

nbc_segments

March 21, 2022

```
[92]: filename = "data/Analysis_Colorado_Fire_12_30_21All_NBC.docx"
```

```
[93]: import sys
      sys.path.append('../')

      from helpers.utils import read_docx_to_dict
```

```
[95]: data = read_docx_to_dict(filename)
```

```
[96]: import pandas as pd
      pd.options.display.max_rows = 500

      # create dataframe
      df = pd.DataFrame.from_dict(data)
```

```
[97]: import sys
      sys.path.append('../')

      from helpers.utils import check_text_likeness
```

```
[98]: df['matches'] = df.apply(lambda row: check_text_likeness(df, row['text']),
      ↪axis=1)
```

```
[99]: from helpers.utils import fetch_biggest_text, mark_use_row
```

```
[119]: df['row_to_use'] = df.apply(lambda row: fetch_biggest_text(df, row['matches']),
      ↪axis=1)
```

```
[120]: mark_use_row(df)
```

```
[120]: 'done'
```

```
[136]: df['text'] = df['text'].str.lower()

      df['words'] = df['text'].str.lower().str.replace(',', ' ').str.replace('>', ' ').
      ↪str.replace('.', ' ').str.replace('\n', ' ').str.replace("'", '"').str.replace(
```

```

    '!', ')).str.replace('?', ')).str.replace('%', ')).str.replace(')', ')).str.
↪replace('(', ')).str.replace('_', ')).str.replace(':', ')).str.strip().str.
↪split(' ')

```

/Users/loren/.pyenv/versions/3.7.4/lib/python3.7/site-packages/ipykernel_launcher.py:3: FutureWarning: The default value of regex will change from True to False in a future version. In addition, single character regular expressions will **not** be treated as literal strings when regex=True.

This is separate from the ipykernel package so we can avoid doing imports until

/Users/loren/.pyenv/versions/3.7.4/lib/python3.7/site-packages/ipykernel_launcher.py:4: FutureWarning: The default value of regex will change from True to False in a future version. In addition, single character regular expressions will **not** be treated as literal strings when regex=True.

after removing the cwd from sys.path.

```

[137]: import sys
sys.path.append('../')

from helpers.utils import parse_words
df['clean_words'] = df.apply(lambda row: parse_words(row['words']), axis=1)

```

```

[138]: class Color:
    PURPLE = '\033[95m'
    CYAN = '\033[96m'
    DARKCYAN = '\033[36m'
    BLUE = '\033[94m'
    GREEN = '\033[92m'
    YELLOW = '\033[93m'
    RED = '\033[91m'
    BOLD = '\033[1m'
    UNDERLINE = '\033[4m'
    END = '\033[0m'

    def highlight_word(word, text):
        highlighted_word = Color.BOLD + Color.RED + Color.UNDERLINE + word + Color.
↪END

        if word in text:
            return text.replace(word, highlighted_word)
        else:
            return ''

    def highlight_words_found(climate_words, text):
        if not climate_words:
            return ''

```

```
return [highlight_word(word, text) for word in climate_words][0]
```

```
[139]: from helpers.utils import fetch_climate_words_in_words, \
        ↪ fetch_climate_phrases_in_text
df['climate_phrases_found'] = df.apply(lambda row: \
        ↪ fetch_climate_phrases_in_text(row['text']), axis=1)
df['climate_words_found'] = df.apply(lambda row: \
        ↪ fetch_climate_words_in_words(row['clean_words']), axis=1)
```

```
[140]: df['highlighted_words'] = df.apply(lambda row: \
        ↪ highlight_words_found(row['climate_words_found'], row['text']), axis=1)
df['highlighted_phrases'] = df.apply(lambda row: \
        ↪ highlight_words_found(row['climate_phrases_found'], row['text']), axis=1)
```

```
[141]: df.head()
```

```
[141]:      time    location station \
0  2021-12-30 6:26 PM    Denver    KUSA
1  2021-12-30 6:25 PM    Denver    KUSA
2  2021-12-30 6:24 PM    Denver    KUSA
3  2021-12-30 6:23 PM    Denver    KUSA
4  2021-12-30 6:22 PM  Nashville    WSMV

      text matches  row_to_use \
0  that are being lost further into the neighborh...    [0]          0
1  would be a different story. yeah, can we show ...    [1]          1
2  you go three houses down and the house is stil...    [2]          2
3  well, and i think tom, one thing that becomes ...    [3]          3
4  year. breaking news at 6, thousands of people ...    [4]          4

      use_row      words \
0    True  [that, are, being, lost, further, into, the, n...
1    True  [would, be, a, different, story, yeah, can, we...
2    True  [you, go, three, houses, down, and, the, house...
3    True  [well, and, i, think, tom, one, thing, that, b...
4    True  [year, breaking, news, at, 6, thousands, of, p...

      clean_words  climate_phrases_found \
0  [that, are, being, lost, further, into, the, n...    []
1  [would, be, a, different, story, yeah, can, we...    []
2  [you, go, three, houses, down, and, the, house...    []
3  [well, and, i, think, tom, one, thing, that, b...    []
4  [year, breaking, news, at, thousands, of, peop...    []

      climate_words_found      highlighted_words \
0  [high, wildfire]  that are being lost further into the neighborh...
1                   []
```

```
2          []
3          []
4          []
```

```
highlighted_phrases
0
1
2
3
4
```

```
[142]: # save data to csv
df.to_csv('reports/nbc_all.csv', encoding='utf-8')
df.to_excel('reports/nbc_all.xlsx', engine='xlsxwriter', encoding='utf-8')

# https://stackoverflow.com/questions/50495463/
# unable-to-change-font-color-in-excel-using-python-xlsxwriter
```

```
[143]: unique_df = df[df['use_row']]
```

```
[144]: total_words = unique_df['clean_words'].str.len().sum()
total_words
```

```
[144]: 16442
```

```
[145]: from helpers.utils import words_found_master_list

words_found = words_found_master_list(unique_df['clean_words'])
len(words_found)
```

```
[145]: 16442
```

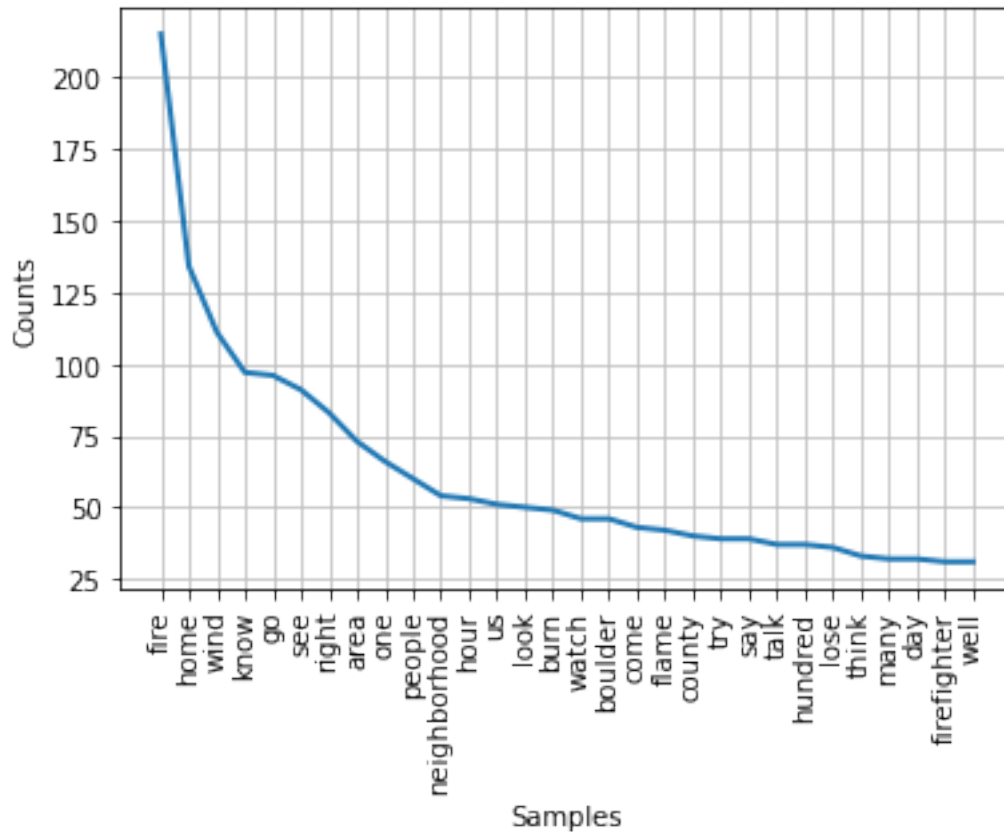
```
[146]: from helpers.utils import clean_lemmatized_words, lemmatize_words
clean_lemma_words = clean_lemmatized_words(lemmatize_words(words_found))
```

```
[147]: from nltk.probability import FreqDist

lfdist = FreqDist(clean_lemma_words)
lfdist
```

```
[147]: FreqDist({'fire': 215, 'home': 134, 'wind': 111, 'know': 97, 'go': 96, 'see':
91, 'right': 83, 'area': 73, 'one': 66, 'people': 60, ...})
```

```
[148]: import matplotlib.pyplot as plt
lfdist.plot(30, cumulative=False)
plt.show()
```



```
[149]: from wordcloud import WordCloud
from wordcloud import ImageColorGenerator
from wordcloud import STOPWORDS
import matplotlib.pyplot as plt

from helpers.utils import master_stopwords_list

wordcloud = WordCloud(width = 3000, height = 2000, random_state=1,
    ↳background_color='black', colormap='Set2', collocations=False, stopwords =
    ↳master_stopwords_list()).generate_from_frequencies(lfdist)

# Plot
plt.figure( figsize=(15,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()

#plt.savefig('word_cloud.png')
```



```

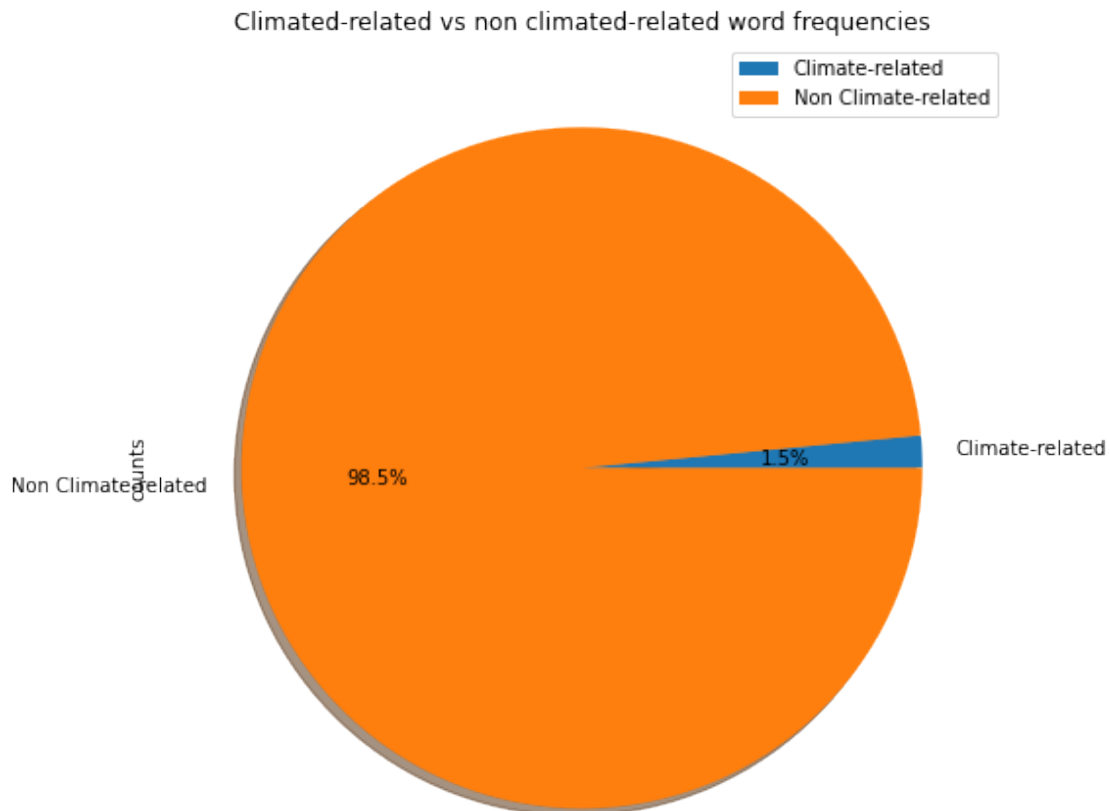
comparison_df = pd.DataFrame({'Words': ['Climate-related', 'Non-Climate-related'],
                               'counts': [climate_words_count, non_climate_words_count]})
comparison_df.set_index('Words', inplace=True)
print(comparison_df)

plot = comparison_df.plot.pie(y='counts', title="Climate-related vs non-Climate-related word frequencies",
                               legend=True, autopct='%1.1f%%', shadow=True, figsize=(8, 8))

fig = plot.get_figure()
fig.savefig("comparison.png")

```

Words	counts
Climate-related	117
Non Climate-related	7571

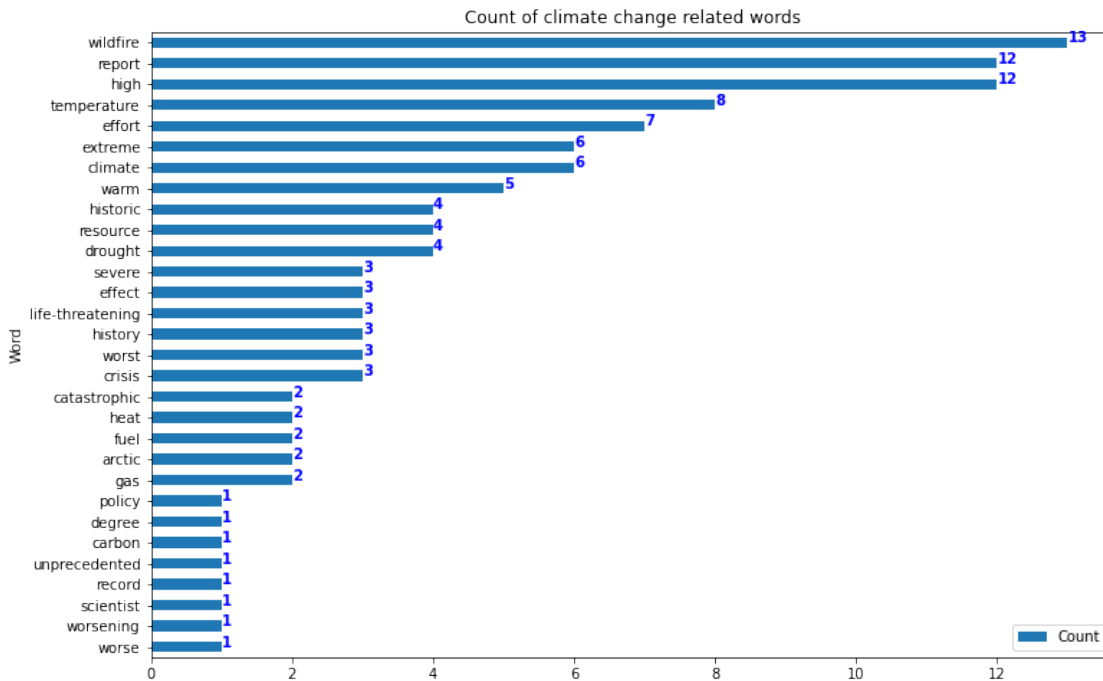


```
[152]: # find climate related word frequencies

# set figure size
fig, ax = plt.subplots(figsize=(12, 8))
# plot horizontal bar plot
climate_change_words_df.sort_values(by='Count').plot.barh(x="Word", y="Count",
    ↪ax=ax)
# set the title
plt.title("Count of climate change related words")

for i, v in enumerate(climate_change_words_df['Count'].sort_values()):
    ax.text(v, i, str(v),
            color = 'blue', fontweight = 'bold')

plt.show()
# plt.savefig('climate-related-words-breakdown.png', transparent=False)
```



```
[153]: # find segments
climate_change_words_found = list(climate_change_words_df['Word'].unique())
climate_change_words_found
```

```
[153]: ['wildfire',
        'report',
        'high',
```



```

'temperature',
'effort',
'extreme',
'climate',
'warm',
'historic',
'resource',
'drought',
'effect',
'life-threatening',
'severe',
'history',
'worst',
'crisis',
'catastrophic',
'heat',
'fuel',
'arctic',
'gas',
'degree',
'carbon',
'unprecedented',
'record',
'scientist',
'worsening',
'policy',
'worse']

```

```

[154]: unique_df[unique_df["climate_words_found"].str.len() != 0].to_csv('reports/
↳ abc_final.csv', encoding='utf-8')

```

```

[155]: # total segments
total_segments = len(df)
total_segments

```

```

[155]: 120

```

```

[156]: # unique segments
unique_segments = len(df[df['use_row'] == True])
unique_segments

```

```

[156]: 98

```

```

[157]: # how many segments had climate related words/phrases - %
possible_climate_related_segments = len(df[(df["climate_words_found"].str.len()
↳ != 0) & (df["use_row"] == True)])
possible_climate_related_segments

```

```
f'{possible_climate_related_segments / unique_segments * 100.0} %'  
# https://stackoverflow.com/questions/50495463/  
↪unable-to-change-font-color-in-excel-using-python-xlsxwriter
```

```
[157]: '43.87755102040816 %'
```