

```
import pandas as pd
import numpy as np
import nltk
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
import string
from string import punctuation
import os
from os import listdir
from collections import Counter
from tensorflow.keras.preprocessing.text import Tokenizer
from nltk.tokenize import word_tokenize
from wordcloud import WordCloud
from textblob import TextBlob
from plotnine import *
from pandas.api.types import CategoricalDtype
from sklearn import decomposition
from sklearn import preprocessing
from sklearn import metrics

def file_to_string(filename):
    '''Opens the input text file and
    returns a string of all its text.'''
    file = open(filename, 'r')
    text = file.read()
    file.close()
    text = text.replace('\n', ' ')
    text = text.replace(' ', ' ')
    return text

#cd ..

pwd

'/Users/kamilapalys/Desktop/school/data450/capstone'

# example of using the function

filepath = 'data/text/cnn_trump.txt'
test_txt = file_to_string(filepath)
```

```
test_txt
```

"Donald Trump faces more than 30 counts related to business fraud in an indictment from a Manhattan district attorney's office over the first time in American history that a current or former president has faced criminal charges, which has never seen one of its ex-leaders confronted with criminal charges, let alone while in office. Trump released a statement in response to the indictment claiming it was a 'shock' to him. 'We are united and strong - will first defeat Alvin Bragg, and then we will defeat Joe Biden, and we will be ready for whatever comes our way. "Is this a shock today? Hell yes," the person said, speaking on a condition of anonymity. The person, who has not been identified, said he or she as well as his 2024 GOP rivals - have condemned the Manhattan district attorney's office over the handling of the case."

```
def clean_text(text, stem):
    '''Takes in a string of text cleans it by converting
    to lowercase, removing punctuation, and removing stopwords.
    Also takes in a binary value to indicate if stemming should
    be performed. Returns the new string.'''
    if stem not in [0, 1]:
        raise ValueError("Stem must be a binary value (0 or 1)")
    ps = PorterStemmer()
    stemmed_dict = {}
    # create list of stopwords
    stopwords_list = stopwords.words('english')
    # make the text lowercase
    text = text.lower()
    text = text.replace('-', ' ')
    text = text.replace('u.s.', 'us')
    # convert to ascii characters
    text = text.encode("ascii", "ignore").decode()
    for chr in text:
        # only keep characters in the string that are not punctuation symbols
        if (chr in string.punctuation or chr in string.digits):
            text = text.replace(chr, ' ')
    text = text.replace(' ', ' ')
    # stem the tokens within the text
    tokens = text.split()
    new_tokens = []
    for token in tokens[:-2]:
        # only include new token in the cleaned list if not a stopword
        if token not in stopwords_list:
            if stem == 1:
                stemmed_word = ps.stem(token)
                new_tokens.append(stemmed_word)
```

```

        # to be able to map each token to the resulting stemmed word
        if token not in stemmed_dict:
            stemmed_dict[token] = stemmed_word
        else:
            new_tokens.append(token)
    new_tokens.append(tokens[-2])
    new_tokens.append(tokens[-1])
    cleaned_text = " ".join(new_tokens)
    cleaned_text = cleaned_text.replace(' ', ' ')
    return cleaned_text

# looping through all text files to apply preprocessing functions
article_docs = []
dir = os.listdir('data/text/')
dir.sort()
for filename in dir:
    filepath = os.path.join('data/text/', filename)
    if filename.split(".")[-1] == "txt":
        article_string = file_to_string(filepath)
        new_string = clean_text(article_string, 1)
        article_docs.append(new_string)

# convert the list of article strings into a binary-value dataframe
t = Tokenizer()
t.fit_on_texts(article_docs)
print(t)
encoded_docs = t.texts_to_matrix(article_docs, mode='binary')
words = [x for x in t.word_index.keys()]
binary_df = pd.DataFrame(data = encoded_docs[:, 1:], columns=words)
# List of conditions
source_conditions = [
    binary_df['abcarticle'] == 1
    , binary_df['bbcarticle'] == 1
    , binary_df['cnnarticle'] == 1
    , binary_df['foxarticle'] == 1
    , binary_df['nbcarticle'] == 1
    , binary_df['nyparticle'] == 1
    , binary_df['nytarticle'] == 1
    , binary_df['wparticle'] == 1
    , binary_df['wsjarticle'] == 1
]

```

```

# List of values to return
source_choices = [
    "ABC News"
    , "BBC"
    , "CNN"
    , "Fox News"
    , "NBC News"
    , "New York Post"
    , "The New York Times"
    , "The Washington Post"
    , "The Wall Street Journal"
]

# List of conditions
topic_conditions = [
    binary_df['affirmativearticle'] == 1
    , binary_df['balloonarticle'] == 1
    , binary_df['bidenarticle'] == 1
    , binary_df['hamasarticle'] == 1
    , binary_df['pentagonarticle'] == 1
    , binary_df['santosarticle'] == 1
    , binary_df['tanksarticle'] == 1
    , binary_df['trumparticle'] == 1
]

# List of values to return
topic_choices = [
    "Supreme Court Ruling on Affirmative Action"
    , "Chinese Surveillance Balloon"
    , "Biden's Low Approval Rates in Polls"
    , "The Deadliest Attack by Hamas"
    , "Pentagon Documents Leak"
    , "George Santos' Expulsion from Congress"
    , "U.S. and Germany Send Tanks to Ukraine"
    , "Trump's Indictment"
]

# create a new source column
binary_df["article_source"] = np.select(source_conditions, source_choices, "ERROR")

# create a new topic column
binary_df["article_topic"] = np.select(topic_conditions, topic_choices, "ERROR")

```

```
binary_df.head()
```

```
<keras.src.preprocessing.text.Tokenizer object at 0x179012e50>
```

	said	us	trump	biden	offici	mr	israel	presid	tank	hous	...	messr	overlap	vs	convert
0	1.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0
1	1.0	1.0	0.0	1.0	1.0	0.0	0.0	1.0	0.0	1.0	...	0.0	0.0	0.0	0.0
2	1.0	0.0	1.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0	...	0.0	0.0	0.0	0.0
3	1.0	1.0	0.0	1.0	1.0	0.0	1.0	1.0	0.0	1.0	...	0.0	0.0	0.0	0.0
4	1.0	1.0	0.0	0.0	1.0	0.0	1.0	1.0	1.0	0.0	...	0.0	0.0	0.0	0.0

```
# check what word a given stemmed word represents  
value = {i for i in stemmed_dict if stemmed_dict[i]=="u"}  
print("key by value:",value)
```

```
key by value: {'u'}
```

```
article_docs
```

```
['suprem court thursday set new limit affirm action program case involv whether public priva  
'massiv spi balloon believ china seen montana track fli across continent unit state presid  
'two day dismal new poll number presid joe biden show public view unfavor rival donald trump  
'hundr peopl israel report dead thousand injur hama milit fire rocket gaza launch ground inc  
'post social media appear sever highli classifi u intellig document might begin could turn s  
'hous repres friday vote expel republican rep georg santo histor move happen year santo scan  
'major increas u support ukrain presid joe biden sign send abram tank war torn countri conce  
'manhattan grand juri indict former presid donald trump make first current former presid fac  
'us suprem court rule race longer consid factor univers admiss landmark rule upend decad old  
'us track suspect chines surveil balloon spot fli sensit site recent day defenc offici said  
'us presid joe biden run elect next year nation opinion poll weak job approv rate suggest vo  
'least peopl report kill wound israel palestinian milit group hama launch biggest attack yea  
'document includ detail account train provid ukrain foreign power make dozen classifi us de  
'us hous repres expel congressman georg santo follow damn ethic report dozen crimin charg he  
'us send power battl tank ukrain join germani send vehicl support fight russia invas decis  
'former us presid donald trump charg hush money payment made porn star presidenti elect deta  
'suprem court say colleg univers longer take race consider specif basi grant admiss landmark  
'us track suspect chines high altitud surveil balloon continent unit state defens offici sa
```

'one third regist voter approv presid joe biden handl isra palestinian conflict new poll ne  
'gaza jerusalem cnn israel prime minist benjamin netanyahu declar countri war saturday pale  
'man arrest fbi connect massiv us classifi document leak charg boston friday unauthor retent  
'hous vote friday expel gop rep georg santo histor vote make new york congressman sixth law  
'leader unit state germani announcc wednesday send conting tank ukraine revers longstand trep  
'donald trump face count relat busi fraud indict manhattan grand juri accord two sourc famili  
'u suprem court hand major rule affirm action thursday reject use race factor colleg admiss  
'u govern monitor suspect chines surveil balloon move northern state past sever day pentagon  
'biden battl rough poll number fox news white hous correspond peter dooci latest presid elec  
'hama widespread coordin attack may suggest outsid help trey yingst foreign correspond trey  
'bret baier intel suspect year old govern worker special report bret baier anchor question :  
'hous repres vote expel scandal plagu rep georg santo r n friday make first hous lawmak exp  
'german chancellor olaf scholz formal announcc wednesday week stall frustrat negoti berlin ag  
'former presid donald trump indict part manhattan district attorney offic year long investig  
'washington suprem court thursday struck affirm action program univers north carolina harvard  
'u militari monitor suspect chines surveil balloon hover northern u past day militari defens  
'differ poll agre presid joe biden polit stand lower barack obama point elect even lower don  
'ashkelon israel israel plung chao saturday palestinian milit group hama launch deadli land  
'dozen leak defens depart classifi document post onlin reveal detail u spi russia war machin  
'washington hous vote overwhelmingli expel indict rep georg santo friday pull curtain temper  
'ukrain set receiv battl tank germani western countri fierc debat expos fissur among alli al  
'grand juri new york citi vote thursday indict donald trump first time former u presid face  
'suprem court struck affirm action program harvard univers univers north carolina thursday  
'us track chines spi balloon float northern part countri day pentagon offici announcc thursda  
'presid biden readi ring new year see number year old command chief end lower approv rate se  
'jerusalem ap hama milit fire thousand rocket sent dozen fighter isra town near gaza strip u  
'massachusett air nation guardsman jack teixeira leader discord group dozen sensit us intellig  
'bye georg lie long island rep georg santo r ny becam sixth member ever expel us hous repres  
'week excus foot drag german final said ye transfer limit number leopard tank ukraine agre w  
'donald trump indict manhattan grand juri thursday hush money payment made ahead elect mark  
'chief justic john g robert jr finish read major opinion suprem court chamber thursday hush  
'helena mont larri mayer newspap photograph point camera sky wednesday began snap pictur app  
'democrat battleground state grow increasingli anxiou presid biden low approv rate worri vot  
'israel news site compil list dead miss funer take place around countri weekend attack confu  
'washington would year old nation guardsman posit access top secret document begin dramat a  
'georg santo new york republican congressman whose tapestri lie scheme made figur nation rid  
'precis militari drill first germani unit state announcc wednesday agre provid battl tank hel  
'manhattan grand juri indict donald j trump thursday role pay hush money porn star accord po  
'suprem court thursday held race consciou admiss program harvard univers north carolina viol  
'chines surveil balloon collect intellig continent unit state right u offici disclos thursda  
'night presid biden depart washington celebr thanksgiv nantucket mass gather closest aid me  
'sderot israel israel formal declar war palestinian milit group hama sunday reel surpris att  
'saturday u offici foreign alli scrambl understand dozen classifi intellig document end int

'hous vote friday expel rep georg santo r n congress action chamber previous taken five time  
'biden administr announc wednesday send premier battl tank ukrain follow agreement germani c  
'new york manhattan grand juri vote indict former presid donald trump make first person u h  
'thursday decis forc rework admiss criteria throughout american higher educ decad pursuit d  
'washington u track offici describ chines reconnaiss balloon continent state week would agg  
'mr biden low stand head head gener elect poll reflect voter dour apprais perform presid su  
'tel aviv isra prime minist benjamin netanyahu said countri war hama milit group forc pour a  
'crimin case unfold u govern scrambl protect secret unauthor disclosur appear provid detail  
'lawmak vote remov two third hous supermajor requir constitut almost democrat mani republica  
'u germani outlin plan wednesday send dozen modern battl tank ukrain mark signific new infus  
'grand juri return indict mr trump vote thursday kick process former presid expect come new

```
encoded_docs_freq = t.texts_to_matrix(article_docs, mode='count')
freq_df = pd.DataFrame(data = encoded_docs_freq[:, 1:], columns=words)
# List of conditions
source_conditions = [
    freq_df['abcarticle'] == 1
    , freq_df['bbcarticle'] == 1
    , freq_df['cnnarticle'] == 1
    , freq_df['foxarticle'] == 1
    , freq_df['nbcarticle'] == 1
    , freq_df['nyparticle'] == 1
    , freq_df['nytarticle'] == 1
    , freq_df['wparticle'] == 1
    , freq_df['wsjarticle'] == 1
]
# List of values to return
source_choices = [
    "ABC News"
    , "BBC"
    , "CNN"
    , "Fox News"
    , "NBC News"
    , "New York Post"
    , "The New York Times"
    , "The Washington Post"
    , "The Wall Street Journal"
]
# List of conditions
```

```

topic_conditions = [
    freq_df['affirmativearticle'] == 1
    , freq_df['balloonarticle'] == 1
    , freq_df['bidenarticle'] == 1
    , freq_df['hamasarticle'] == 1
    , freq_df['pentagonarticle'] == 1
    , freq_df['santosarticle'] == 1
    , freq_df['tanksarticle'] == 1
    , freq_df['trumparticle'] == 1
]
# List of values to return
topic_choices = [
    "Supreme Court Ruling on Affirmative Action"
    , "Chinese Surveillance Balloon"
    , "Biden's Low Approval Rates in Polls"
    , "The Deadliest Attack by Hamas"
    , "Pentagon Documents Leak"
    , "George Santos' Expulsion from Congress"
    , "U.S. and Germany Send Tanks to Ukraine"
    , "Trump's Indictment"
]
# create a new source column
freq_df["article_source"] = np.select(source_conditions, source_choices, "ERROR")

# create a new topic column
freq_df["article_topic"] = np.select(topic_conditions, topic_choices, "ERROR")

freq_df.head()

```

	said	us	trump	biden	offici	mr	israel	presid	tank	hous	...	messr	overlap	vs	conver
0	5.0	4.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0
1	30.0	15.0	0.0	5.0	28.0	0.0	0.0	4.0	0.0	3.0	...	0.0	0.0	0.0	0.0
2	9.0	0.0	17.0	27.0	0.0	0.0	0.0	8.0	0.0	4.0	...	0.0	0.0	0.0	0.0
3	17.0	10.0	0.0	3.0	6.0	0.0	28.0	2.0	0.0	1.0	...	0.0	0.0	0.0	0.0
4	9.0	20.0	0.0	0.0	5.0	0.0	2.0	2.0	3.0	0.0	...	0.0	0.0	0.0	0.0

```

# create dataframe with tf-idf values

encoded_docs_tfidf = t.texts_to_matrix(article_docs, mode='tfidf')
tfidf_df = pd.DataFrame(data = encoded_docs_tfidf[:, 1:], columns=words)

```

```

# List of conditions
source_conditions = [
    tfidf_df['abcarticle'] != 0
    , tfidf_df['bbcarticle'] != 0
    , tfidf_df['cnnarticle'] != 0
    , tfidf_df['foxarticle'] != 0
    , tfidf_df['nbcarticle'] != 0
    , tfidf_df['nyparticle'] != 0
    , tfidf_df['nytarticle'] != 0
    , tfidf_df['wparticle'] != 0
    , tfidf_df['wsjarticle'] != 0
]

# List of values to return
source_choices = [
    "ABC News"
    , "BBC"
    , "CNN"
    , "Fox News"
    , "NBC News"
    , "New York Post"
    , "The New York Times"
    , "The Washington Post"
    , "The Wall Street Journal"
]

# List of conditions
topic_conditions = [
    tfidf_df['affirmativearticle'] != 0
    , tfidf_df['balloonarticle'] != 0
    , tfidf_df['bidenarticle'] != 0
    , tfidf_df['hamasarticle'] != 0
    , tfidf_df['pentagonarticle'] != 0
    , tfidf_df['santosarticle'] != 0
    , tfidf_df['tanksarticle'] != 0
    , tfidf_df['trumparticle'] != 0
]

# List of values to return
topic_choices = [
    "Supreme Court Ruling on Affirmative Action"
    , "Chinese Surveillance Balloon"
]

```

```

        , "Biden's Low Approval Rates in Polls"
        , "The Deadliest Attack by Hamas"
        , "Pentagon Documents Leak"
        , "George Santos' Expulsion from Congress"
        , "U.S. and Germany Send Tanks to Ukraine"
        , "Trump's Indictment"
    ]
# create a new source column
tfidf_df["article_source"] = np.select(source_conditions, source_choices, "ERROR")

# create a new topic column
tfidf_df["article_topic"] = np.select(topic_conditions, topic_choices, "ERROR")

tfidf_df.head()

```

	said	us	trump	biden	offici	mr	israel	presid	tank	hous	...
0	1.827036	1.839130	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	...
1	3.081563	2.857814	0.000000	2.267700	4.139471	0.0	0.000000	1.881491	0.000000	1.871958	...
2	2.238584	0.000000	5.086179	3.733245	0.000000	0.0	0.000000	2.428008	0.000000	2.128570	...
3	2.683881	2.545320	0.000000	1.823774	2.667558	0.0	6.446654	1.334974	0.000000	0.891998	...
4	2.238584	3.079532	0.000000	0.000000	2.493348	0.0	2.519533	1.334974	3.689062	0.000000	...

```

# create a list of strings where each string is all articles from one source (for dendogram)
source_docs = []

j = 0

for i in range(9):
    # looping through each article of each source, and only taking up until and not includ
    # token, bc last two tokens represent the name of the source and topic of article
    # for last article, we index up until and not including only the last token bc we don'
    # the topic of the article, but we do want to include the source name, which is the se
    source = " ".join(article_docs[j].split()[:-2]) + " " + ".join(article_docs[j+1].spli
        + " ".join(article_docs[j+2].split()[:-2]) + " " + " ".join(article_docs[j+3].spli
        + " ".join(article_docs[j+4].split()[:-2]) + " " + " ".join(article_docs[j+5].spli
        + " ".join(article_docs[j+6].split()[:-2]) + " " + " ".join(article_docs[j+7].spli
    source_docs.append(source)
    j += 8

source_docs

```

```
['suprem court thursday set new limit affirm action program case involv whether public priva  
'us suprem court rule race longer consid factor univers admiss landmark rule upend decad old  
'suprem court say colleg univers longer take race consider specif basi grant admiss landmark  
'us suprem court hand major rule affirm action thursday reject use race factor colleg admiss  
'washington suprem court thursday struck affirm action program univers north carolina harvard  
'suprem court struck affirm action program harvard univers univers north carolina thursday r  
'chief justic john g robert jr finish read major opinion suprem court chamber thursday hush  
'suprem court thursday held race consciou admiss program harvard univers north carolina vio  
'thursday decis forc rework admiss criteria throughout american higher educ decad pursuit d
```

```
for source in source_docs:  
    for i in range(len(source)-2):  
        if source[i] + source[i+1] + source[i+2] == ' u ':  
            print(source[i-10:i+100], source)
```

en countri u k also commit challeng battl tank poland ask germani permiss send leopard stock  
gia meloni u k prime minist rishi sunak grata chancellor scholz provid german leopard tank l  
ear accord u n mideast envoy tor wennesland israel say raid aim milit stone throw protest pe  
gust brief u n secur council surpris hama offens come year day israel arab neighbor launch c  
h mile ton u k l e mm rifl gun two mm machin gun us offici said abram would enough ukrainian

```
# create a dataframe of token tf-idf's with each row representing all articles of one sourc  
  
# convert the list of article strings into a tf-idf-value dataframe  
t = Tokenizer()  
t.fit_on_texts(source_docs)  
print(t)  
encoded_source_docs = t.texts_to_matrix(source_docs, mode='tfidf')  
words = [x for x in t.word_index.keys()]  
tfidf_source_df = pd.DataFrame(data = encoded_source_docs[:, 1:], columns=words)  
# List of conditions  
source_conditions = [  
    tfidf_source_df['abcarticle'] != 0  
    , tfidf_source_df['bbcarticle'] != 0  
    , tfidf_source_df['cnnarticle'] != 0  
    , tfidf_source_df['foxarticle'] != 0  
    , tfidf_source_df['nbcarticle'] != 0  
    , tfidf_source_df['nyparticle'] != 0  
    , tfidf_source_df['nytarticle'] != 0  
    , tfidf_source_df['wparticle'] != 0
```

```

        , tfidf_source_df['wsjarticle'] != 0
    ]

# List of values to return
source_choices = [
    "ABC News"
    , "BBC"
    , "CNN"
    , "Fox News"
    , "NBC News"
    , "New York Post"
    , "The New York Times"
    , "The Washington Post"
    , "The Wall Street Journal"
]

# create a new source column
tfidf_source_df["article_source"] = np.select(source_conditions, source_choices, "ERROR")
tfidf_source_df.set_index('article_source', inplace=True) # removes article_source column
tfidf_source_df.drop(['abcarticle', 'bbcarticle', 'cnnarticle', 'foxarticle',
                     'nbcarticle', 'nyparticle', 'nytarticle', 'wparticle',
                     'wsjarticle'], axis=1, inplace=True)
tfidf_source_df

```

<keras.src.preprocessing.text.Tokenizer object at 0x178aaaa60>

article_source	said	us	trump	biden	offici	mr	israel	presid
ABC News	3.591249	3.331000	2.959536	3.099282	3.139834	1.173600	2.824926	2.8861
BBC	3.280434	3.165511	2.780642	2.335743	2.460363	3.283902	2.460363	2.6543
CNN	3.698474	3.446401	3.025425	2.976653	3.331000	0.693147	3.225542	2.9933
Fox News	3.055995	2.866350	3.350161	2.824926	2.421451	0.693147	2.866350	2.9052
NBC News	3.670441	3.213976	3.311249	2.866350	2.654383	1.654053	3.099282	2.5317
New York Post	3.340652	3.202199	3.126599	2.757300	2.335743	1.654053	3.085175	3.0095
The New York Times	3.584733	2.052151	3.269825	2.993325	2.654383	4.166255	1.347002	2.8249
The Washington Post	3.785551	3.202199	3.269825	3.269825	3.225542	0.000000	2.976653	3.0992
The Wall Street Journal	3.849334	3.421552	3.213976	2.866350	3.126599	3.919027	2.905262	2.8031

```

# create a dataframe of token binary values with each row representing all articles of one

# convert the list of article strings into a binary-value dataframe
t = Tokenizer()
t.fit_on_texts(source_docs)
print(t)
encoded_source_docs = t.texts_to_matrix(source_docs, mode='binary')
words = [x for x in t.word_index.keys()]
binary_source_df = pd.DataFrame(data = encoded_source_docs[:, 1:], columns=words)
# List of conditions
source_conditions = [
    binary_source_df['abcarticle'] == 1
    , binary_source_df['bbcarticle'] == 1
    , binary_source_df['cnnarticle'] == 1
    , binary_source_df['foxarticle'] == 1
    , binary_source_df['nbcarticle'] == 1
    , binary_source_df['nyparticle'] == 1
    , binary_source_df['nytarticle'] == 1
    , binary_source_df['wparticle'] == 1
    , binary_source_df['wsjarticle'] == 1
]
# List of values to return
source_choices = [
    "ABC News"
    , "BBC"
    , "CNN"
    , "Fox News"
    , "NBC News"
    , "New York Post"
    , "The New York Times"
    , "The Washington Post"
    , "The Wall Street Journal"
]
# create a new source column
binary_source_df["article_source"] = np.select(source_conditions, source_choices, "ERROR")
binary_source_df.set_index('article_source', inplace=True)
binary_source_df.drop(['abcarticle', 'bbcarticle', 'cnnarticle', 'foxarticle',
                      'nbcarticle', 'nyparticle', 'nytarticle', 'wparticle',
                      'wsjarticle'], axis=1, inplace=True)

```

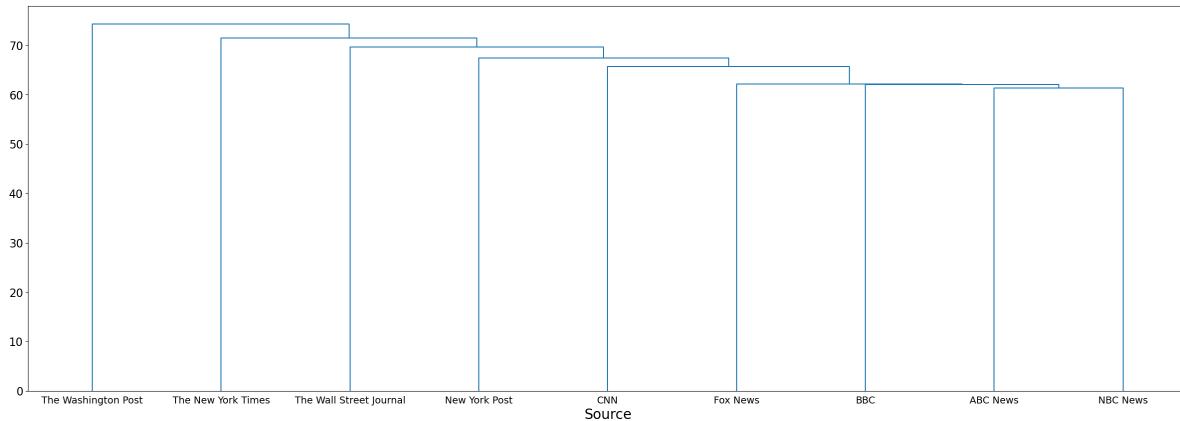
```
binary_source_df
```

```
<keras.src.preprocessing.text.Tokenizer object at 0x178d7c490>
```

article_source	said	us	trump	biden	offici	mr	israel	presid	tank	hou	...	squelch
ABC News	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	0.0
BBC	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	0.0
CNN	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	0.0
Fox News	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	0.0
NBC News	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	0.0
New York Post	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	0.0
The New York Times	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	0.0
The Washington Post	1.0	1.0	1.0	1.0	1.0	0.0	1.0	1.0	1.0	1.0	...	0.0
The Wall Street Journal	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	...	1.0

```
# create dendrogram with average linkage from tf_idf scores df, whose rows each represent  
  
from scipy.cluster.hierarchy import dendrogram, linkage  
import matplotlib.pyplot as plt  
  
Z = linkage(tfidf_source_df, 'average')  
fig = plt.figure(figsize=(30, 10))  
fig.suptitle("Average Linkage", fontsize=24)  
plt.xlabel('Source', fontsize=20)  
plt.yticks(fontsize = 16)  
dn = dendrogram(Z, labels=tfidf_source_df.index)  
plt.xticks(fontsize = 14)  
plt.show()
```

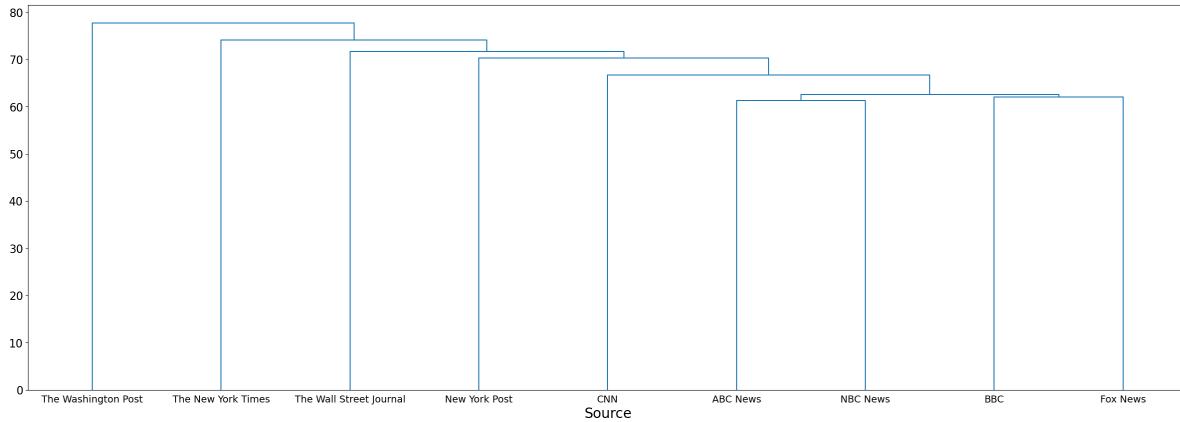
Average Linkage



```
# create dendrogram with complete linkage from tf_idf scores df, whose rows each represent
```

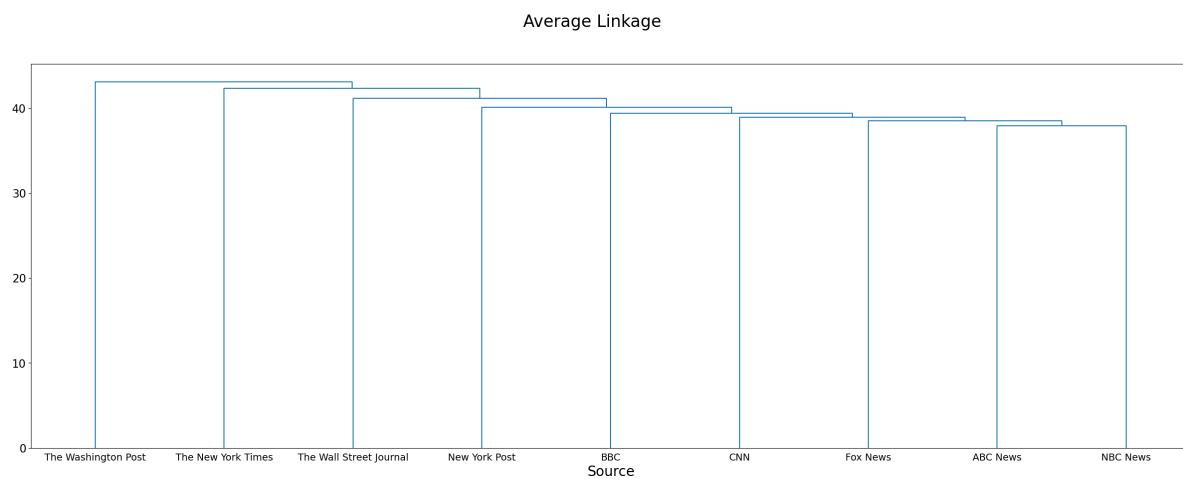
```
Z = linkage(tfidf_source_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=tfidf_source_df.index)
plt.xticks(fontsize = 14)
plt.show()
```

Complete Linkage



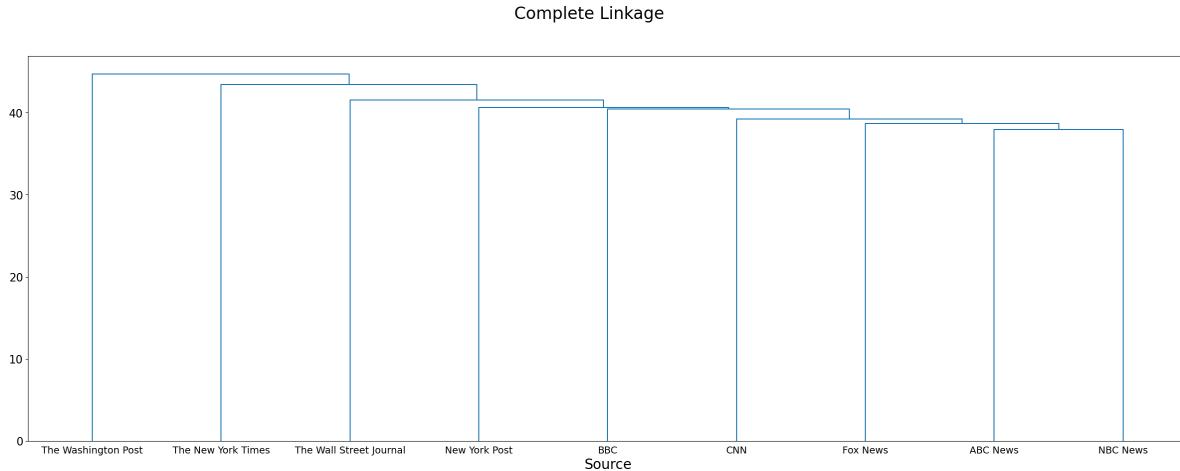
```
# create dendrogram with average linkage from binary scores df, whose rows each represent

Z = linkage(binary_source_df, 'average')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Average Linkage", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=binary_source_df.index)
plt.xticks(fontsize = 14)
plt.show()
```



```
# create dendrogram with complete linkage from binary scores df, whose rows each represent

Z = linkage(binary_source_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=binary_source_df.index)
plt.xticks(fontsize = 14)
plt.show()
```



```
# create dataframes, each containing articles of one topic for dendograms for each topic

import warnings
warnings.filterwarnings('ignore')

tfidf_df_dropped = tfidf_df.drop(['abcarticle', 'bbcarticle', 'cnnarticle', 'foxarticle',
                                   'nbcarticle', 'nyparticle', 'nytarticle', 'wparticle',
                                   'wsjarticle', 'affirmativearticle', 'balloonarticle',
                                   'bidenarticle', 'hamasarticle', 'pentagonarticle',
                                   'santosarticle', 'tanksarticle', 'trumparticle'],
                                   axis=1, inplace=False)

affirm_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "Supreme Court Rul"
affirm_tfidf_df.set_index('article_source', inplace=True)
affirm_tfidf_df.drop(['article_topic'], axis=1, inplace=True)

balloon_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "Chinese Surveill"
balloon_tfidf_df.set_index('article_source', inplace=True)
balloon_tfidf_df.drop(['article_topic'], axis=1, inplace=True)

biden_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "Biden's Low Approv"
biden_tfidf_df.set_index('article_source', inplace=True)
biden_tfidf_df.drop(['article_topic'], axis=1, inplace=True)

hamas_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "The Deadliest Atta
```

```

pentagon_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "Pentagon Document"]
pentagon_tfidf_df.set_index('article_source', inplace=True)
pentagon_tfidf_df.drop(['article_topic'], axis=1, inplace=True)

 santos_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "George Santos' Executive Order"]
 santos_tfidf_df.set_index('article_source', inplace=True)
 santos_tfidf_df.drop(['article_topic'], axis=1, inplace=True)

 tanks_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "U.S. and Germany Sanctions"]
 tanks_tfidf_df.set_index('article_source', inplace=True)
 tanks_tfidf_df.drop(['article_topic'], axis=1, inplace=True)

 trump_tfidf_df = tfidf_df_dropped[tfidf_df_dropped['article_topic'] == "Trump's Indictment"]
 trump_tfidf_df.set_index('article_source', inplace=True)
 trump_tfidf_df.drop(['article_topic'], axis=1, inplace=True)

affirm_tfidf_df.head()

```

article_source	said	trump	biden	offici	mr	u	israel	presid	tank	hous
ABC News	1.827036	0.000000	0.000000	0.0	0.000000	2.313275	0.0	0.000000	0.0	0.000000
BBC	2.379087	1.326871	0.869038	0.0	2.295628	0.000000	0.0	2.057431	0.0	0.891000
CNN	2.641434	2.784588	0.000000	0.0	0.000000	0.000000	0.0	1.334974	0.0	0.891000
Fox News	1.827036	0.000000	0.869038	0.0	0.000000	1.641338	0.0	0.788457	0.0	0.000000
NBC News	2.156116	2.246588	2.073780	0.0	0.000000	1.641338	0.0	1.881491	0.0	0.000000

```

# create dendrograms with complete linkage from tfidf scores df's, each one representing a topic

Z = linkage(affirm_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage Clustering of Affirmative Action Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=affirm_tfidf_df.index)
plt.xticks(fontsize = 14)
plt.show()

Z = linkage(balloon_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))

```

```

fig.suptitle("Complete Linkage Clustering of Chinese Surveillance Balloon Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=balloon_tfidf_df.index)
plt.xticks(fontsize = 14)
plt.show()

Z = linkage(biden_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage Clustering of Biden's Low Ratings Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=biden_tfidf_df.index)
plt.xticks(fontsize = 14)
plt.show()

Z = linkage(hamas_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage Clustering of Hamas Attack Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=hamas_tfidf_df.index)
plt.xticks(fontsize = 14)
plt.show()

Z = linkage(pentagon_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage Clustering of Pentagon Document Leak Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=pentagon_tfidf_df.index)
plt.xticks(fontsize = 14)
plt.show()

Z = linkage(santos_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage Clustering of George Santos Expulsion Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=santos_tfidf_df.index)
plt.xticks(fontsize = 14)

```

```

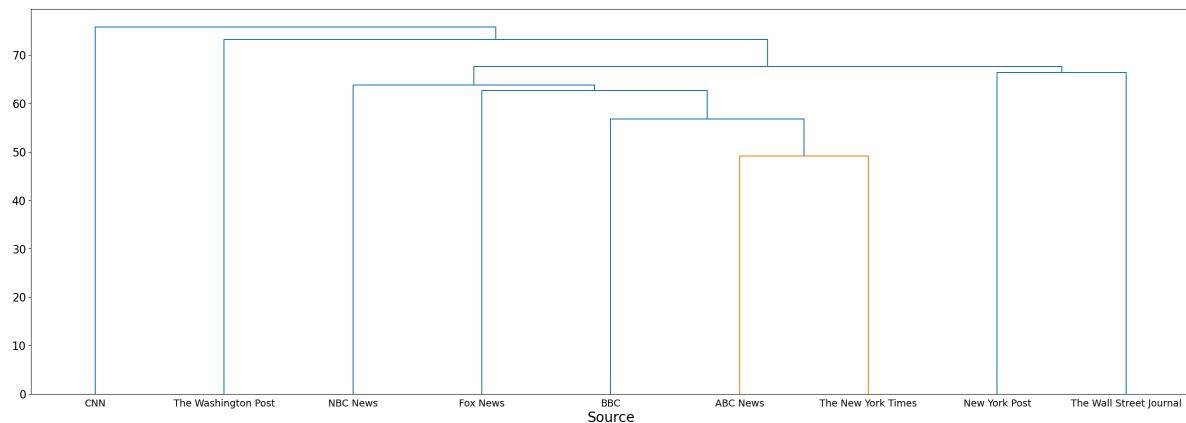
plt.show()

Z = linkage(tanks_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage Clustering of Sending Ukraine Tanks Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=tanks_tfidf_df.index)
plt.xticks(fontsize = 14)
plt.show()

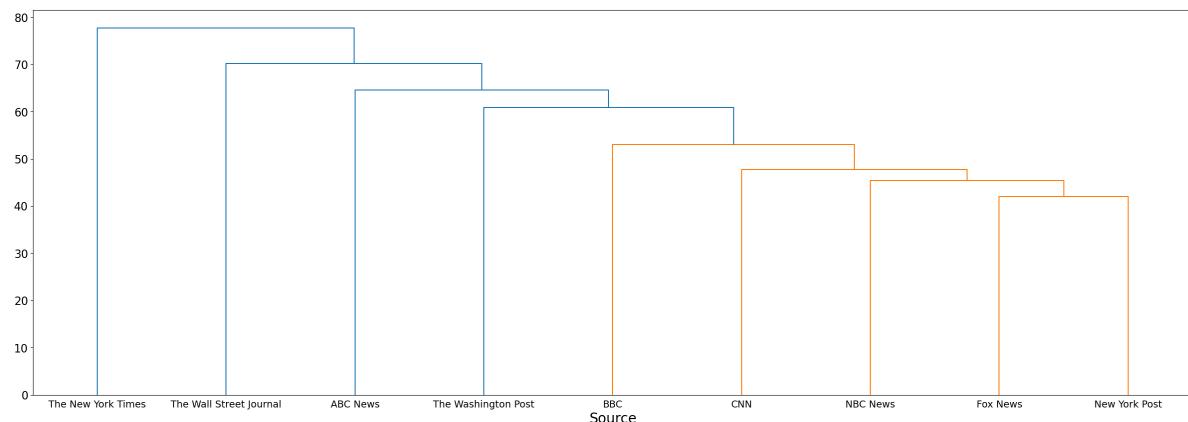
Z = linkage(trump_tfidf_df, 'complete')
fig = plt.figure(figsize=(30, 10))
fig.suptitle("Complete Linkage Clustering of Trump's Indictment Articles", fontsize=24)
plt.xlabel('Source', fontsize=20)
plt.yticks(fontsize = 16)
dn = dendrogram(Z, labels=trump_tfidf_df.index)
plt.xticks(fontsize = 14)
plt.show()

```

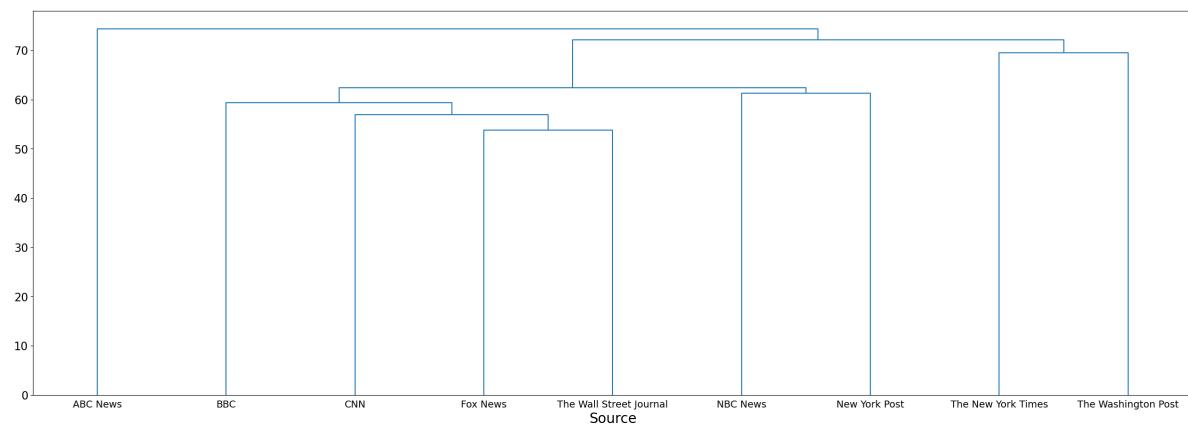
Complete Linkage Clustering of Affirmative Action Articles



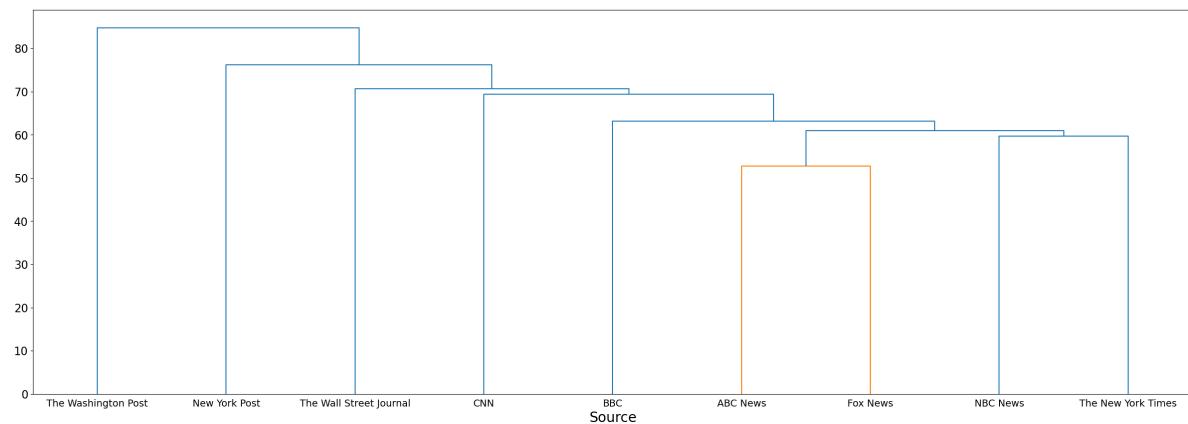
Complete Linkage Clustering of Chinese Surveillance Balloon Articles



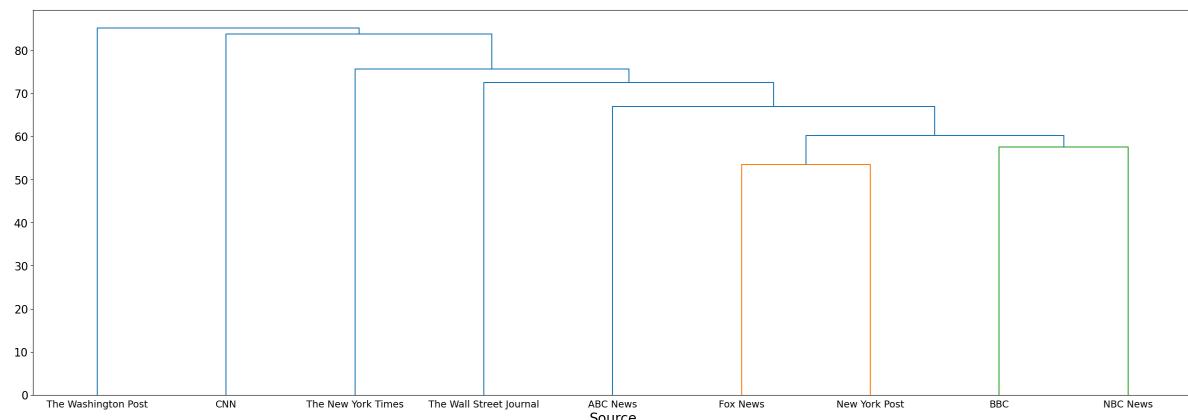
Complete Linkage Clustering of Biden's Low Ratings Articles



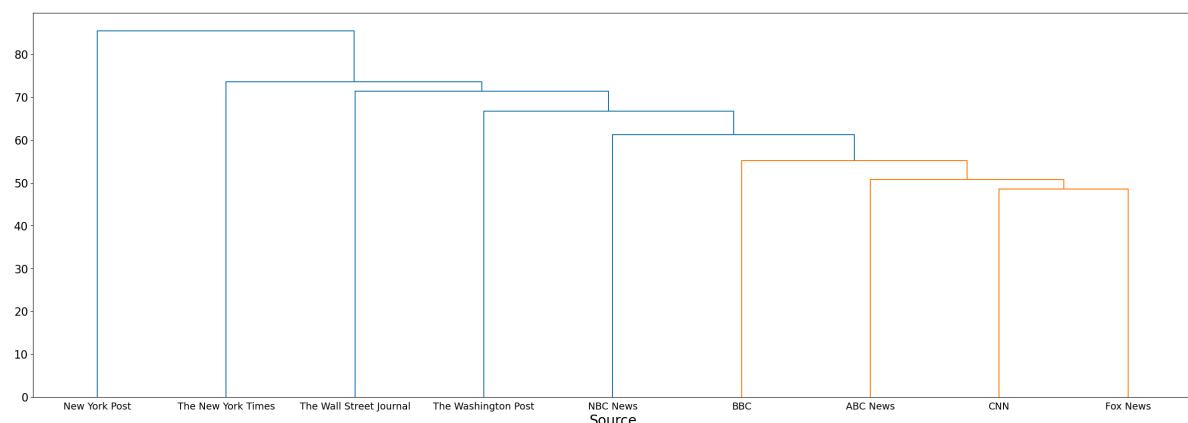
Complete Linkage Clustering of Hamas Attack Articles



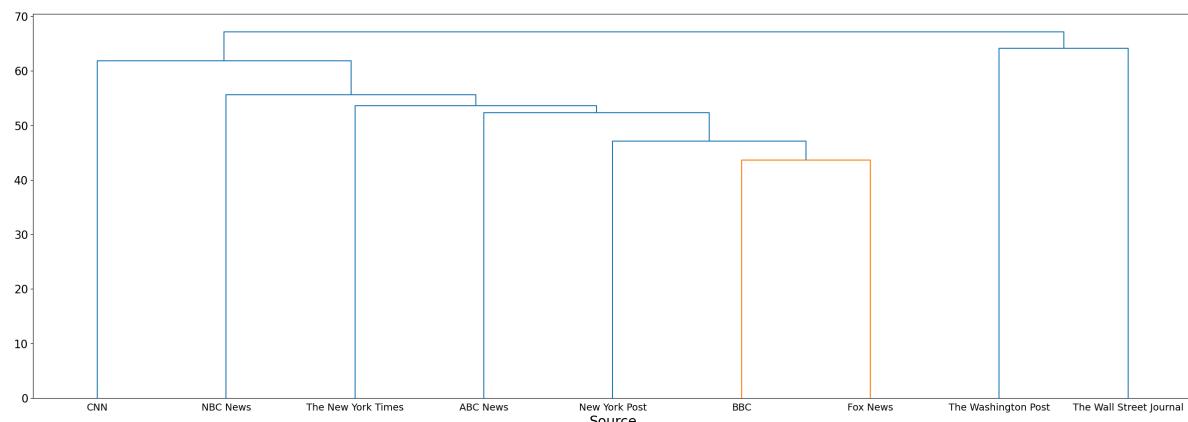
Complete Linkage Clustering of Pentagon Document Leak Articles



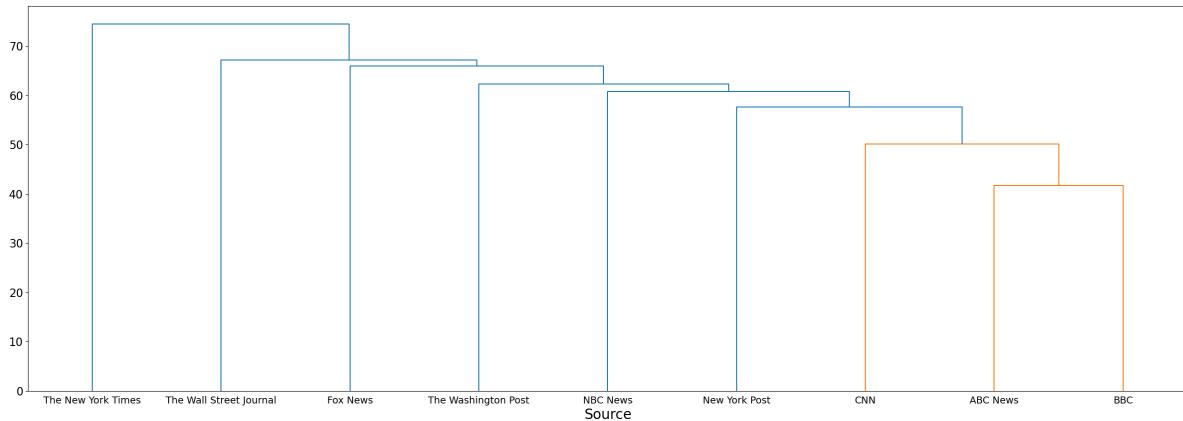
Complete Linkage Clustering of George Santos Expulsion Articles



Complete Linkage Clustering of Sending Ukraine Tanks Articles



Complete Linkage Clustering of Trump's Indictment Articles

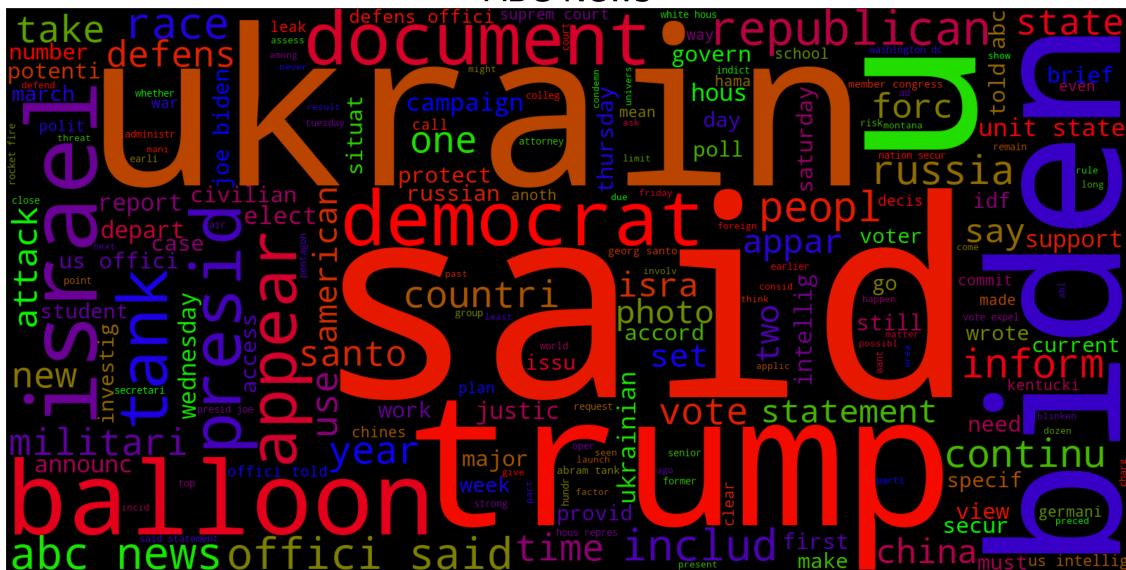


```
# other color themes i liked: magma, PiYG

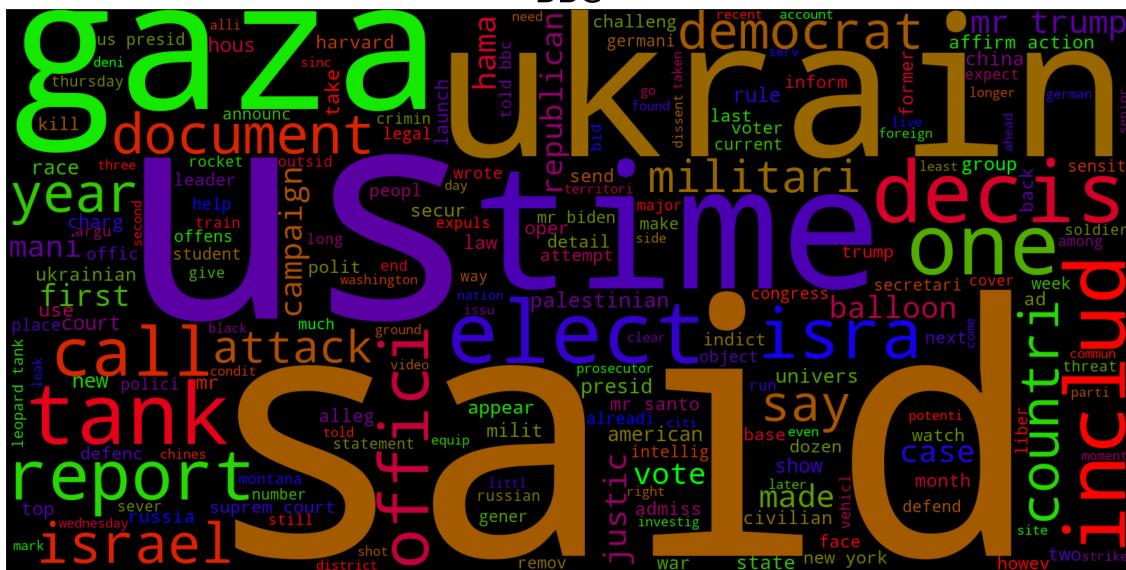
def create_wordcloud(text, title):
    '''Given a string of all text and a string
    for the title, creates a wordcloud.'''
    plt.figure(figsize = (20,10))
    wc = WordCloud(colormap = "brg", width=1600, height=800).generate(text)
    plt.title(title, fontsize=40)
    plt.axis("off")
    title_snake = title.lower().replace(" ", "_")
    wc.to_file(f'wordcloud_{title_snake}.png')
    #plt.savefig(f'wordcloud_{title_snake}.png', dpi = 1000)
    plt.tight_layout(pad=0)
    plt.imshow(wc)

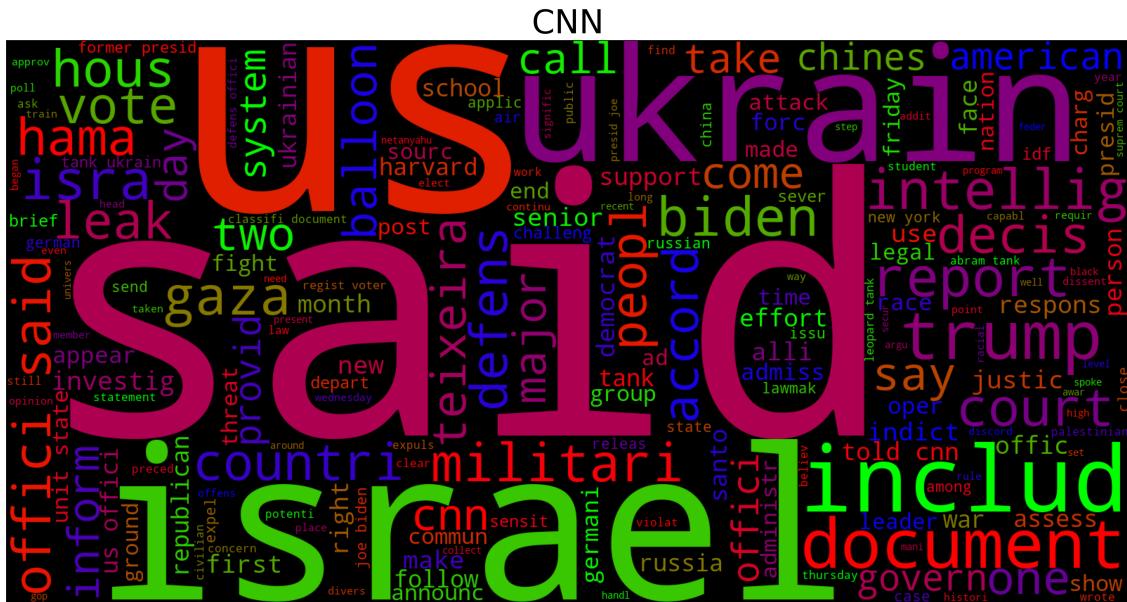
for i in range(len(source_docs)):
    create_wordcloud(source_docs[i], tfidf_source_df.index[i])
```

ABC News

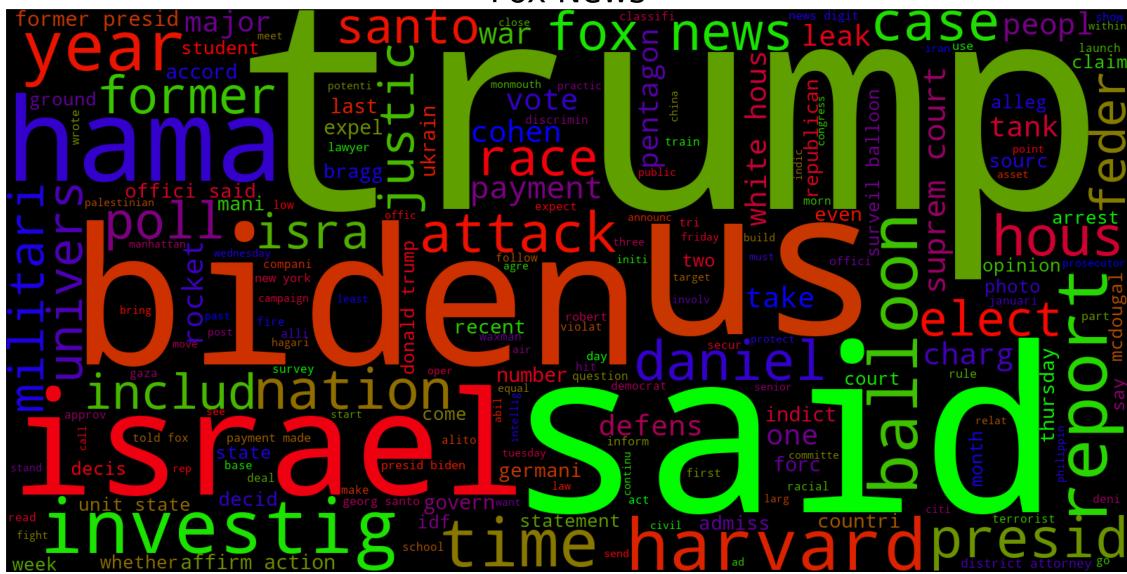


BBC

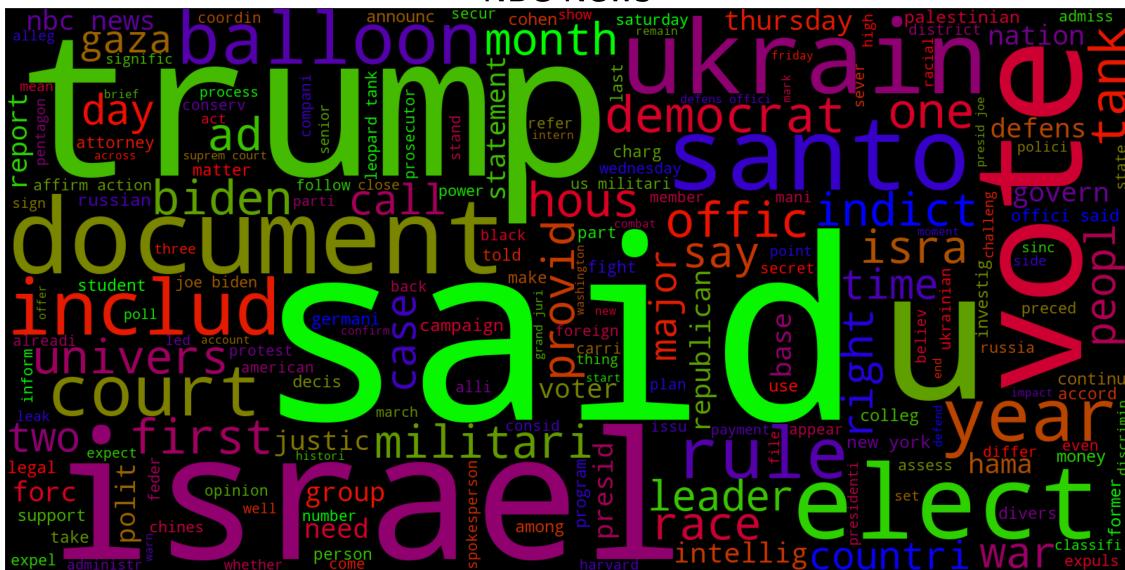




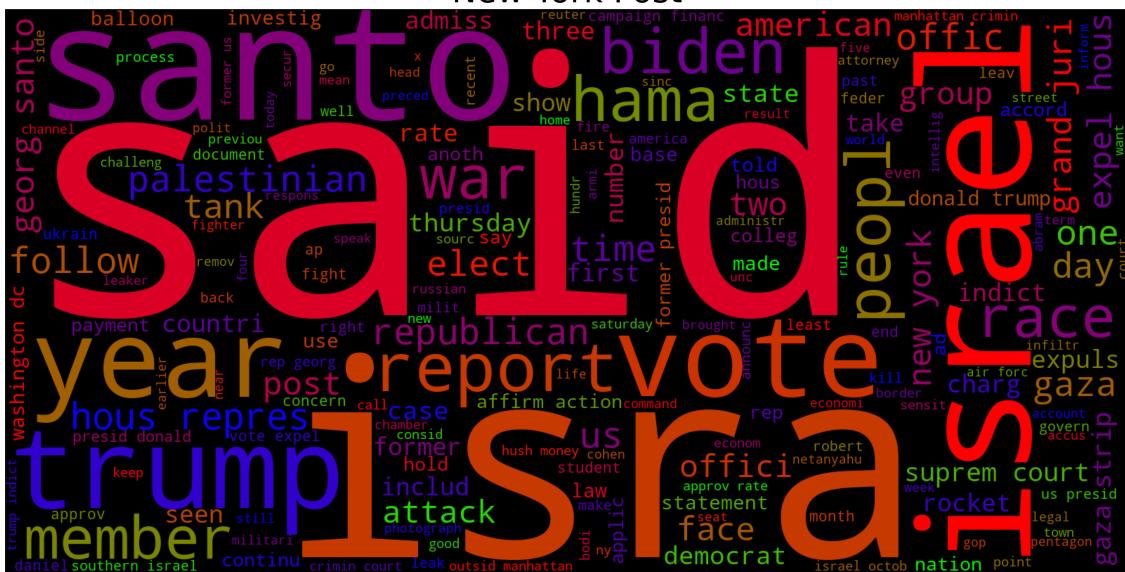
Fox News



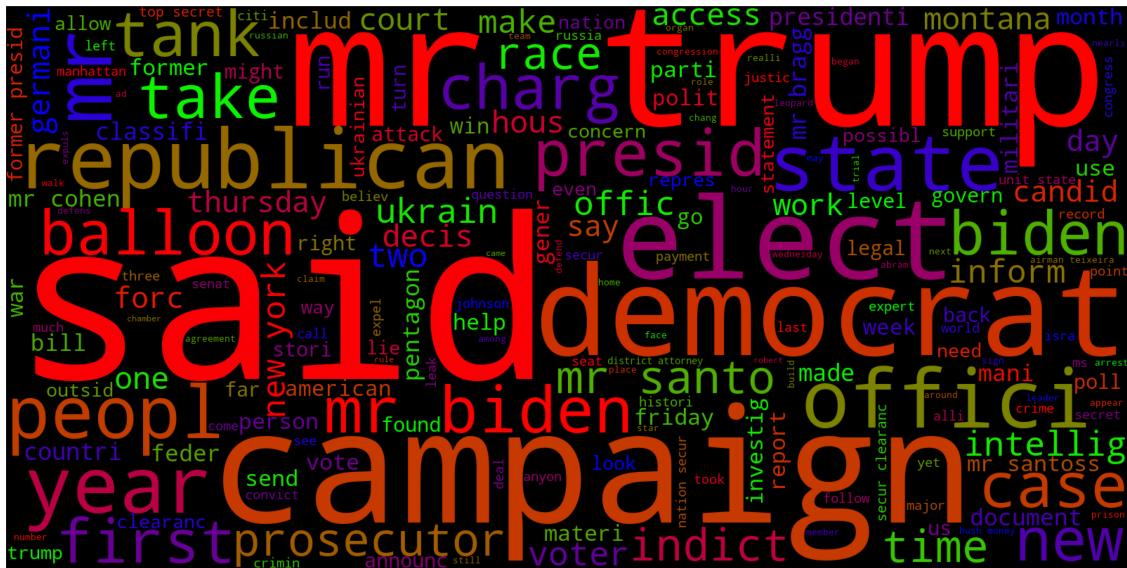
NBC News



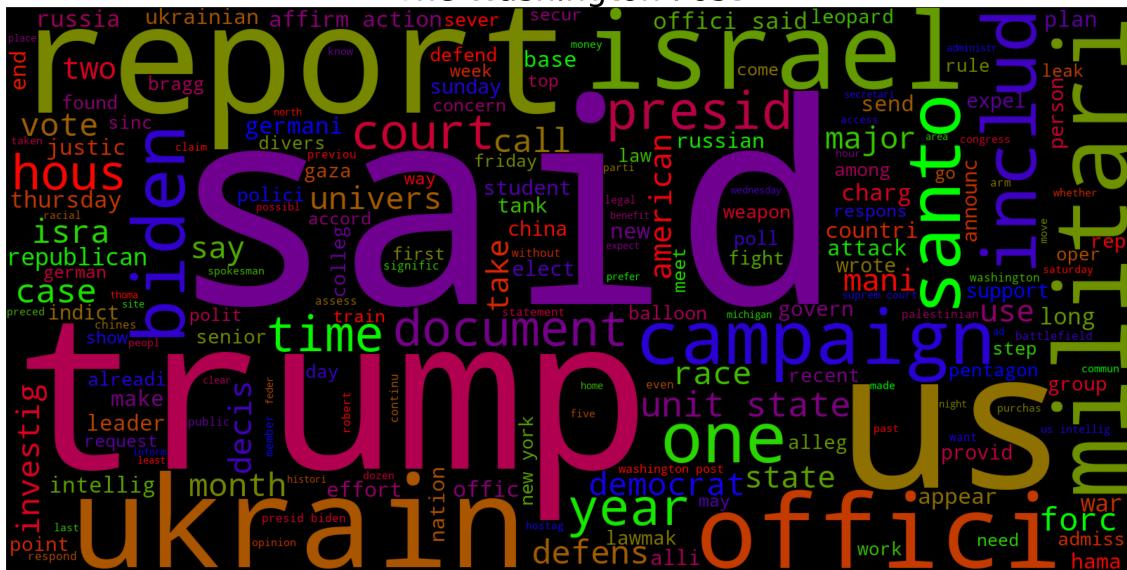
New York Post



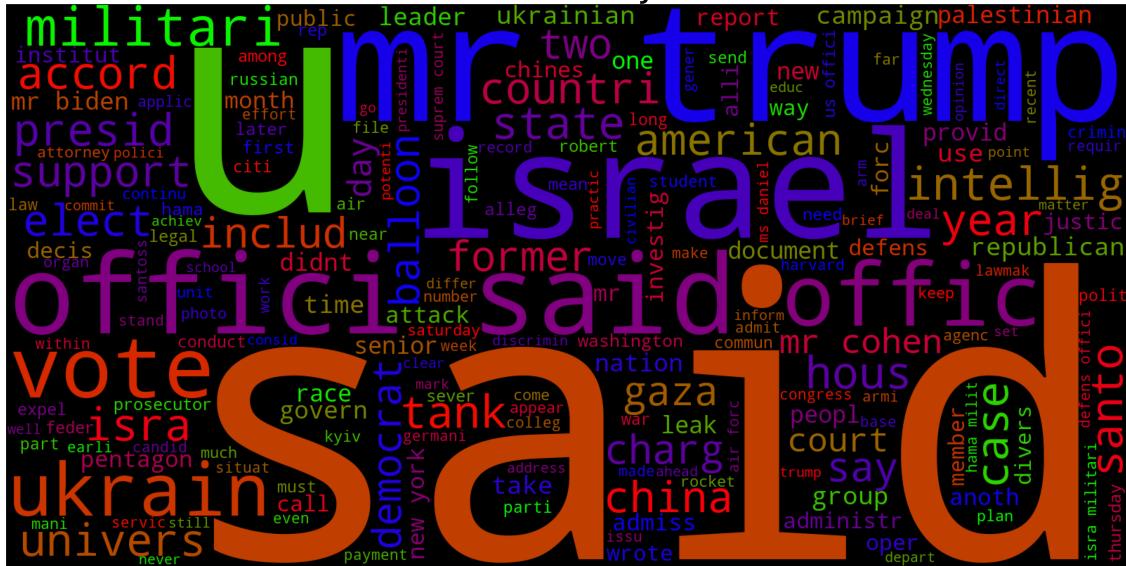
The New York Times



The Washington Post



## The Wall Street Journal



```
# create a list of unstemmed article document texts
# looping through all text files to apply preprocessing functions
article_docs_unstemmed = []
dir = os.listdir('data/text/')
dir.sort()
for filename in dir:
    filepath = os.path.join('data/text/', filename)
    if filename.split('.')[ -1] == "txt":
        article_string = file_to_string(filepath)
        new_string = clean_text(article_string, 0)
        article_docs_unstemmed.append(new_string)

article_docs_unstemmed
```

['supreme court thursday set new limits affirmative action programs cases involving whether p  
'massive spy balloon believed china seen montana tracked flies across continental united sta  
'two days dismal new polling numbers president joe biden showed public views unfavorably ri  
'hundreds people israel reported dead thousands injured hamas militants fired rockets gaza  
'posting social media appears several highly classified u intelligence documents might begin  
'house representatives friday voted expel republican rep george santos historic move happened  
'major increase u support ukraine president joe biden signed sending abrams tanks war torn  
'manhattan grand jury indicted former president donald trump making first current former pr  
'us supreme court ruled race longer considered factor university admissions landmark ruling  
'us tracking suspected chinese surveillance balloon spotted flying sensitive sites recent da

'us president joe biden running election next year national opinion polls weak job approval  
'least people reported killed wounded israel palestinian militant group hamas launched biggest  
'documents include detailed accounts training provided ukraine foreign powers make dozens  
'us house representatives expelled congressman george santos following damning ethics report  
'us send powerful battle tanks ukraine joining germany sending vehicles support fight russia  
'former us president donald trump charged hush money payments made porn star presidential election  
'supreme court says colleges universities longer take race consideration specific basis grants  
'us tracking suspected chinese high altitude surveillance balloon continental united states  
'one third registered voters approve president joe bidens handling israeli palestinian conflict  
'gaza jerusalem cnn israels prime minister benjamin netanyahu declared country war saturday  
'man arrested fbi connection massive us classified documents leak charged boston friday unaccounted  
'house voted friday expel gop rep george santos historic vote makes new york congressman single  
'leaders united states germany announced wednesday send contingents tanks ukraine reversing  
'donald trump faces counts related business fraud indictment manhattan grand jury according  
'u supreme court handed major ruling affirmative action thursday rejecting use race factor  
'u government monitoring suspected chinese surveillance balloon moving northern states past  
'biden battles rough poll numbers fox news white house correspondent peter doocy latest press  
'hamas widespread coordinated attack may suggest outside help trey yingst foreign correspondents  
'bret baier intel suspect year old government worker special report bret baier anchor questions  
'house representatives voted expel scandal plagued rep george santos r n friday making first  
'german chancellor olaf scholz formally announced wednesday weeks stalling frustrating negotiations  
'former president donald trump indicted part manhattan district attorney office years long  
'washington supreme court thursday struck affirmative action programs university north carolina  
'u military monitoring suspected chinese surveillance balloon hovering northern u past days  
'different polls agree president joe bidens political standing lower barack obamas point election  
'ashkelon israel israel plunged chaos saturday palestinian militant group hamas launched dead  
'dozens leaked defense department classified documents posted online reveal details u spying  
'washington house voted overwhelmingly expel indicted rep george santos friday pulling curtain  
'ukraine set receive battle tanks germany western countries fierce debate exposed fissures  
'grand jury new york city voted thursday indict donald trump first time former u president  
'supreme court struck affirmative action programs harvard university university north carolina  
'us tracking chinese spy balloon floating northern part country days pentagon officials announced  
'president biden ready ring new year seeing numbers year old commander chief ends lower approach  
'jerusalem ap hamas militants fired thousands rockets sent dozens fighters israeli towns near  
'massachusetts air national guardsman jack teixeira leader discord group dozens sensitive us  
'bye george lying long island rep george santos r ny became sixth member ever expelled us house  
'weeks excuses foot dragging germans finally said yes transferring limited number leopard tanks  
'donald trump indicted manhattan grand jury thursday hush money payments made ahead election  
'chief justice john g roberts jr finished reading majority opinion supreme court chamber this  
'helena mont larry mayer newspaper photographer pointed camera sky wednesday began snapping  
'democrats battleground states growing increasingly anxious president bidens low approval rates  
'israels news sites compiling lists dead missing funerals taking place around country weekend  
'washington would year old national guardsman position access top secret documents begin drawn

'george santos new york republican congressman whose tapestry lies schemes made figure nation  
'precision military drill first germany united states announced wednesday agreed provide bat  
'manhattan grand jury indicted donald j trump thursday role paying hush money porn star acc  
'supreme court thursday held race conscious admissions programs harvard university north car  
'chinese surveillance balloon collecting intelligence continental united states right u offi  
'night president biden departed washington celebrate thanksgiving nantucket mass gathered c  
'sderot israel israel formally declared war palestinian militant group hamas sunday reeled s  
'saturday u officials foreign allies scrambled understand dozens classified intelligence doc  
'house voted friday expel rep george santos r n congress action chamber previously taken fir  
'biden administration announced wednesday send premier battle tanks ukraine following agree  
'new york manhattan grand jury voted indict former president donald trump making first perso  
'thursdays decision force reworking admissions criteria throughout american higher education  
'washington u tracked officials described chinese reconnaissance balloon continental states  
'mr bidens low standing head head general election polls reflects voters dour appraisal per  
'tel aviv israeli prime minister benjamin netanyahu said country war hamas militant groups :  
'criminal case unfolding u government scrambles protect secrets unauthorized disclosures app  
'lawmakers voted remove two thirds house supermajority required constitution almost democrati  
'u germany outlined plans wednesday send dozens modern battle tanks ukraine marking signific  
'grand jury returned indictment mr trump vote thursday kicking process former president exp

```
# not using this for anything as of 3/14, don't remember why i created it
# create a list of strings where each string is all articles from one source (unstemmed)
source_docs_unstemmed = []

j = 0

for i in range(9):
    source = " ".join(article_docs_unstemmed[j].split()[:-2]) + " " + ".join(article_docs_
        + " ".join(article_docs_unstemmed[j+2].split()[:-2]) + " " + ".join(article_docs_
        + " ".join(article_docs_unstemmed[j+4].split()[:-2]) + " " + ".join(article_docs_
        + " ".join(article_docs_unstemmed[j+6].split()[:-2]) + " " + ".join(article_docs_
    source_docs_unstemmed.append(source)
    j += 8

source_docs_unstemmed

['supreme court thursday set new limits affirmative action programs cases involving whether p
'us supreme court ruled race longer considered factor university admissions landmark ruling
'supreme court says colleges universities longer take race consideration specific basis gran
'u supreme court handed major ruling affirmative action thursday rejecting use race factor o
'washington supreme court thursday struck affirmative action programs university north caro
```

```

'supreme court struck affirmative action programs harvard university university north carolina
'chief justice john g roberts jr finished reading majority opinion supreme court chamber the
'supreme court thursday held race conscious admissions programs harvard university north carolina
'thursdays decision force reworking admissions criteria throughout american higher education

first = article_docs_unstemmed[0]
" ".join(first.split()[:-2])

'supreme court thursday set new limits affirmative action programs cases involving whether programs

# calculate polarity and subjectivity scores, create dataframe

scores_df = pd.DataFrame({'source':tfidf_df['article_source'], 'topic':tfidf_df['article_topic']})

polarity_scores = []
subjectivity_scores = []

for article in article_docs_unstemmed:
    polarity_scores.append(round(TextBlob(" ".join(article.split()[:-2])).sentiment.polarity))
    subjectivity_scores.append(round(TextBlob(" ".join(article.split()[:-2])).sentiment.subjectivity))

scores_df['polarity_score'] = polarity_scores
scores_df['subjectivity_score'] = subjectivity_scores

average_topic_polarity_scores = []
average_source_polarity_scores = []

for topic in scores_df['topic'].value_counts().index:
    mean_score = round(scores_df[scores_df['topic'] == topic]['polarity_score'].mean(), 2)
    average_topic_polarity_scores.append(mean_score)

for source in scores_df['source'].value_counts().index:
    mean_score = round(scores_df[scores_df['source'] == source]['polarity_score'].mean(), 2)
    for i in range(8):
        average_source_polarity_scores.append(mean_score)

scores_df['average_polarity_for_topic'] = average_topic_polarity_scores * 9
scores_df['average_polarity_for_source'] = average_source_polarity_scores

average_topic_subjectivity_scores = []

```

```

average_source_subjectivity_scores = []

for topic in scores_df['topic'].value_counts().index:
    mean_score = round(scores_df[scores_df['topic'] == topic]['subjectivity_score'].mean())
    average_topic_subjectivity_scores.append(mean_score)

for source in scores_df['source'].value_counts().index:
    mean_score = round(scores_df[scores_df['source'] == source]['subjectivity_score'].mean())
    for i in range(8):
        average_source_subjectivity_scores.append(mean_score)

scores_df['average_subjectivity_for_topic'] = average_topic_subjectivity_scores * 9
scores_df['average_subjectivity_for_source'] = average_source_subjectivity_scores

scores_df['polarity_diff_from_topic_mean'] = scores_df['polarity_score'] - scores_df['average_polarity']
scores_df['subjectivity_diff_from_topic_mean'] = scores_df['subjectivity_score'] - scores_df['average_subjectivity']
scores_df['polarity_diff_from_source_mean'] = scores_df['polarity_score'] - scores_df['average_polarity']
scores_df['subjectivity_diff_from_source_mean'] = scores_df['subjectivity_score'] - scores_df['average_subjectivity']

sources_polarity_dev = []

for source in scores_df['source'].value_counts().index:
    total_dev = scores_df[scores_df['source'] == source]['polarity_diff_from_topic_mean'].sum()
    for i in range(8):
        sources_polarity_dev.append(total_dev/8)

sources_subjectivity_dev = []

for source in scores_df['source'].value_counts().index:
    total_dev = scores_df[scores_df['source'] == source]['subjectivity_diff_from_topic_mean'].sum()
    for i in range(8):
        sources_subjectivity_dev.append(total_dev/8)

scores_df['average_source_polarity_deviation_from_topic_mean'] = sources_polarity_dev
scores_df['average_source_subjectivity_deviation_from_topic_mean'] = sources_subjectivity_dev

scores_df.head(10)

```

	source	topic	polarity_score	subjectivity_score	average_
0	ABC News	Supreme Court Ruling on Affirmative Action	0.13	0.42	0.11
1	ABC News	Chinese Surveillance Balloon	0.05	0.36	0.04

	source	topic	polarity_score	subjectivity_score	average_
2	ABC News	Biden's Low Approval Rates in Polls	0.06	0.43	0.08
3	ABC News	The Deadliest Attack by Hamas	0.10	0.39	0.02
4	ABC News	Pentagon Documents Leak	0.03	0.32	0.05
5	ABC News	George Santos' Expulsion from Congress	0.03	0.49	0.02
6	ABC News	U.S. and Germany Send Tanks to Ukraine	0.08	0.40	0.07
7	ABC News	Trump's Indictment	0.03	0.40	0.04
8	BBC	Supreme Court Ruling on Affirmative Action	0.16	0.37	0.11
9	BBC	Chinese Surveillance Balloon	0.06	0.38	0.04

```
# create a df for each source and its polarity score total deviations from the topic means

polarity_devs = pd.DataFrame({'source':scores_df['source'].value_counts().index,
                               'polarity_dev':scores_df[scores_df['topic'] == 'Chinese Surveillance Balloon']['average_source_polarity_deviation_from_topic_mean']})
polarity_devs['sign'] = np.where(polarity_devs['polarity_dev'] > 0, 'pos', 'neg')

polarity_devs = polarity_devs.reindex(polarity_devs['polarity_dev'].abs().sort_values(ascending=True))
polarity_devs.reset_index(drop=True, inplace=True)

source_order = CategoricalDtype(
    ["The New York Times", "The Wall Street Journal", "ABC News", "NBC News",
     "CNN", "New York Post", "The Washington Post", "Fox News", "BBC"],
    ordered=False
)
polarity_devs['source'] = polarity_devs['source'].astype(source_order)

polarity_devs
```

	source	polarity_dev	sign
0	The New York Times	-0.01625	neg
1	ABC News	0.01000	pos
2	The Wall Street Journal	-0.01000	neg
3	NBC News	0.00750	pos
4	CNN	0.00750	pos
5	New York Post	-0.00375	neg
6	The Washington Post	-0.00375	neg
7	Fox News	0.00375	pos
8	BBC	0.00125	pos

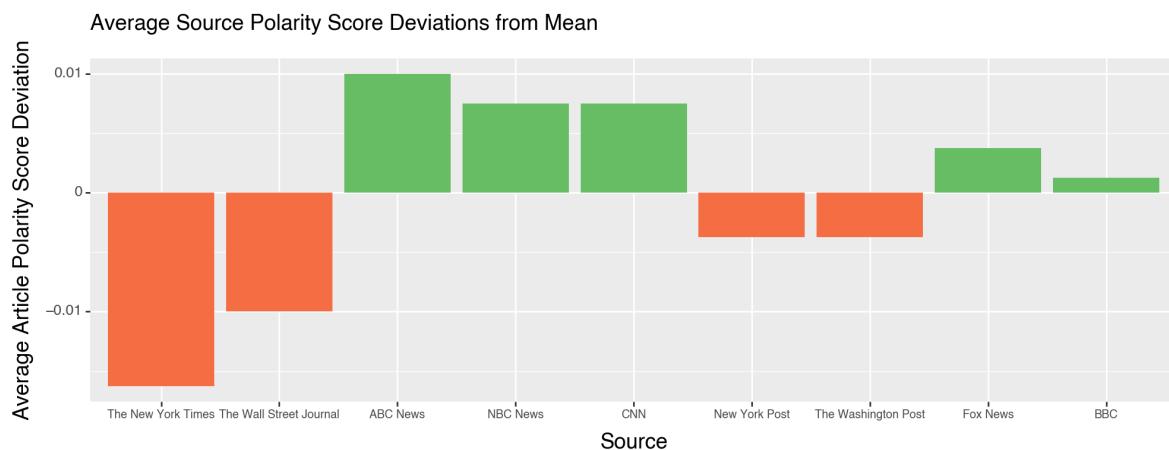
```

# create bar chart representing how much, on average, a source's article
# deviates away from that topic's mean polarity score
# intended to show if a source tends to be more negative/positive than average,
# even accounting for the nature of the topic

colors = {"pos": '#66bd63', "neg": '#f46d43'}

(
ggplot(polarity_devs, aes(x="source", y="polarity_dev", fill="sign"))
+ geom_col(stat="identity")
+ labs(x="Source", y='Average Article Polarity Score Deviation', title='Average Source Pol')
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=7, face='bold'),
       title = element_text(size=13),legend_position = "none")
+ scale_fill_manual(values = colors)
)

```

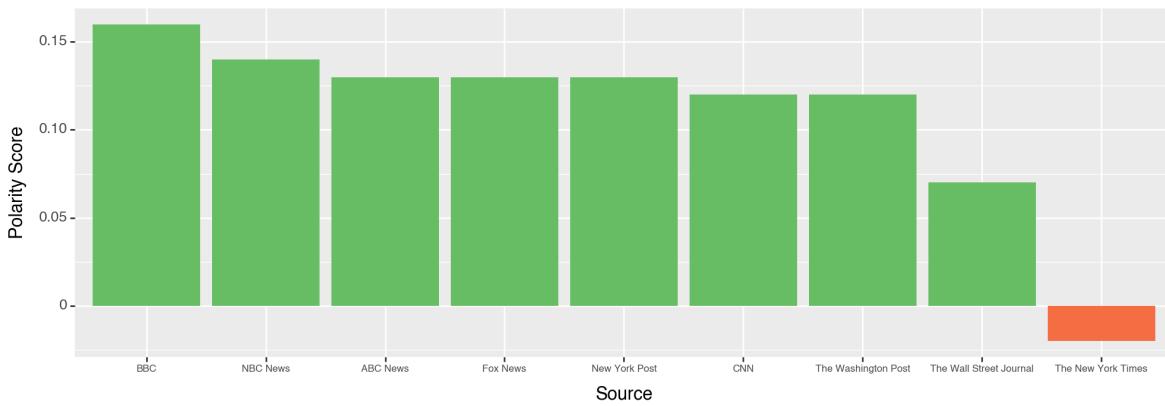


<Figure Size: (1000 x 400)>

```

(
ggplot(scores_df[scores_df['topic']=="Supreme Court Ruling on Affirmative Action"],
       aes(x="reorder(source, -polarity_score)", y="polarity_score", fill="polarity_score"))
+ geom_col(stat="identity")
+ labs(x='Source', y='Polarity Score', title='Polarity Scores for Affirmative Action Artic')
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_posit
+ scale_fill_manual(values = colors)
)
```

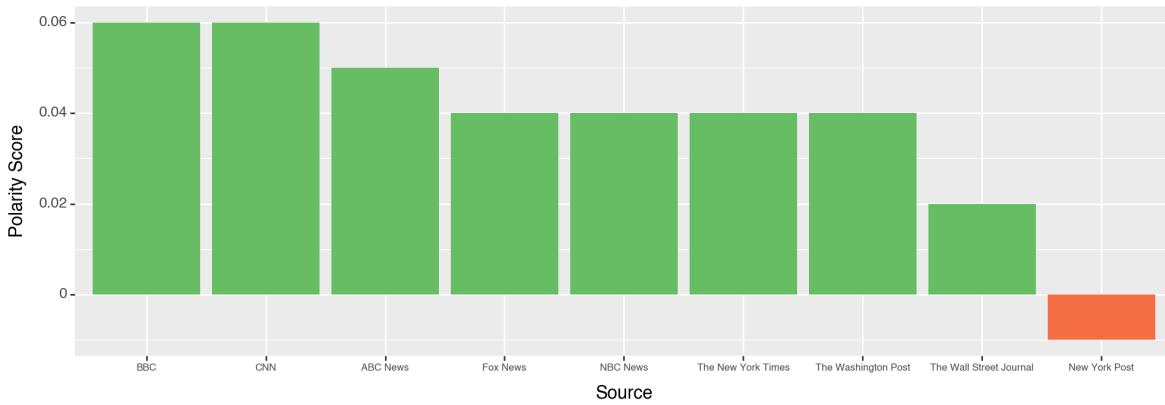
Polarity Scores for Affirmative Action Articles



<Figure Size: (1000 x 400)>

```
(  
  ggplot(scores_df[scores_df['topic']=="Chinese Surveillance Balloon"],  
         aes(x=reorder(source, -polarity_score), y=polarity_score, fill=polarity_score_<  
+ geom_col(stat="identity")  
+ labs(x='Source', y='Polarity Score', title='Polarity Scores for Chinese Surveillance Ballo<  
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_posit<  
+ scale_fill_manual(values = colors)  
)
```

Polarity Scores for Chinese Surveillance Balloon Articles

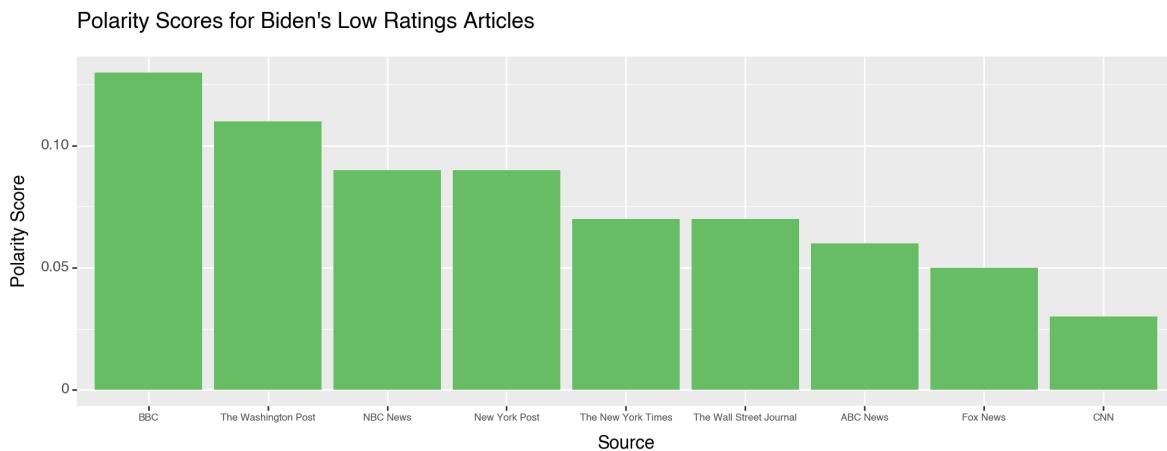


<Figure Size: (1000 x 400)>

```

(
ggplot(scores_df[scores_df['topic']=="Biden's Low Approval Rates in Polls"],
       aes(x="reorder(source, -polarity_score)", y="polarity_score", fill="polarity_score",
+ geom_col(stat="identity")
+ labs(x='Source', y='Polarity Score', title="Polarity Scores for Biden's Low Ratings Articles",
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_position="none"),
+ scale_fill_manual(values = colors)
)

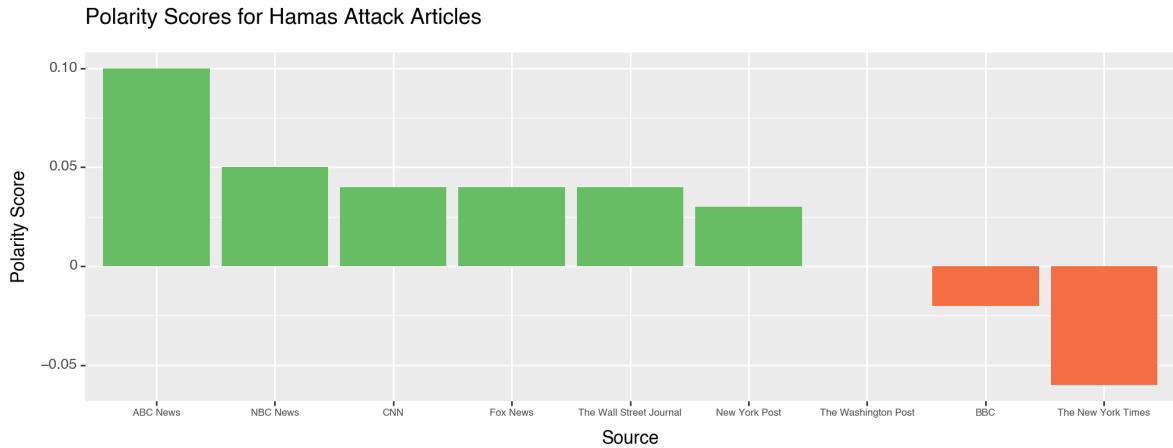
```



<Figure Size: (1000 x 400)>

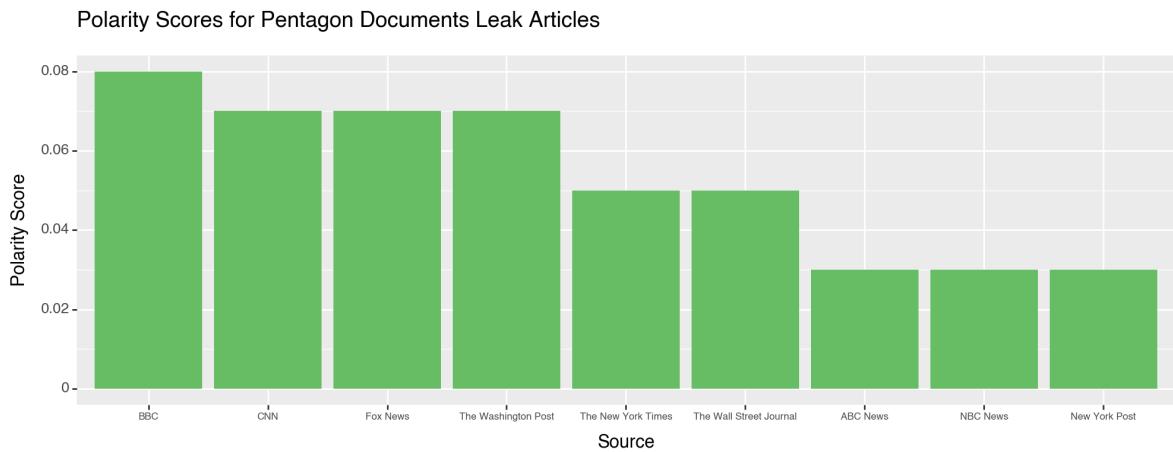
```

(
ggplot(scores_df[scores_df['topic']=="The Deadliest Attack by Hamas"],
       aes(x="reorder(source, -polarity_score)", y="polarity_score", fill="polarity_score",
+ geom_col(stat="identity")
+ labs(x='Source', y='Polarity Score', title="Polarity Scores for Hamas Attack Articles"),
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_position="none"),
+ scale_fill_manual(values = colors)
)
```



<Figure Size: (1000 x 400)>

```
(  
  ggplot(scores_df[scores_df['topic']=="Pentagon Documents Leak"],  
         aes(x=reorder(source, -polarity_score), y=polarity_score, fill=polarity_score_<  
+ geom_col(stat="identity")  
+ labs(x='Source', y='Polarity Score', title="Polarity Scores for Pentagon Documents Leak  
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_posit  
+ scale_fill_manual(values = colors)  
)
```

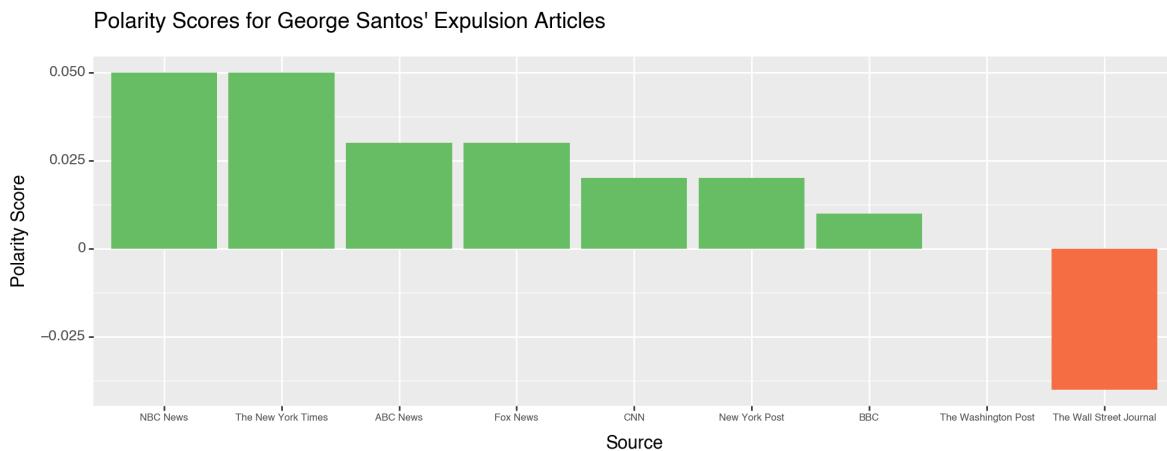


<Figure Size: (1000 x 400)>

```

(
ggplot(scores_df[scores_df['topic']=="George Santos' Expulsion from Congress"],
       aes(x="reorder(source, -polarity_score)", y="polarity_score", fill="polarity_score_"
+ geom_col(stat="identity")
+ labs(x='Source', y='Polarity Score', title="Polarity Scores for George Santos' Expulsion"
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_posit
+ scale_fill_manual(values = colors)
)

```

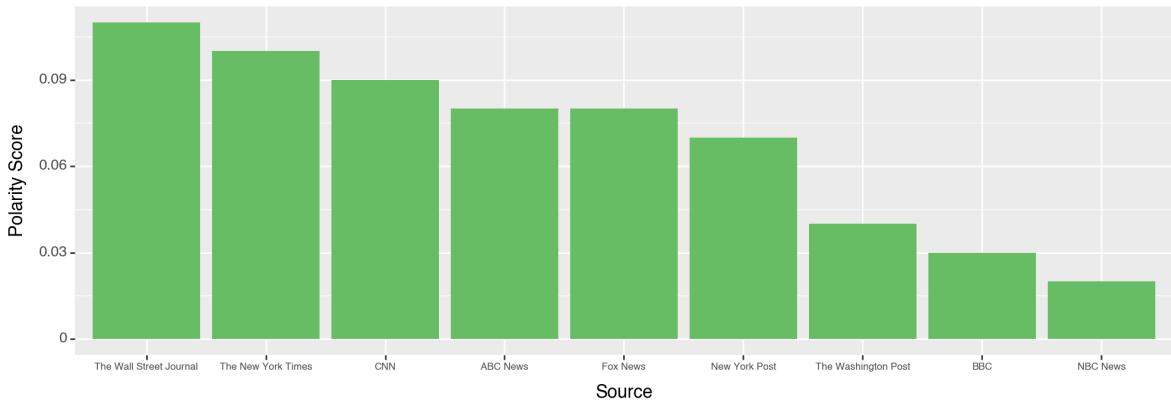


<Figure Size: (1000 x 400)>

```

(
ggplot(scores_df[scores_df['topic']=="U.S. and Germany Send Tanks to Ukraine"],
       aes(x="reorder(source, -polarity_score)", y="polarity_score", fill="polarity_score_"
+ geom_col(stat="identity")
+ labs(x='Source', y='Polarity Score', title="Polarity Scores for Sending Tanks to Ukraine"
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_posit
+ scale_fill_manual(values = colors)
)
```

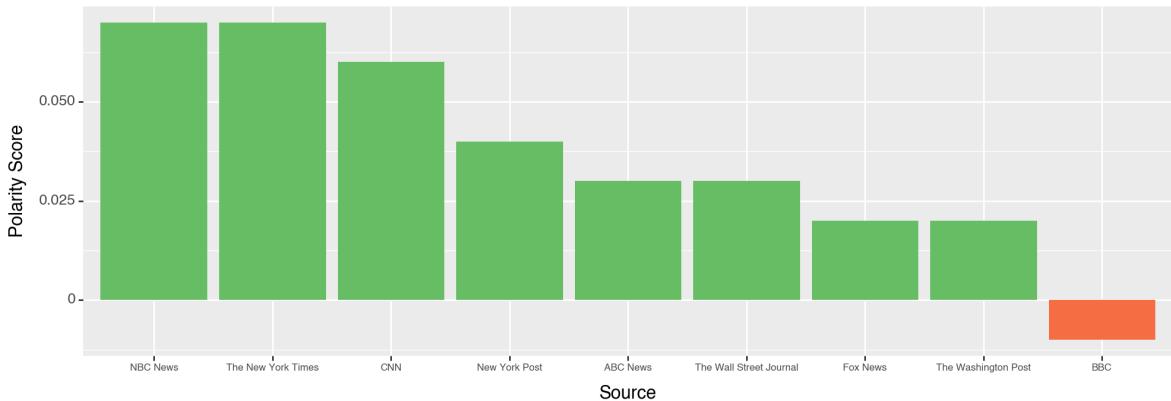
Polarity Scores for Sending Tanks to Ukraine Articles



<Figure Size: (1000 x 400)>

```
(  
  ggplot(scores_df[scores_df['topic']=="Trump's Indictment"],  
         aes(x=reorder(source, -polarity_score), y=polarity_score, fill=polarity_score_  
+ geom_col(stat="identity")  
+ labs(x='Source', y='Polarity Score', title="Polarity Scores for Trump's Indictment Artic  
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=6, face='bold'), legend_posit  
+ scale_fill_manual(values = colors)  
)
```

Polarity Scores for Trump's Indictment Articles



<Figure Size: (1000 x 400)>

Do the bar charts above mean for example that ABC News was portraying the attack by hamas in a positive light? no, this is just supposed to give insight into the source's word choices.

```
# check if normally distributed

# hypothesis test for difference in mean polarity score deviations

import statsmodels.api as sm
from statsmodels.formula.api import ols

# formula = 'mood_gain ~ drug'

# model = ols(formula, data=df).fit()

# aov_table = sm.stats.anova_lm(model, typ=2)
# aov_table

source_avg_subj = pd.DataFrame({'source':scores_df['source'].value_counts().index,
                                'avg_subj':scores_df[scores_df['topic'] == 'Chinese Surveillance']['average_subjectivity_for_source']})
source_avg_subj.sort_values(by='avg_subj', ascending=True, inplace=True)
source_avg_subj.reset_index(drop=True, inplace=True)

source_order = CategoricalDtype(
    ["Fox News", "The Wall Street Journal", "CNN",
     "The Washington Post", "BBC", "New York Post", "ABC News", "NBC News", "The New York
     ordered=False
)
source_avg_subj['source'] = source_avg_subj['source'].astype(source_order)

source_avg_subj
```

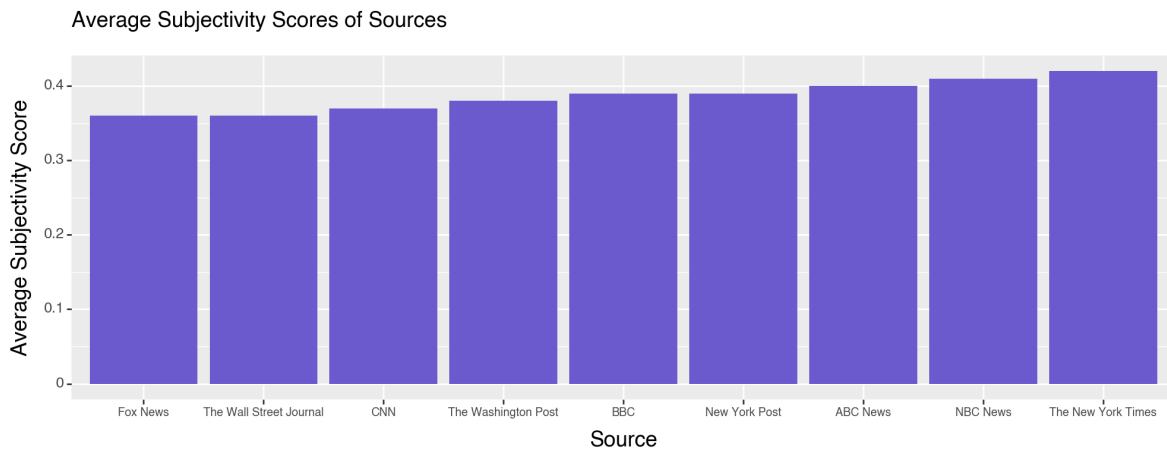
	source	avg_subj
0	Fox News	0.36
1	The Wall Street Journal	0.36
2	CNN	0.37
3	The Washington Post	0.38
4	BBC	0.39
5	New York Post	0.39
6	ABC News	0.40
7	NBC News	0.41

source	avg_subj
8 The New York Times	0.42

```

(
ggplot(source_avg_subj, aes(x="source", y="avg_subj"))
+ geom_col(stat="identity", fill="slateblue")
+ labs(x='Source', y='Average Subjectivity Score', title='Average Subjectivity Scores of Sources')
+ theme(figure_size = (10, 4), axis_text_x=element_text(size=7, face='bold'),
       title = element_text(size=13), legend_position = "none")
+ scale_fill_manual(values = colors)
)

```



<Figure Size: (1000 x 400)>

```

# create a list of tuples that contain articles that have an outlier polarity score,
# compared to other articles about the same topic

articles = []

for topic in scores_df['topic'].value_counts().index:
    topic_df = scores_df[scores_df['topic'] == topic]
    Q1 = topic_df['polarity_score'].describe()[4]
    Q3 = topic_df['polarity_score'].describe()[6]
    IQR = Q3-Q1
    lower = Q1-(1.5*IQR)

```

```
upper = Q3+(1.5*IQR)
print(topic, "lower bound:", round(lower, 2), "upper bound:", round(upper,2))
for row in topic_df.itertuples():
    if row.polarity_score < lower or row.polarity_score > upper:
        articles.append((row.source, row.topic, row.polarity_score))

print(articles)
print(len(articles))
```

Supreme Court Ruling on Affirmative Action lower bound: 0.1 upper bound: 0.15  
Chinese Surveillance Balloon lower bound: 0.02 upper bound: 0.06  
Biden's Low Approval Rates in Polls lower bound: 0.02 upper bound: 0.14  
The Deadliest Attack by Hamas lower bound: -0.06 upper bound: 0.1  
Pentagon Documents Leak lower bound: -0.03 upper bound: 0.13  
George Santos' Expulsion from Congress lower bound: -0.02 upper bound: 0.06  
U.S. and Germany Send Tanks to Ukraine lower bound: -0.03 upper bound: 0.16  
Trump's Indictment lower bound: -0.04 upper bound: 0.12  
[('BBC', 'Supreme Court Ruling on Affirmative Action', 0.16), ('The New York Times', 'Supreme  
6

By the outlier detection above, there are 6 articles that are considered to have an outlier polarity score for their topic. Three of those are articles from the Wall Street Journal. Also, three of the outlier articles are about the Supreme Court Ruling on Affirmative Action.

```
freq_standardized.head()
```

	said	trump	biden	offici	mr	u	israel	presid	tank
0	-0.917117	-0.451784	-0.576257	-0.630037	-0.323434	-0.026389	-0.395549	-0.828123	-0.367273
1	2.263636	-0.451784	0.076109	3.440971	-0.323434	1.581300	-0.395549	-0.014132	-0.367273
2	-0.408197	0.747746	2.946523	-0.630037	-0.323434	-0.611003	-0.395549	0.799860	-0.367273
3	0.609644	-0.451784	-0.184837	0.242322	-0.323434	0.850532	2.253711	-0.421127	-0.367273
4	-0.408197	-0.451784	-0.576257	0.096929	-0.323434	2.312067	-0.206316	-0.421127	-0.076683

```
# standardize the frequency df before doing PCA
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
scaler = StandardScaler()

freq_standardized = pd.DataFrame(scaler.fit_transform(freq_df.iloc[:, :-2]), columns = freq_df.columns)
#freq_standardized = preprocessing.scale(freq_df.iloc[:, :-2], axis=0) # standardize data

## # Instantiate PCA estimator
# pca = decomposition.PCA()
## # fit : Run PCA
# tokens_pc = pca.fit(freq_standardized)
# vars(tokens_pc)

# Defining the number of principal components to generate
n = freq_standardized.shape[1]

# Finding principal components for the data
pca = PCA(n_components = 72, random_state = 1) # Apply the PCA algorithm with random_state

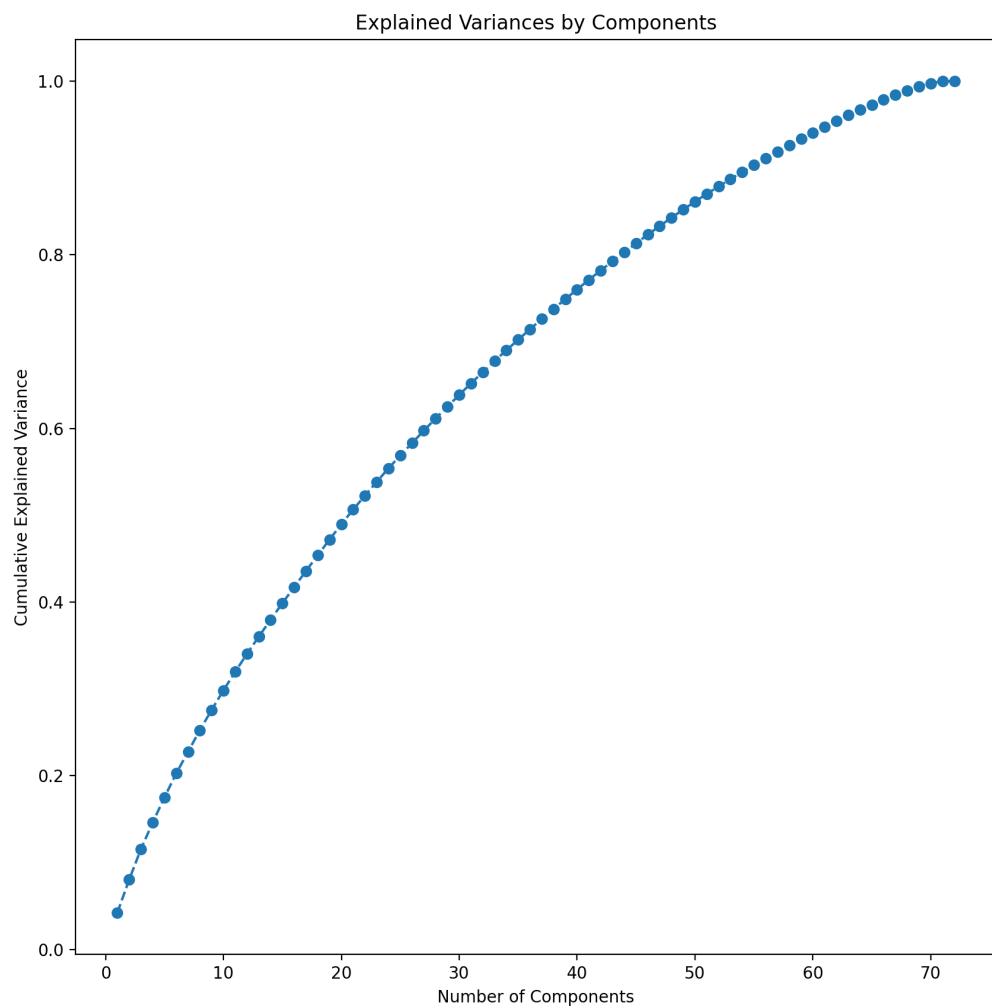
data_pca1 = pd.DataFrame(pca.fit_transform(freq_standardized)) # Fit and transform the p

# The percentage of variance explained by each principal component
exp_var = pca.explained_variance_ratio_

# Visualize the explained variance by individual components
plt.figure(figsize = (10, 10))

plt.plot(range(1, 73), exp_var.cumsum(), marker = 'o', linestyle = '--')
```

```
plt.title("Explained Variances by Components")  
plt.xlabel("Number of Components")  
plt.ylabel("Cumulative Explained Variance")  
plt.show()
```



```
# Finding the least number of components that can explain more than 90% variance
sum = 0

for ix, i in enumerate(exp_var):
    sum = sum + i
    if(sum>0.90):
        print("Number of PCs that explain at least 90% variance: ", ix + 1)
        break

Number of PCs that explain at least 90% variance:  55

PCs = pd.DataFrame(tokens_pc.components_.T[:10],
                    columns=['PC1', 'PC2','PC3','PC4', 'PC5', 'PC6', 'PC7','PC8','PC9', 'PC10'])

PCs           # PC loadings saved with Scikit-learn

ValueError: Shape of passed values is (10, 72), indices imply (10, 10)
```