you will have questions to answer. The questions will be marked by a "Q:" and will have a corresponding "A:" spot for you. Make sure to answer every question marked with a Q: for full credit. In [1]: import os import re import emoji import pandas as pd from collections import Counter, defaultdict from nltk.corpus import stopwords from string import punctuation from wordcloud import WordCloud from sklearn.feature extraction.text import TfidfTransformer, CountVectorizer In [2]: # Use this space for any additional import statements you need In [31]: # Place any additional functions or constants you need here. # Some punctuation variations punctuation = set(punctuation) # speeds up comparison tw punct = punctuation - {"#"} # Stopwords sw = stopwords.words("english") # Two useful regex whitespace_pattern = re.compile($r"\s+"$) hashtag pattern = re.compile($r"^{\#}[0-9a-zA-Z]+"$) # It's handy to have a full set of emojis all language emojis = set() for country in emoji.EMOJI DATA : for em in emoji.EMOJI DATA[country] : all language emojis.add(em) # and now our functions def descriptive stats(tokens, num tokens = 5, verbose=True) : Given a list of tokens, print number of tokens, number of unique tokens, number of characters, lexical diversity (https://en.wikipedia.org/wiki/Lexical diversity), and num tokens most common tokens. Return a list with the number of tokens, number of unique tokens, lexical diversity, and number of characters. # Fill in the correct values here. num tokens = len(tokens) num_unique_tokens = len(set(tokens)) lexical diversity = num unique tokens / num tokens num_characters = len("".join(tokens)) if verbose : print(f"There are {num tokens} tokens in the data.") print(f"There are {num unique tokens} unique tokens in the data.") print(f"There are {num characters} characters in the data.") print(f"The lexical diversity is {lexical diversity:.3f} in the data.") # print the five most common tokens return([num tokens, num unique tokens, lexical diversity, num characters]) def contains emoji(s): emojis = [ch for ch in s if emoji.is emoji(ch)] return(len(emojis) > 0) def remove stop(tokens) : # modify this function to remove stopwords return[i for i in tokens if i not in sw] def remove punctuation(text, punct set=tw punct) : return("".join([ch for ch in text if ch not in punct set])) def tokenize(text) : """ Splitting on whitespace rather than the book's tokenize function. That function will drop tokens like '#hashtag' or '2A', which we need for Twitter. """ return([item.lower() for item in whitespace pattern.split(text)]) def prepare(text, pipeline) : tokens = str(text) for transform in pipeline : tokens = transform(tokens) return (tokens) **Data Ingestion** Use this section to ingest your data into the data structures you plan to use. Typically this will be a dictionary or a pandas In [32]: # Feel fre to use the below cells as an example or read in the data in a way you prefer data location = "/Users/ruddysimonpour/Desktop/University of Sandiego - Curriculum/ADS509 - Applied text analys twitter folder = "twitter/" lyrics_folder = "lyrics/" artist files = {'cher':'cher followers data.txt', 'robyn':'robynkonichiwa followers data.txt'} twitter data = pd.read csv(data location + twitter folder + artist files['cher'], sep="\t", quoting=3) twitter data['artist'] = "cher" In [34]: | twitter_data_2 = pd.read_csv(data_location + twitter_folder + artist_files['robyn'], $sep="\t",$ quoting=3) twitter_data_2['artist'] = "robyn" twitter_data = pd.concat([twitter_data, twitter_data_2]) del (twitter_data_2) In [35]: # read in the lyrics here lyrics data = {} artists = ['cher', 'robyn'] for artist in artists: lyrics data[artist] = {} song lists = os.listdir(data location + lyrics folder + str(artist)) for song in song lists: # Get the song titles and artist name by splitting the filename at the underscore song title = os.path.splitext(song)[0].split(' ')[1] with open(data location + lyrics folder + str(artist) + '/' + str(song), 'r', encoding='utf-8') as file: lyrics data[artist][song title] = file.read() In [36]: data = []for artist, songs in lyrics_data.items(): for song title, lyrics in songs.items(): data.append([artist, song_title.replace(' ', ' '), lyrics]) lyrics_data = pd.DataFrame(data, columns=['artist_name', 'song_title', 'lyrics']) lyrics_data.head() Out [36]: artist_name song_title "Come And Stay With Me"\n\n\n\nI'll send away ... cher comeandstaywithme cher pirate "Pirate"\n\n\nHe'll sail on with the summer ... 2 cher "Stars"\n\n\n\nI was never one for saying what... stars 3 cher thesedays "These Days" $\n\n\$ lovesohigh "Love So High"\n\n\nEvery morning I would wa... cher **Tokenization and Normalization** In this next section, tokenize and normalize your data. We recommend the following cleaning. Lyrics Remove song titles Casefold to lowercase Remove stopwords (optional) Remove punctuation Split on whitespace Removal of stopwords is up to you. Your descriptive statistic comparison will be different if you include stopwords, though TF-IDF should still find interesting features for you. Note that we remove stopwords before removing punctuation because the stopword set includes punctuation. **Twitter Descriptions** Casefold to lowercase Remove stopwords Remove punctuation other than emojis or hashtags Split on whitespace Removing stopwords seems sensible for the Twitter description data. Remember to leave in emojis and hashtags, since you analyze those. # apply the `pipeline` techniques from BTAP Ch 1 or 5 In [37]: my_pipeline = [str.lower, remove_punctuation, tokenize, remove_stop] lyrics_data["tokens"] = lyrics_data["lyrics"].apply(prepare,pipeline=my_pipeline) lyrics data["num tokens"] = lyrics data["tokens"].map(len) twitter_data["tokens"] = twitter_data["description"].apply(prepare,pipeline=my_pipeline) twitter_data["num_tokens"] = twitter_data["tokens"].map(len) In [39]: twitter_data['has_emoji'] = twitter_data["description"].apply(contains_emoji) Let's take a quick look at some descriptions with emojis. twitter data[twitter data.has emoji].sample(10)[["artist","description","tokens"]] In [40]: Out[40]: artist description tokens 3034606 28, Lab Scientist 👮 ,homo 🌈 , LSU 🐺 . Thats w... [28, lab, scientist 👮, homo, 🌈, lsu, 😽, that... cher 3467056 [#captainmarvel, #roguexmen, #scarletwitch, #p... cher #CaptainMarvel #RogueXmen #ScarletWitch #Potte... 2987624 cher Once Upon A Time, An Angel And A Devil Fell In ... [upon, timean, angel, devil, fell, love, ♥♥, ♥♥] 2776387 Lahey Clinic Burlington Single Mom, Master Est... [lahey, clinic, burlington, single, mom, maste... cher 1498158 cher 📸 snap: kybug9752 Instagram kyleigh_nelms🤭 [🚵, snap, kybug9752, instagram, kyleighnelms😁] 238933 cher Producer Feminist https://t.co/NRxSyLG8Fb [producer, w, feminist, httpstconrxsylg8fb] 2053736 y justo cuando la oruga pensaba q el mundo se ... cher [justo, cuando, la, oruga, pensaba, q, el, mun... 1829538 [live, life, dare, 🐇 💗 instagram, joeyheide] cher I live my life, because I dare. 🤞 💗 Instagram:... 2434013 aspiring youtuber \odot i will film some videos s... [aspiring, youtuber, @, film, videos, soon, h... cher 273434 cher Why are we reliant on the economy aka the gove... [reliant, economy, aka, government, 💯 👇 , blog, ... With the data processed, we can now start work on the assignment questions. Q: What is one area of improvement to your tokenization that you could theoretically carry out? (No need to actually do it; let's not make perfect the enemy of good enough.) A: Emojis and special characters like 👮 and 🌈 are being treated as individual tokens. While this may be desirable in some cases, in other cases, you might want to remove emojis or encode them differently, perhaps converting them to a word representation of their meaning, especially if they do not contribute to the analysis. Calculate descriptive statistics on the two sets of lyrics and compare the results. In [41]: import nltk import pandas as pd nltk.download('punkt') artists = ['cher', 'robyn'] for artist in artists: artist data = lyrics data[lyrics data['artist name'] == artist] all lyrics = " ".join(artist data['lyrics']) tokens = nltk.word tokenize(all lyrics) print(f"Descriptive Statistics for {artist.capitalize()}:") descriptive stats(tokens) print("\n") [nltk data] Downloading package punkt to [nltk_data] /Users/ruddysimonpour/nltk_data...
[nltk_data] Package punkt is already up-to-date! Descriptive Statistics for Cher: There are 79284 tokens in the data. There are 4676 unique tokens in the data. There are 274999 characters in the data. The lexical diversity is 0.059 in the data. Descriptive Statistics for Robyn: There are 33812 tokens in the data. There are 2675 unique tokens in the data. There are 117506 characters in the data. The lexical diversity is 0.079 in the data. Q: what observations do you make about these data? A: Cher has a substantially higher volume of content compared to Robyn, with 79,284 tokens and 274,999 characters. Cher utilizes 4,676 unique tokens, which is notably more than Robyn's 2,675 unique tokens. The lexical diversity (the ratio of unique tokens to the total number of tokens) is higher for Robyn at 0.079 compared to Cher at 0.059. This suggests that, relative to the volume of content, Robyn has a more varied vocabulary. Cher's lower lexical diversity score indicates a higher degree of word repetition compared to Robyn, which can be interpreted as Cher using a more focused or limited set of vocabulary in her content. Find tokens uniquely related to a corpus Typically we would use TF-IDF to find unique tokens in documents. Unfortunately, we either have too few documents (if we view each data source as a single document) or too many (if we view each description as a separate document). In the latter case, our problem will be that descriptions tend to be short, so our matrix would be too sparse to support analysis. To avoid these problems, we will create a custom statistic to identify words that are uniquely related to each corpus. The idea is to find words that occur often in one corpus and infrequently in the other(s). Since corpora can be of different lengths, we will focus on the concentration of tokens within a corpus. "Concentration" is simply the count of the token divided by the total corpus length. For instance, if a corpus had length 100,000 and a word appeared 1,000 times, then the concentration would be $\frac{1000}{100000}=0.01$. If the same token had a concentration of 0.005 in another corpus, then the concentration ratio would be $\frac{0.01}{0.005}=2$. Very rare words can easily create infinite ratios, so you will also add a cutoff to your code so that a token must appear at least n times for you to return it. An example of these calculations can be found in this spreadsheet. Please don't hesitate to ask questions if this is confusing. In this section find 10 tokens for each of your four corpora that meet the following criteria: 1. The token appears at least n times in all corpora 2. The tokens are in the top 10 for the highest ratio of appearances in a given corpora vs appearances in other corpora. You will choose a cutoff for yourself based on the side of the corpus you're working with. If you're working with the Robyn-Cher corpora provided, n=5 seems to perform reasonably well. In [42]: def calculate token statistics(tokens, min count=5): count corpora = {'token': [], 'count': [], 'concentration': []} counter = Counter(tokens) for t, count in counter.items(): if count >= min count: count corpora['token'].append(t) count corpora['count'].append(count) concentration = count / len(tokens) count corpora['concentration'].append(concentration) return pd.DataFrame(count corpora) def flatten nested list(nested list): return [element for inner list in nested list for element in inner list] def get tokens by artist(lyrics data, artist): tokens = lyrics data.loc[lyrics data['artist name'] == artist, 'tokens'] tokens = flatten nested list(list(tokens)) return tokens def main(): cher tokens = get tokens by artist(lyrics data, 'cher') cher df = calculate token statistics(cher tokens) cher df.columns = ['token', 'corpus1 count', 'concentration1'] robyn tokens = get tokens by artist(lyrics data, 'robyn') robyn df = calculate token statistics(robyn tokens) robyn df.columns = ['token', 'corpus2 count', 'concentration2'] # Merge the two DataFrames final df = cher df.merge(robyn df, on='token') return final df main() In [43]: Out [43]: token corpus1 count concentration1 corpus2 count concentration2 0 come 270 0.007452 72 0.004696 1 stay 78 0.002153 15 0.000978 2 ill 170 0.004692 73 0.004762 3 196 0.005410 39 0.002544 awav 4 life 124 0.003422 39 0.002544 6 0.000166 8 0.000522 373 anytime 374 5 0.000138 5 0.000326 makin 375 45 0.001242 6 0.000391 boys **376** darkness 0.000248 0.000326 377 8 0.000221 16 0.001044 378 rows × 5 columns Q: What are some observations about the top tokens? Do you notice any interesting items on the list? A: Words like "come", "stay", "ill", "away", and "life" are common to both corpora, reflecting possibly common themes or topics discussed in the contents. The word "come" has a higher concentration in corpus1 than in corpus2, which might indicate different usage or importance of this word in the respective contents. Words like "ill" and "life" have similar concentrations in both corpora, suggesting a similar level of emphasis or relevance in both The variety in the concentration of different words between the two corpora suggests different styles, topics, or thematic elements. For instance, higher occurrences of words like "stay" and "away" in corpus1 may imply different narrative focuses or stylistic preferences compared to corpus2.

Build word clouds for all four corpora.

wc = WordCloud(width=800, height=400,

filter stop words in frequency counter

wc.generate from frequencies(counter)

plt.imshow(wc, interpolation='bilinear')

process tokens and update counter

transform counter into data frame

create counter and run through all data

freq df = freq df[freq df['freq'] >= min freq]

return freq df.sort values('freq', ascending=False)

counter.update(tokens)

freq df.index.name = 'token'

!pip install --upgrade Pillow

plt.figure(figsize=(10,5))

plt.title(f'{artist} Twitter Cloud')

producerque

convert data frame into dict if type(word freq) == pd.Series:

counter = word freq

if stopwords is not None:

plt.title(title)

plt.axis("off")

def update(doc):

counter = Counter() df[column].map(update)

In [45]: # !pip install --upgrade pip

In [46]: **from** wordcloud **import** WordCloud

for artist in artists:

for twitter

plt.axis('off') plt.show()

fashion

for artist in artists: # for lyrics

> plt.axis('off') plt.show()

gonha want

believe

thing

make

artists.

comeright

ant

betterkilling ^{yeah}

plt.figure(figsize=(10,5))

plt.title(f'{artist} Lyrics Cloud')

In [48]: min freq = 5

min freq = 5

artists = ['cher', 'robyn']

In [44]: from matplotlib import pyplot as plt

you should absolutely clone the repository that accompanies the book.

def wordcloud(word freq, title=None, max words=200, stopwords=None):

counter = Counter(word freq.fillna(0).to dict())

def count words(df, column='tokens', preprocess=None, min freq=2):

tokens = doc if preprocess is None else preprocess (doc)

background color= "black", colormap="Paired",

max font size=150, max words=max words)

counter = {token:freq for (token, freq) in counter.items()

if token not in stopwords}

freq df = pd.DataFrame.from dict(counter, orient='index', columns=['freq'])

twitter_data_wcloud = twitter_data[twitter_data['artist'] == artist]

twitter wc = wc.generate from frequencies(word freq['freq'])

cher Twitter Cloud

robyn Twitter Cloud

lyrics data wcloud = lyrics data[lyrics data['artist name'] == artist] word_freq = count_words(lyrics_data_wcloud, column='tokens', min_freq=5)

lyrics_wc = wc.generate_from_frequencies(word_freq['freq'])

cher Lyrics Cloud

find find

make neverevery Wa

heart dance see

A: We observe that certain words such as "Love," "Life," "Music," and "Follow" recur in the lyrics of both artists, suggesting that these terms might be prevalent in the music industry. While some words are common between the two, there are also distinctly different words and phrases evident in their music. Similarly, we notice the usage of common words in the Twitter posts of both

lifeSa\

time

plt.imshow(lyrics wc, interpolation='bilinear')

1ve

ill

cause

robyn Lyrics Cloud

go

cause back

Q: What observations do you have about these (relatively straightforward) wordclouds?

е

plt.imshow(twitter wc, interpolation='bilinear')

word freq = count words(twitter data wcloud, column='tokens', min freq=min freq) wc = WordCloud(background color='white', max words=50, width=800, height=400)

For building wordclouds, we'll follow exactly the code of the text. The code in this section can be found here. If you haven't already,

ADS 509 Module 3: Group Comparison

• Calculate descriptive statistics on the two sets of lyrics and compare the results.

For each of the four corpora, find the words that are unique to that corpus.

learned in reading and lecture.

you are asked to do the following:

Build word clouds for all four corpora.

guide, please include a cell with a link.

data and do your cleaning (tokenization and normalization).

General Assignment Instructions

your token set.

cell after reading it.

The task of comparing two groups of text is fundamental to textual analysis. There are innumerable applications: survey

respondents from different segments of customers, speeches by different political parties, words used in Tweets by different constituencies, etc. In this assignment you will build code to effect comparisons between groups of text data, using the ideas

This assignment asks you to analyze the lyrics and Twitter descriptions for the two artists you selected in Module 1. If the results from that pull were not to your liking, you are welcome to use the zipped data from the "Assignment Materials" section. Specifically,

Read in the data, normalize the text, and tokenize it. When you tokenize your Twitter descriptions, keep hashtags and emojis in

Each one of the analyses has a section dedicated to it below. Before beginning the analysis there is a section for you to read in the

These instructions are included in every assignment, to remind you of the coding standards for the class. Feel free to delete this

One sign of mature code is conforming to a style guide. We recommend the Google Python Style Guide. If you use a different style

Your code should be relatively easy-to-read, sensibly commented, and clean. Writing code is a messy process, so please be sure to

Make use of non-code cells for written commentary. These cells should be grammatical and clearly written. In some of these cells

edit your final submission. Remove any cells that are not needed or parts of cells that contain unnecessary code. Remove

inessential import statements and make sure that all such statements are moved into the designated cell.