

## Statistics/Plots

We have analyzed the frequencies and distributions of the existence and existence\_confidence columns. The distribution of the scores for existence\_confidence are most frequent around 0.6 to 0.75 and highest at 1.0. We have an imbalance where there are more tweets believing in climate change compared to those that don't. This is further confirmed by the existence distribution histogram. The top twenty words do not change much before and after preprocessing. One interesting pattern is that the words "climate", "change", and "global" stay as the top twenty either way. One interesting finding from the word clouds is that the words pertaining to links (bit.ly, tinyurl) are showing up as the most impactful. After preprocessing, words/phrases such as greenhouse, conscious, and oil spill show up as the most impactful. This confirms that preprocessing helped in bringing out the most relevant words in the analysis.

```
In [ ]: 1 import numpy as np
        2 import pandas as pd
        3 import spacy
        4 from nltk.stem import PorterStemmer
        5 import matplotlib.pyplot as plt
        6 from itertools import chain
        7 from wordcloud import WordCloud, STOPWORDS
        8
        9 !pip install xlrdd==1.2.0
```

Looking in indexes: <https://pypi.org/simple>, (<https://pypi.org/simple>,) <https://us-python.pkg.dev/colab-wheels/public/simple/> (<https://us-python.pkg.dev/colab-wheels/public/simple/>)  
Requirement already satisfied: xlrdd==1.2.0 in /usr/local/lib/python3.7/dist-packages (1.2.0)

```
In [ ]: 1 import cufflinks as cf
        2 cf.go_offline()
        3 cf.set_config_file(offline=False, world_readable=True)
```

```
In [ ]: 1 from google.colab import drive
        2 drive.mount('/content/gdrive')
```

Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive.mount("/content/gdrive", force\_remount=True).

```
In [ ]: 1 data = pd.read_excel('/content/gdrive/My Drive/CMPE257-ML-Group-13/glob
2 data.head()
```

```
Out[3]:
```

	tweet	existence	existence.confidence
0	Global warming report urges governments to act...	Yes	1.0000
1	Fighting poverty and global warming in Africa ...	Yes	1.0000
2	Carbon offsets: How a Vatican forest failed to...	Yes	0.8786
3	Carbon offsets: How a Vatican forest failed to...	Yes	1.0000
4	URUGUAY: Tools Needed for Those Most Vulnerabl...	Yes	0.8087

```
In [ ]: 1 data['word_count'] = data['tweet'].apply(lambda x: len(x.split(" ")))
2 data.loc[data['existence'] == 'Y', 'existence'] = "Yes"
3 data.loc[data['existence'] == 'N', 'existence'] = "No"
4 data.dropna()
5 print(data['existence'].value_counts())
6 data.head()
```

```
Yes    3111
No     1114
Name: existence, dtype: int64
```

```
Out[4]:
```

	tweet	existence	existence.confidence	word_count
0	Global warming report urges governments to act...	Yes	1.0000	18
1	Fighting poverty and global warming in Africa ...	Yes	1.0000	8
2	Carbon offsets: How a Vatican forest failed to...	Yes	0.8786	12
3	Carbon offsets: How a Vatican forest failed to...	Yes	1.0000	12
4	URUGUAY: Tools Needed for Those Most Vulnerabl...	Yes	0.8087	11

```
In [ ]: 1 tweets = data["tweet"]
2 tweets = tweets.drop_duplicates()
3 tweets
```

```
Out[5]: 0 Global warming report urges governments to act...
1 Fighting poverty and global warming in Africa ...
2 Carbon offsets: How a Vatican forest failed to...
4 URUGUAY: Tools Needed for Those Most Vulnerabl...
5 RT @sejorg: RT @JaymiHeimbuch: Ocean Saltiness...
...
6085 @bloodless_coup "The phrase 'global warming' s...
6086 Virginia to Investigate Global Warming Scienti...
6087 Global warming you tube parody you will enjoy ...
6088 One-Eyed Golfer: Don't dare tell me about glob...
6089 man made global warming a hair brained theory ...
Name: tweet, Length: 5541, dtype: object
```

For Preprocessing we have used

```
In [ ]: 1 ### Source: https://spacy.io/usage/linguistic-features
2
3 def spacyPipeline(tweets):
4     ps = PorterStemmer()
5     nlp = spacy.load('en_core_web_sm')
6
7     preprocessed_tweets = []
8     for t in tweets:
9         doc = nlp(t)
10        filtered_tweet = []
11
12        for token in doc:
13            if (not token.is_stop) and token.is_alpha:
14                filtered_tweet.append(ps.stem(str(token)))
15
16        preprocessed_tweets.append(filtered_tweet)
17
18    return preprocessed_tweets
```

```
In [ ]: 1 def preprocess(tweets):
2     # Convert all to lowercase
3     tweets = [t.lower() for t in tweets]
4
5     # Process tweets through spaCy pipeline
6     tweets = spacyPipeline(tweets)
7
8     # Filter out words
9     tweets = [list(filter(lambda w: w != 'link', t)) for t in tweets]
10
11    # Remove words less than length 2
12    tweets = [list(filter(lambda w: len(w) > 2, t)) for t in tweets]
13
14    print(tweets)
15    return tweets
```

```
In [ ]: 1 def generateWordCloud(tweets):
2     allwords = " ".join(set(chain.from_iterable(tweets)))
3     wordcloud = WordCloud(width = 800, height = 800,
4                            background_color = 'white',
5                            stopwords = set(STOPWORDS),
6                            min_font_size = 10).generate(allwords)
7
8     plt.axis("off")
9     plt.tight_layout(pad = 0)
10
11    plt.figure(figsize = (7, 7), facecolor = 'white', edgecolor='blue')
12    plt.imshow(wordcloud)
13
14    plt.show()
```

```
In [ ]: 1 preprocessed_tweets = preprocess(tweets)

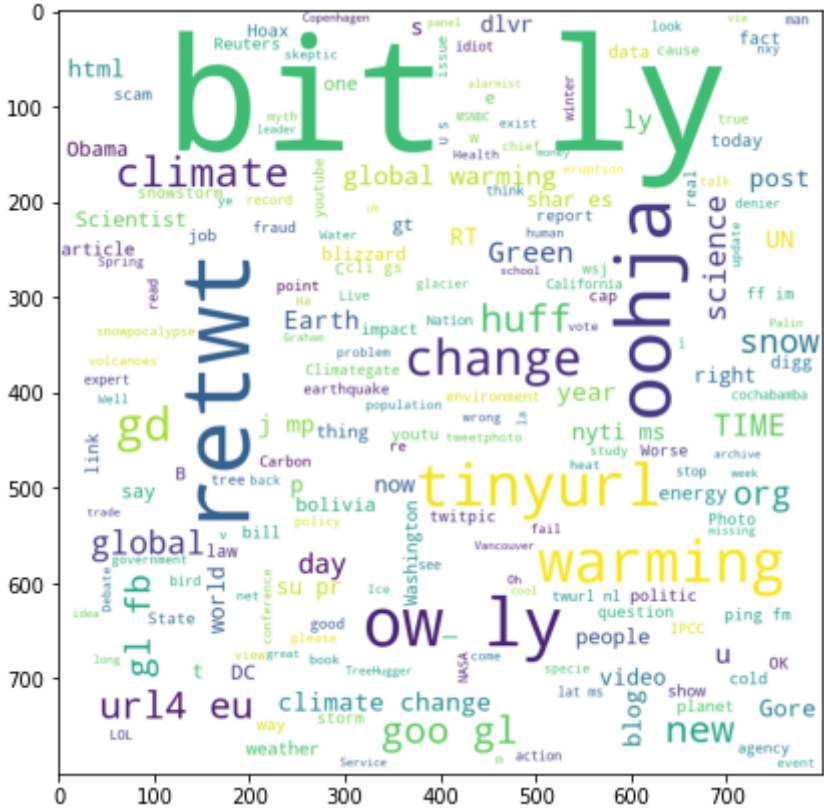
[['global', 'warm', 'report', 'urg', 'govern', 'belgium', 'world', 'face', 'increas', 'hunger'], ['fight', 'poverti', 'global', 'warm', 'africa'], ['carbon', 'offset', 'vatican', 'forest', 'fail', 'reduc', 'global', 'warm'], ['uruguay', 'tool', 'need', 'vulner', 'climat', 'chang'], ['ocean', 'salti', 'show', 'global', 'warm', 'intensifi', 'water', 'cyclicl'], ['global', 'warm', 'evid', 'messag', 'global', 'warm', 'denier', 'doubter', 'look'], ['migratori', 'bird', 'new', 'climat', 'chang', 'strategi', 'stay', 'home'], ['southern', 'africa', 'compet', 'limpopo', 'water', 'climat', 'chang', 'bring', 'higher', 'temperatur', 'south'], ['global', 'warm', 'impact', 'wheat', 'rice', 'product', 'apr', 'scarciti', 'water'], ['solv', 'global', 'warm', 'thing'], ['preliminari', 'analysisi', 'suggest', 'natur', 'contribut', 'far', 'global', 'warm', 'previous'], ['ecoton', 'climat', 'chang', 'popul', 'perspect'], ['climat', 'chang', 'blame', 'coastal', 'whale', 'migrat', 'dwindl', 'ventura', 'counti'], ['spring', 'storm', 'season', 'start', 'littl', 'late', 'year', 'global', 'warm'], ['govern', 'report', 'say', 'global', 'warm', 'caus', 'cancer', 'mental', 'ill'], ['earthday', 'global', 'warm', 'affect', 'patient', 'symptom'], ['wait', 'idea', 'natur', 'climat', 'chang', 'human', 'induc', 'global', 'warm'], ['epa', 'issu', 'report', 'clima
```

## VISUALIZATIONS

<https://towardsdatascience.com/a-complete-exploratory-data-analysis-and-visualization-for-text-data-29fb1b96fb6a> (<https://towardsdatascience.com/a-complete-exploratory-data-analysis-and-visualization-for-text-data-29fb1b96fb6a>)

<https://github.com/yrtnsari/Sentiment-Analysis-NLP-with-Python/blob/main/wordcloud.ipynb> (<https://github.com/yrtnsari/Sentiment-Analysis-NLP-with-Python/blob/main/wordcloud.ipynb>)

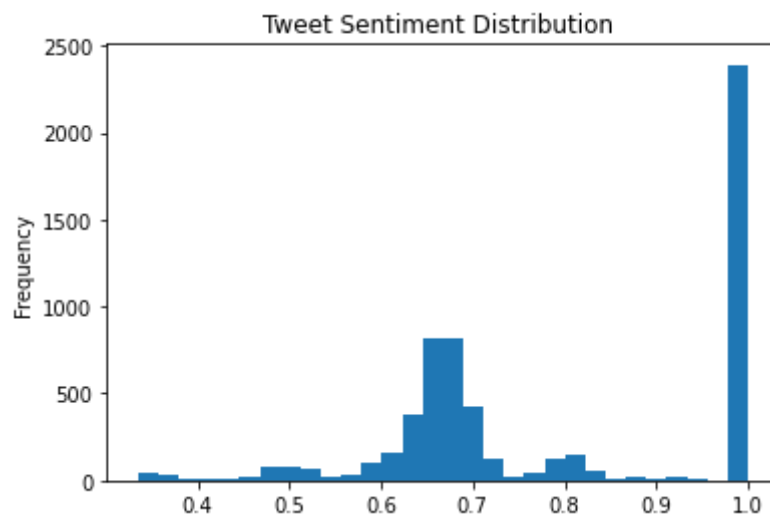
```
1 generateWordCloud([t.split(' ') for t in tweets])
2 generateWordCloud(preprocessed_tweets)
```





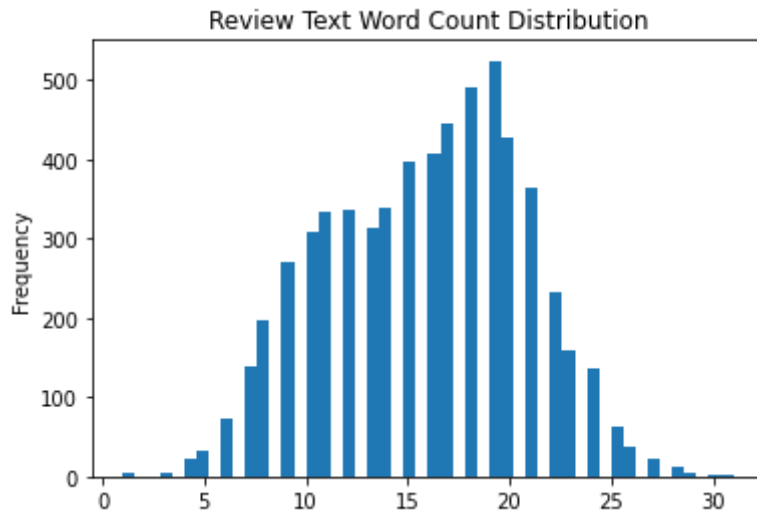
```
In [ ]: 1 data['existence.confidence'].plot(  
2       kind='hist',  
3       bins=30,  
4       title='Tweet Sentiment Distribution ')
```

Out[11]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7ff753e6f690>



```
In [ ]: 1 data['word_count'].plot(
        2     kind='hist',
        3     # xTitle = "Word Count",
        4     bins=50,
        5     title='Tweet Word Count Distribution')
```

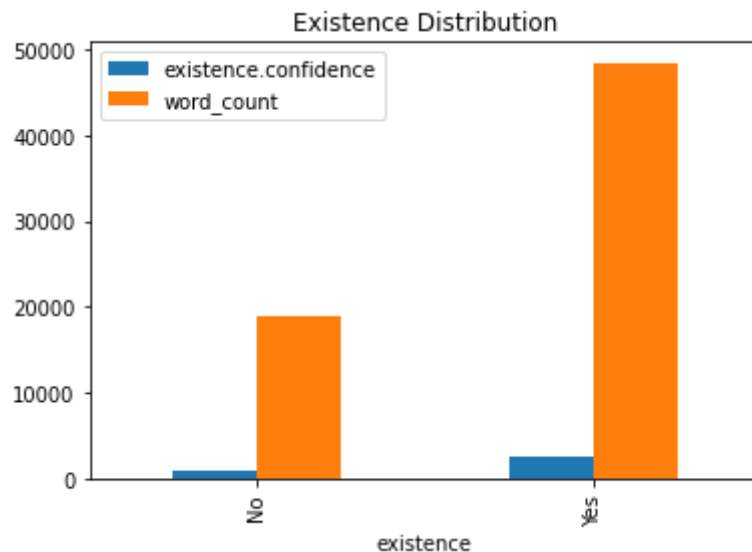
Out[27]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7ff744da9090>



```
In [ ]: 1 print(data.groupby(by=["existence"]).sum())
        2 data.groupby(by=["existence"]).sum().plot(kind = 'bar',title = "Existen
```

	existence.confidence	word_count
existence		
No	849.8694	19016
Yes	2555.6852	48464

Out[53]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7ff74331d950>





```

In [ ]: 1 # The distribution of top unigrams before removing stop words
2 from sklearn.feature_extraction.text import CountVectorizer
3
4 def get_top_n_words(corpus, n=None):
5     vec = CountVectorizer().fit(corpus)
6     bag_of_words = vec.transform(corpus)
7     sum_words = bag_of_words.sum(axis=0)
8     words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary.items()]
9     # words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary.items()]
10    words_freq = sorted(words_freq, key = lambda x: x[1], reverse=True)
11    return words_freq[:n]
12
13 common_words = get_top_n_words(data['tweet'], 20)
14 for word, freq in common_words:
15     print(word, freq)
16 df1 = pd.DataFrame(common_words, columns = ['tweet' , 'word_count'])
17 display(df1.head())
18
19 df1.groupby('tweet').sum()['word_count'].sort_values(ascending=False).plot(
20     kind='bar', title='Top 20 words in review before preprocessing')
21
22

```

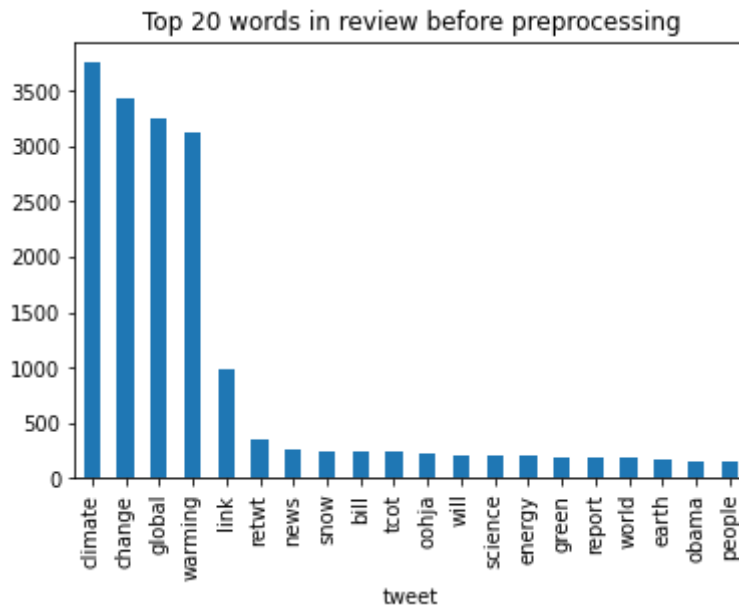
```

climate 3756
change 3438
global 3247
warming 3134
link 980
retwt 358
news 266
snow 248
bill 247
tcot 242
oohja 217
will 212
science 210
energy 204
green 192
report 179
world 178
earth 173
obama 155
people 149

```

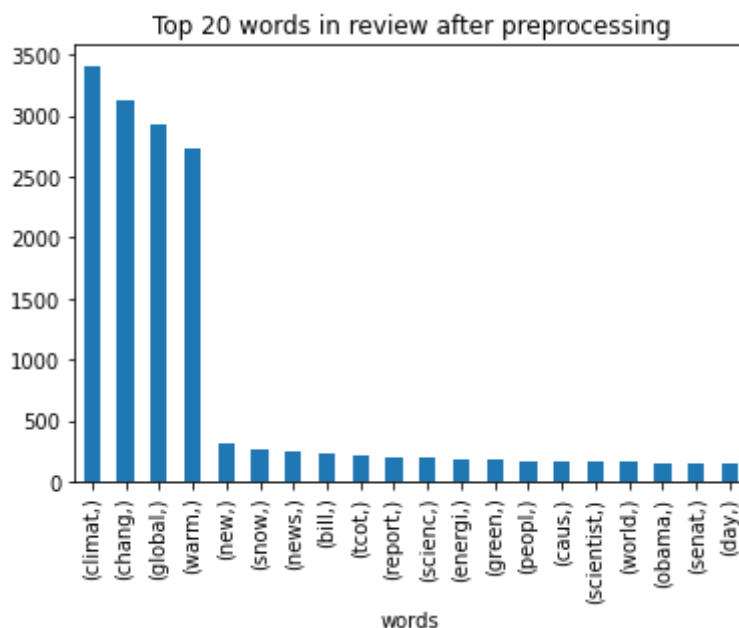
	tweet	word_count
0	climate	3756
1	change	3438
2	global	3247
3	warming	3134
4	link	980

Out[54]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7ff7432e5510>



```
In [ ]: 1 allwords = list(chain.from_iterable(preprocessed_tweets))
2 df2 = pd.DataFrame(allwords, columns = ["words"])
3 df2.value_counts()[:20].plot(kind='bar', title='Top 20 words in review
4 # allwords
5 # preprocessed_tweets
```

Out[51]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7ff73bbc4bd0>



```
In [ ]: 1
```

In [ ]:

1