer_secret as secret
28 from keys import access_token as authtoken
29 from keys import access_token_secret as tokensecret
30
31 api = create_api(key, secret, authtoken, tokensecret)
```

Next we defined code to collect tweets from `UArizonaCALS` to store into an array. This included the timestamp of when the tweet was created, the user (in case we wanted to include others), and the full text of the tweet. Debugging code was provided at the end to verify the number of tweets we were able to retrieve.

```

In [3]: 1 # -----
2 # Retrieve Tweets
3 # Run this if you want to repopulate the tweets List
4 # -----
5 # Author: Jake
6
7 def get_tweets(tweeter, count=20):
8     # Input:  Twitter handle of user, how many tweets to count (default 20 to prevent Timeout)
9     # Output: Array of tweets' timestamp, user, and full text
10    arr = []
11
12    # Define what kind of tweets you want selected from timeline
13    timeline = tweepy.Cursor(api.user_timeline,
14                             screen_name=tweeter,
15                             tweet_mode='extended', # full text
16                             exclude_replies=True, # reply tweets
17                             include_rts=False # retweets
18                             ).items()
19
20    for status in timeline:
21        if count > 0:
22            arr.append([status.full_text, status.favorite_count, status.created_at])
23            count -= 1
24        else:
25            break
26    return arr
27
28 def print_tweets(arr, count=2):
29     # Input:  Array of tweets, how many tweets to print (default 2)
30     # Output: None
31     # Display: Tweet counter, timestamp, user, and full text of tweet
32     x = count
33     for tweet in arr:
34         if count > 0:
35             print(str(x-count+1), tweet[0], tweet[1], tweet[2], sep='\t')
36             count -= 1
37         else:
38             break
39     #print('\nTotal of', x, 'tweets displayed from', Len(arr))
40
41 tweeter = 'UArizonaCALs'
42 limit = 5000 # Limits the number of tweets to retrieve
43 tweets = get_tweets(tweeter, count=limit)
44 #try:
45     #print_tweets(tweets) # Debugging only
46 #except:
47     #print('There was an error printing the tweets.')

```

To avoid timeout from the API, we saved what we collected to a CSV file so we can pull the data from the CSV without having to re-run the above code.

```

In [4]: 1 # -----
2 # Save tweets to text file so we can run code
3 # without running API commands again. (Avoids Timeout)
4 # -----
5 # Author: Jake
6
7 import pandas as pd
8
9 def write_to_csv(tweets, filename='output.csv'):
10     # Input:  Array of tweets
11     # Output: None
12
13     # Create dataframe
14     col = ['Tweet', 'FavoritesCount', 'Timestamp']
15     data = []
16     try:
17         for tweet in tweets:
18             tweet[0] = tweet[0].replace('\n', ' ') # Replace newline characters
19             tweet[0] = tweet[0].replace(',', ' ') # Replace commas
20             data.append(tweet)
21         df = pd.DataFrame(data, columns=col)
22         df.to_csv(filename, sep=',')
23         #print('Tweets written to file', filename)
24     except:
25         print('There was an error in writing the tweets to csv file.')
26
27 write_to_csv(tweets, 'tweets.csv')

```

Step Two: Creating Classes for Categories and Words

We needed to keep track of the like count for each word in each category. To do this we implemented dictionaries to look for specific words to each category and keep track of the favorites count. The next step was to define what words we wanted to search for and what to categorize the words under. Since this list can change as the project evolves, we elected to use simple TXT files to store the words for each category. These files will be included in the submission for this project.

```

In [1]: 1 # -----
2 # Defining classes for categories and words
3 # -----
4 # Author: Katelyn
5
6 # class Cat is the categories
7
8 class Cat():
9     def __init__(self, name):
10         self.words = {} # Dictionary of words
11         self.name = name # Name of category
12         self.totalLikes = 0 # Total Likes for category
13
14     def getWords(self):
15         return self.words
16
17     def getName(self):
18         return self.name
19
20     def getTotalLikes(self):
21         return self.totalLikes
22
23     def addWord(self, key, w):
24         self.words[key] = w
25
26     def inclLikes(self, wLikes):
27         self.totalLikes += wLikes
28
29     def __str__(self):
30         return f'Words: {self.words}'
31
32 # class Word is the individual word
33
34 class Word():
35     def __init__(self, word, cat='none'):
36         self.word = word # Word
37         self.likeCount = 0 # Total Likes for word
38         self.myCategory = cat # Category the word belongs to ~ Jake
39
40     def setlikeCount(self, likes):
41         self.likeCount += likes
42
43     def getWord(self):
44         return self.word
45
46     def getlikeCount(self):
47         return self.likeCount
48
49     def __str__(self):
50         return f'Words: {self.word}\nLike Count: {self.likeCount}'
51
52     def getCategory(self):
53         # Return which category the word belongs to ~ Jake
54         return self.myCategory
55
56 # Verify classes were created successfully
57 #print('Category and Word classes created.')
58
59 # -----
60 # Creating dictionary of categories
61 # -----
62 # Author: Katelyn
63
64 # Load file
65 def loadFile(filepath, display=False):
66     try:
67         f = open(filepath, 'r', encoding="utf-8")
68         plist = f.readlines()
69         f.close()
70         #if display != False:
71             #print(len(plist), ' Lines read from ', filepath, '.', sep='')
72         return plist
73
74     except FileNotFoundError:
75         print("File", filepath, "could not be located. Halting...")
76         exit()
77
78 # Create category list
79 def buildCategories(catlist):
80     dictionary = {}
81
82     for category in catlist:
83         c = Cat(category)
84         dictionary[c.getName()] = c
85
86     return dictionary
87
88 categories = ['Climate',\
89 'Employment',\
90 'Congratulations',\
91 'Research',\
92 'Scholarships',\
93 'WomenInStem']
94 dictionary = buildCategories(categories)
95 #print('Dictionary created.')
96

```

```

97
98 # Create word List
99 def buildWords(wordFiles, dictionary):
100     categories = list(dictionary.keys())
101
102     for f in range(len(wordFiles)):
103         fl = loadFile(wordFiles[f], display=True)
104
105         # Add words to dict
106         for w in range(len(fl)):
107             q = Word(fl[w][:~1], cat=categories[f])
108             dictionary[categories[f]].addWord(q.getWord(), q)
109
110     return dictionary
111
112 # Add words to Dictionary
113 filenames = ['climate.txt',\
114             'employment.txt',\
115             'congratulations.txt',\
116             'research.txt',\
117             'scholarships.txt',\
118             'women_in_stem.txt']
119 dictionary = buildWords(filenames, dictionary)
120 #print('Words added to dictionary.')
121
122 # Open the tweets file
123 t = loadFile('tweets.csv', display=True)
124
125 for i in range(1, len(t)):
126     t2 = t[i].split(',') # Delimiter in CSV file
127     t3 = t2[1].lower().split() # Lowercase and split all words
128
129     # Go through each tweet
130     for word in t3:
131
132         # Go through each section
133         for key in dictionary:
134             wordlist = dictionary[key].getWords() # Get words for section
135
136             # Go through each word for each section
137             for value in wordlist:
138
139                 # Update Like counts
140                 if value in word:
141                     dictionary[key].words[value].setlikeCount(int(t2[2]))
142                     dictionary[key].inclLikes(int(t2[2]))
143                     break
144 #print('Dictionary updated with total Likes in each category:')
145
146 cat_likes = []
147 for cat in dictionary:
148     cat_likes.append(dictionary[cat].gettotalLikes()) #Added cat_list for easier access to Likes for each category ~Cortland
149     #print(f" {cat+'':<16} {dictionary[cat].gettotalLikes():>5}")

```

Step Three: Visualizations of the Data

Next we took the data created in the previous cells and started to graph them for comparisons. The first visualization we created was a bar graph comparing each category by the number of likes associated with it. Next, when comparing like counts within each category of words, we created six bar-subplots comparing the like data for each individual word. For these plots, any word that had zero likes associated with it were dropped from the visualization to avoid clutter and streamline the models.

Finally, the final graph created in this cell displayed the Top Ten most-liked words, no matter which category they came from. Each word is color-coded, corresponding to the category in which it came.

```

In [2]: 1 #-----
2 # Creating visualizations from the data
3 #-----
4 # Author: Cortland
5
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 import numpy as np
9
10 # Bar Graph of Categories and Corresponding Likes
11 cats = list(dictionary.keys())
12 colors = ['mediumblue', 'seagreen', 'gold', 'darkorange', 'firebrick', 'rebeccapurple']
13
14 # Total Likes in each category
15 plt.figure(figsize=(800/72,400/72))
16 plt.grid(True, alpha=0.5)
17 plt.rcParams['axes.axisbelow'] = True
18 plt.ylabel('Likes', fontsize=14)
19 plt.title('Total Likes by Category (Figure 1)', fontsize=20)
20 plt.bar(cats, cat_likes, tick_label=cats, color=colors)
21 #plt.yticks(np.linspace(0,1500, 31))
22
23 # -----
24 # Total Likes for words in each category
25 # -----
26
27 fig, axs = plt.subplots(6, 1, figsize=(800/72, 6*300/72)) # 6 rows, 1 columns
28 i = 0
29 title_count = ['2', '3', '4', '5', '6', '7'] #Added title_count so each subplot iteration is Labeled by number
30 for ax in axs.flat:
31     top20 = []
32     words = dictionary[cats[i]].getWords()
33     #x, y = [], []
34     for word in words:
35         x = dictionary[cats[i]].words[word].getlikeCount()
36         if x > 0:
37             top20.append((word, x))
38
39     # Sort toptwenty array by number of Likes
40     top20 = sorted(top20, key=lambda a: a[1], reverse=True)
41
42     # Keep the top 20 words
43     if len(top20) > 20:
44         top20 = top20[:20]
45
46     # Sort to arrays
47     x, y = [], []
48     for w in top20:
49         y.insert(0, w[0])
50         x.insert(0, w[1])
51
52     # Draw the graph
53     ax.grid(True, alpha=0.5)
54     #ax.set_xticks(np.linspace(0, 1000, 21))
55     ax.set_title('Category: '+cats[i] + ' (Figure '+title_count[i]+' )', fontsize=20)
56     ax.set_xlabel('Likes')
57     ax.barh(y, x, color=colors[i])
58     i += 1
59
60 fig.tight_layout()
61
62 #-----
63 # Histogram of most Liked words
64 #-----
65
66 # Finding the top ten most Liked words
67 topten = []
68 i = 0
69 for c in categories:
70     # categories were defined in Katelyn's code
71     for word in dictionary[c].words:
72         y = dictionary[c].words[word].getlikeCount()
73         x = word
74         z = colors[i]
75         topten.append((x, y, z))
76     i += 1
77 # Sort the words by most Liked
78 topten = sorted(topten, key=lambda a: a[1], reverse=True)
79 # Keep the top ten words
80 topten = topten[:10]
81
82 # Make the graph
83 plt.figure(figsize=(800/72,400/72))
84 plt.grid(True, alpha=0.5)
85 plt.rcParams['axes.axisbelow'] = True
86 plt.xlabel('Likes', fontsize=14)
87 plt.xticks(np.linspace(0, 1000, 21))
88 plt.title('Top Ten Most Liked Words (Figure 8)', fontsize=20)
89
90 # Sort tuple elements into arrays
91 x, y, z = [], [], []
92 for word in topten:
93     x.insert(0, word[0])
94     y.insert(0, word[1])
95     z.insert(0, word[2])
96 plt.barh(x, y, color=z)

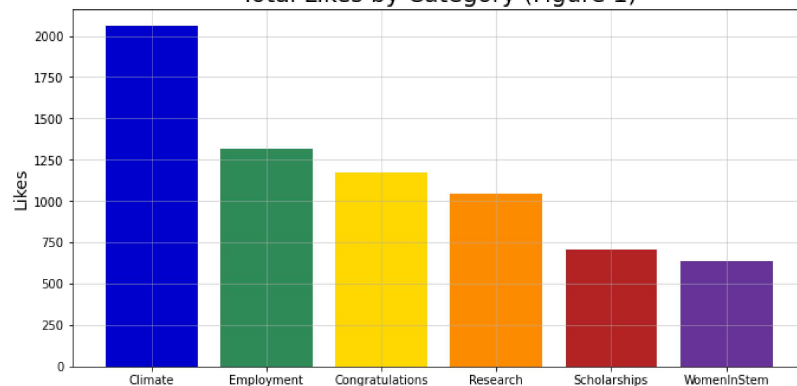
```

```

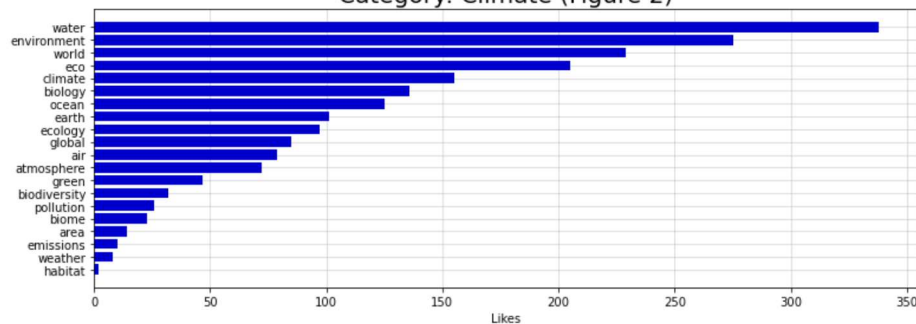
97
98 # Add the Legend before showing the graph
99 h = []
100 i = 0
101 for c in cats:
102     line1, = ax.plot(0, label=c, color=colors[i])
103     h.append(line1)
104     i += 1
105 plt.legend(handles=h)
106
107
108 plt.show()

```

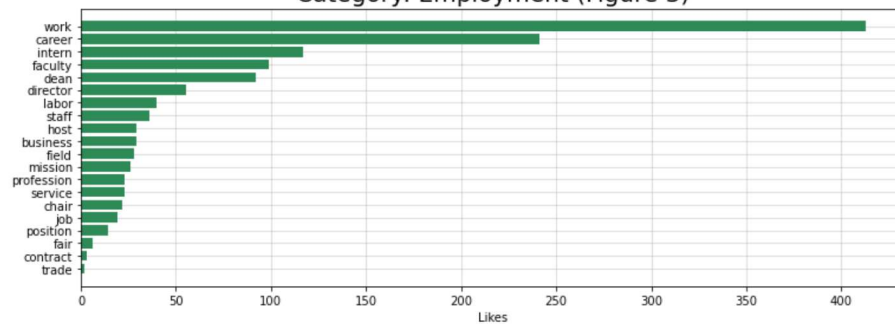
Total Likes by Category (Figure 1)



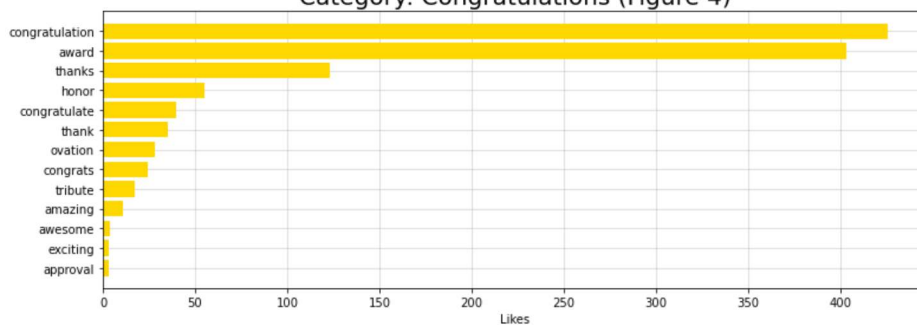
Category: Climate (Figure 2)



Category: Employment (Figure 3)

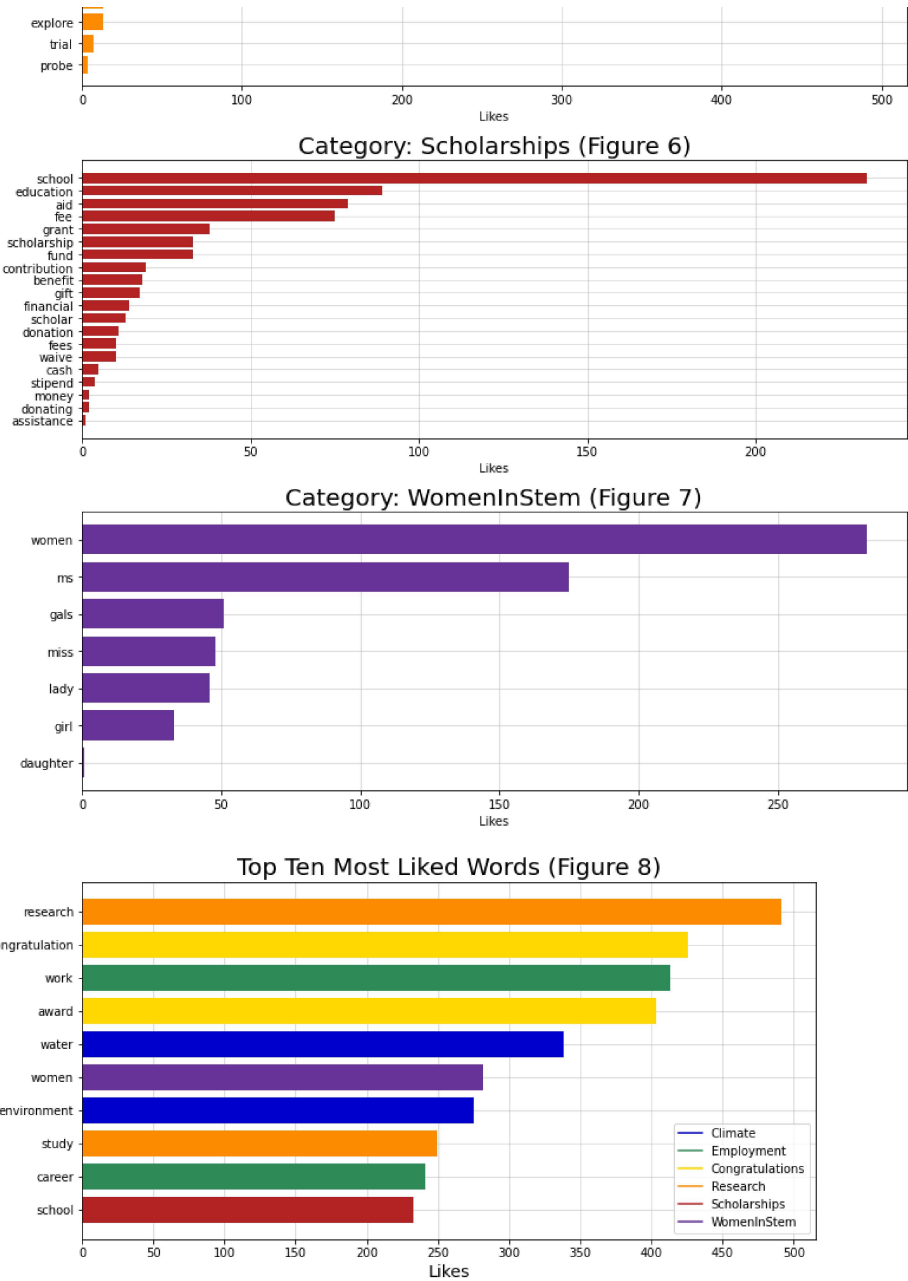


Category: Congratulations (Figure 4)



Category: Research (Figure 5)





Step Four: Creating a Word Cloud

To better visualize the occurrence of words and frequency of likes associated to tweets with such words, we created a word cloud; Essentially a mass of words that scales a word based on frequency, thereby showing an image of different words at different scaled sizes. The first word cloud was created from data from all categories, producing a good model for visualizing the most popular words in the dataset. The subsequent word clouds were based on data gathered from each individual category, showing the most popular words within their group.

```

In [4]: 1 # -----
2 # Creating a Word Cloud
3 # -----
4 # Author: Jake
5
6 import matplotlib.pyplot as plt
7 import wordcloud as wc
8 import numpy as np
9 import random
10 from PIL import Image
11
12 """
13 Theory: To make this code work for our purposes, we'll need to
14 create a string of the words found in each category and append
15 each word by the number of times it occurs.
16
17 x Goal: Generate graph by frequency of words in all categories.
18 Goal: Generate graph by likes for a word in all categories.
19 Ex. Goal: Use different images for word clouds.
20
21 """
22
23 categories = dictionary.keys()
24 string = ''
25 for c in categories:
26     for w in dictionary[c].getWords():
27         n = dictionary[c].words[w].getlikeCount()
28         string += (w+' ')*n
29
30 def word_cloud(string, filename=None, size=(800/72, 800/72)):
31     """
32     Input: A string and an image file you would like to use. Please use a file that ends in .jpg
33     Output: Wordcloud graph
34     """
35     m = None
36     if filename != None:
37         m = np.array(Image.open(filename))
38
39     cloud = wc.WordCloud(max_font_size=1000,\
40     collocations=False,\
41     mask=m,\
42     background_color='white').generate(string)
43     plt.figure(figsize=size)
44     plt.imshow(cloud, interpolation='bilinear')
45     plt.axis('off')
46     plt.show()
47
48 # Display word cloud for all categories
49 word_cloud(string, filename='lxu.jpg')
50

```



Results and Discussion

Category Comparison Results

The first of the results to be analyzed are the results found in Figure 1, which displays the total likes by category. From our initial data comprising of 1,152 tweets from the University of Arizona CALS twitter account, the category of "Climate" garnered the most amount of likes by a considerable margin. The Climate category had 2,060 likes associated with it, whereas the category with the least amount of likes was the "Women In Stem" category, which garnered only 636 likes total. For the next group of data gathered from the results, we will analyze Figures 2 through 7. These figures display and compare the number of likes associated with each word within their respective categories, allowing for a more concentrated examination of the data.

Category Subplot Results

1) **Figure 2 (Climate):** Figure 2 has the widest variety of data (with 2,060 likes associated with the category), and the most liked word was "Water", which has 340 total likes associated with it. The word with the least amount of likes associated with it was the word "Habitat", which only had 5 likes associated with it.

2) **Figure 3 (Employment):** Figure 3 has a total like count of 1,317, but the data is most skewed with the top two words within the category: "Work", the most popular word with 440 likes, and "Career" with 240 likes. The least-liked word in the category, "Trade," only had 2 likes associated with it.

3) **Figure 4 (Congratulations):** Figure 4 has a total like count of 1,172, and like Figure 3, the likes are mostly concentrated with the top two words: "Congratulation", which has 450 likes, and "Award", which has 408 likes. There are three words tied for the least-liked word, and they are "Awesome", "Exciting", and "Approval", which all have only 6 likes associated with them.

4) **Figure 5 (Research):** Figure 5 has a total like count of 1,041, and it also follows the trend of likes being concentrated around the top two words in the category. The top two words are "Research", which has 490 likes, and "Study", which has 245 likes. The least-liked word is "Probe", which has 8 likes.

5) **Figure 6 (Scholarships):** Figure 6 has a total like count of 706, and it has a slightly more varied distribution than Figures 3-5. The most liked word is "School", which has 240 likes. The least-liked word has hardly any likes associated with it, at only 2 likes.

6) **Figure 7 (Women In Stem):** Figure 7 displays the data from the "WomenInStem" category, which has the least amount of likes associated with it (636 likes). The data is concentrated around the top two words: "Women" at 282 likes, and "Ms" at 175 likes. The word with the least amount of likes is "Daughter", with only 1 like associated with it.

Top Ten Most Liked Words

From Figure 8, we can easily see the top ten most-liked words in the entire dataset. The word with the most likes associated with it within all of the tweets gathered was "Research", at 490 total likes. The words in the top ten were not skewed in the favor of any particular category, with all but two categories having multiple words appear in the graph.

Conclusion

From the data gathered and analyzed within the program, we can make predictions on whether a tweet from the @UArizonaCALS twitter account will garner likes or not. From our results, the safest tweet the account can make if it is seeking engagement is to add words pertaining to the category of "Climate". The Climate category was by far the most-liked category, with a wider variety of words with considerable like counts than the other categories. If seeking to maximize engagement, it would be better in general to avoid words associated with the "Women In Stem" category, which was the category with the least amount of likes associated with it. Although the "Climate" category had the most likes associated with it out of all of the categories, if the account wanted to maximize the amount of engagement they would publish a tweet containing the word "Research", which was the word with the most likes associated with it regardless of category.

While the results of this is not a foolproof way of increasing engagement for @UArizonaCALS, it can provide a rough idea of what topics will be more well-received. There are a multitude of factors that lead to an increased amount of engagement on a tweet, but it is important to analyze which tweets did well in the past and examine which topics and words contributed to that engagement.

References

Matplotlib Development Team. (n.d.). Matplotlib API Overview. DevDocs. Retrieved May 6, 2022, from <https://devdocs.io/matplotlib-3.1/> (<https://devdocs.io/matplotlib-3.1/>)

Rosslein, Joshua. (2022). Tweepy documentation. Tweepy Documentation - tweepy 4.9.0 documentation. Retrieved May 6, 2022, from <https://docs.tweepy.org/en/stable/> (<https://docs.tweepy.org/en/stable/>)