These instructions are included in every assignment, to remind you of the coding standards for the class. Feel free to delete this cell after reading it. One sign of mature code is conforming to a style guide. We recommend the Google Python Style Guide. If you use a different style guide, please include a cell with a link. Your code should be relatively easy-to-read, sensibly commented, and clean. Writing code is a messy process, so please be sure to 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. Make use of non-code cells for written commentary. These cells should be grammatical and clearly written. In some of these cells 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 import numpy as np from collections import Counter, defaultdict from nltk.corpus import stopwords from string import punctuation sw = stopwords.words("english") # Add any additional import statements you need here import collections # I ended up using the M1 Results for Cher and Robyn because my files for Rush and The Killers would not wor # see now from the feedback on my module 1 HW. # change `data location` to the location of the folder on your machine. data location = "./M1 Results/" # These subfolders should still work if you correctly stored the # data from the Module 1 assignment twitter\_folder = "twitter/" lyrics folder = "lyrics/" # Section Complete def descriptive stats(tokens, num tokens = 5, verbose=True) : In [4]: 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 unique tokens, lexical diversity, and number of characters. # Fill in the correct values here. I used the len function to get these correct values. 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. I think when I ran this the first time a few of these funcitons return([num\_tokens, num\_unique\_tokens, num\_characters, lexical diversity # Section Complete text = """here is some example text with other example text here in this text""".split() assert(descriptive\_stats(text, verbose=True)[0] == 13) assert(descriptive\_stats(text, verbose=False)[1] == 9) assert(descriptive\_stats(text, verbose=False)[2] == 55) assert(abs(descriptive\_stats(text, verbose=False)[3] - 0.69) < 0.02)</pre> There are 13 tokens in the data. There are 9 unique tokens in the data. There are 55 characters in the data. The lexical diversity is 0.692 in the data. Q: Why is it beneficial to use assertion statements in your code? A: To perform a sanity check during the development. Data Input Now read in each of the corpora. For the lyrics data, it may be convenient to store the entire contents of the file to make it easier to inspect the titles individually, as you'll do in the last part of the assignment. In the solution, I stored the lyrics data in a dictionary with two dimensions of keys: artist and song. The value was the file contents. A data frame would work equally well. For the Twitter data, we only need the description field for this assignment. Feel free all the descriptions read it into a data structure. In the solution, I stored the descriptions as a dictionary of lists, with the key being the artist. In [6]: # Read in the lyrics data root\_path = "./M1 Results/lyrics" lyric data = dict() for dir name in os.listdir(root path): if dir name != ".DS Store": file\_path = root\_path+'/'+dir\_name lyric\_data[dir\_name] = dict() for file name in os.listdir(file path): full\_file\_path = file\_path+'/'+file\_name with open(full\_file\_path, 'r') as f: text = f.read() title = text.split("\n")[0].strip('"') artist = dir name lyric data[artist][title] = text print(f"Read lyrics: {len(lyric data[artist].keys())} songs of {dir name}") # Section Complete Read lyrics: 313 songs of cher Read lyrics: 93 songs of robyn In [7]: # To make it easier to use further down I read this into a pandas df. lyric df = (pd.DataFrame(lyric data) .reset index() .rename(columns = {'index':'title'})) lyric df.loc[lyric df['robyn'].isnull() ==False, 'artist'] = 'robyn' lyric df.loc[lyric df['cher'].isnull() == False, 'artist'] = 'cher' lyric df.loc[lyric df['robyn'].isnull() == False, 'lyric'] = lyric df['robyn'] lyric df.loc[lyric df['cher'].isnull() == False, 'lyric'] = lyric df['cher'] lyric df.drop(['cher','robyn'],axis = 1, inplace=True) lyric df.head() # Section Complete title artist lyric 88 Degrees cher "88 Degrees"\n\n\n\Stuck in L.A., ain't got n... 1 A Different Kind Of Love Song cher "A Different Kind Of Love Song"\n\n\n\NWhat if... 2 After All "After All"\n\n\nWell, here we are again\nl ... cher Again cher "Again"\n\n\nAgain evening finds me at your ... "Alfie"\n\n\nWhat's it all about, Alfie?\nls... 4 Alfie cher # Read in the twitter data root path = "./M1 Results/twitter" twitter data = list() for file name in os.listdir(root path): if 'data' in file name: file path = root path+'/'+file name columns = ['screen name', 'name', 'id', 'location', 'followers count', 'friends count', 'description'] df = pd.read\_csv(file\_path,sep='\\t',engine='python') df['following'] = file name.split(" followers")[0] twitter data.append(df[['following','description']]) print(f"Read file:{file path}",f"{df.shape[0]} rows") twitter df = pd.concat(twitter data,axis=0).fillna(" ") print(f"Create dataframe: {twitter df.shape[0]} rows in total.") # Section Complete Read file:./M1 Results/twitter/cher followers data.txt 3994803 rows Read file:./M1 Results/twitter/robynkonichiwa followers data.txt 358372 rows Create dataframe: 4353175 rows in total. In [9]: # To validate that I created the twitter df with following as the index of which artist the follower is foll twitter\_df.head() # Section Complete following description 0 cher cher ��ð�š~ð�šžð�š� ð�šœð�šžð�š™ð�š™ð�š°ð�š›ð... 2 cher 163��æ"›ã�∢ã�£ã�·ðŸ′œ26æ³ðŸ�′ 工〇好ã�... cher csu cher Writer @Washinformer @SpelmanCollege alumna #D... **Data Cleaning** Now clean and tokenize your data. Remove punctuation chacters (available in the punctuation object in the string library), split on whitespace, fold to lowercase, and remove stopwords. Store your cleaned data, which must be accessible as an interable for descriptive\_stats , in new objects or in new columns in your data frame. punctuation = set(punctuation) # speeds up comparison print(punctuation) {'\*', ':', '|', ';', '^', '=', '<', '\$', '.', "'", '/', '+', '[', '#', '`', '!', '\\', ',', '>', '?', '@', 18', '(', '&', '{', '}', ')', '\_', '"', '~', ']', '-'} # Create funcitons to help with the cleaning as asked above: def remove\_stop(tokens): stopwords = set(nltk.corpus.stopwords.words('english')) return [t for t in tokens if t.lower() not in stopwords] def tokenize(text): tokenizer = nltk.RegexpTokenizer(r"\w+") tokens = tokenizer.tokenize(text) return [t for t in tokens if t.isalpha()] def remove\_punctuation(tokens): return [t for t in tokens if t not in set(punctuation)] def remove\_letterlike\_symbol(tokens): e.g. Japenese, Korean... return [t for t in tokens if t.isascii()] pipeline = [str.lower, tokenize, remove\_punctuation, remove\_stop,remove\_letterlike\_symbol] def prepare(text, pipeline): tokens = text for transform in pipeline: tokens = transform(tokens) return " ".join(tokens) # Section Complete import nltk In [14]: # create your clean twitter data here twitter\_df['description\_token']=twitter\_df['description'].apply(prepare,pipeline=pipeline) twitter\_df['description\_token'].head() # Section Complete Out[14]: 0 writer washinformer spelmancollege alumna dcna... Name: description token, dtype: object In [15]: # create your clean lyrics data here lyric\_df['lyric\_token']=lyric\_df['lyric'].apply(prepare,pipeline=pipeline) lyric\_df['lyric\_token'].head() Out[15]: 0 degrees stuck 1 got friends hollywood nuts man... different kind love song world crazy sane woul... well guess must fate tried deep inside known b... evening finds door ask could try know quite sa... alfie alfie moment live sort alfie meant take ... Name: lyric token, dtype: object In [16]: | lyric\_df['title\_token']=lyric\_df['title'].apply(prepare,pipeline=pipeline) lyric\_df['title\_token'].head() # Section Complete degrees different kind love song 2 3 alfie Name: title\_token, dtype: object **Basic Descriptive Statistics** Call your descriptive\_stats function on both your lyrics data and your twitter data and for both artists (four total calls). In [17]: # calls to descriptive\_stats here # I'm starting by breaking this up first for lyrics of Robyn: text = " ".join(lyric df[lyric df['artist'] == "robyn"]['lyric token'].tolist()).split() print("Lyrics by Artist: robyn") descriptive\_stats(text, verbose=True) Lyrics by Artist: robyn There are 11877 tokens in the data. There are 2053 unique tokens in the data. There are 58057 characters in the data. The lexical diversity is 0.173 in the data. Out[17]: [11877, 2053, 58057, 0.17285509808874294] In [18]: # Now for cher: text = " ".join(lyric df[lyric df['artist']=="cher"]['lyric token'].tolist()).split() print("Lyrics by Artist: cher") descriptive stats(text, verbose=True) Lyrics by Artist: cher There are 32601 tokens in the data. There are 3527 unique tokens in the data. There are 158050 characters in the data. The lexical diversity is 0.108 in the data. Out[18]: [32601, 3527, 158050, 0.10818686543357565] In [19]: # Now to break down the descriptive stats for followers of Robyn: text = " ".join(twitter df[twitter df['following'] == "robynkonichiwa"]['description token'].tolist()).split() print("Twitter Followers of robynkonichiwa") descriptive stats(text, verbose=True) Twitter Followers of robynkonichiwa There are 1402208 tokens in the data. There are 155185 unique tokens in the data. There are 8128442 characters in the data. The lexical diversity is 0.111 in the data. Out[19]: [1402208, 155185, 8128442, 0.11067188320135102] In [20]: # Descriptive stats for followers of cher: text = " ".join(twitter df[twitter df['following']=="cher"]['description token'].tolist()).split() print("Twitter Followers of cher") descriptive stats(text, verbose=True) Twitter Followers of cher There are 14331010 tokens in the data. There are 617321 unique tokens in the data. There are 82829471 characters in the data. The lexical diversity is 0.043 in the data. Out[20]: [14331010, 617321, 82829471, 0.043075889277866666] Q: How do you think the "top 5 words" would be different if we left stopwords in the data? A: I think the top 5 words would change to include the stop words. I think this would not be useful to include that is why we drop them. Q: What were your prior beliefs about the lexical diversity between the artists? Does the difference (or lack thereof) in lexical diversity between the artists conform to your prior beliefs? A: I don't know too much about these artists so my initial thought would be that they would have similar lexical diversity. I was wrong because we can see that Roby has a higher lexical diversity compared to cher. **Specialty Statistics** The descriptive statistics we have calculated are quite generic. You will now calculate a handful of statistics tailored to these data. 1. Ten most common emojis by artist in the twitter descriptions. 2. Ten most common hashtags by artist in the twitter descriptions. 3. Five most common words in song titles by artist. 4. For each artist, a histogram of song lengths (in terms of number of tokens) We can use the emoji library to help us identify emojis and you have been given a function to help you. def is emoji(s): return(s in emoji.UNICODE EMOJI['en']) assert(is\_emoji("♥♥")) assert(not is\_emoji(":-)")) Emojis 🖨 What are the ten most common emojis by artist in the twitter descriptions? # Your code here, first for robyn text = " ".join(twitter\_df[twitter\_df['following'] == "robynkonichiwa"]['description'].tolist()).split() emojis = [x for x in text if is emoji(x)] counter=collections.Counter(emojis) print("Ten most common emojis by robynkonichiwa") print(counter.most\_common(10)) Ten most common emojis by robynkonichiwa #Emojis for cher: text = " ".join(twitter\_df[twitter\_df['following']=="cher"]['description'].tolist()).split() emojis = [x for x in text if is\_emoji(x)] counter=collections.Counter(emojis) print("Ten most common emojis by cher") print(counter.most\_common(10)) # Section Complete Ten most common emojis by cher [('TM', 1)] Hashtags What are the ten most common hashtags by artist in the twitter descriptions? # Your code here, first for Robyn: In [24]: text = " ".join(twitter\_df[twitter\_df['following'] == "robynkonichiwa"]['description'].tolist()).split() hashtags = [x for x in text if x[0] == '#']counter=collections.Counter(hashtags) print("Ten most common hashtags by robynkonichiwa") print(counter.most\_common(10)) Ten most common hashtags by robynkonichiwa [('#BlackLivesMatter', 311), ('#BLM', 273), ('#blacklivesmatter', 200), ('#1', 187), ('#music', 150), ('#', 150), ('#Music', 93), ('#EDM', 79), ('#blm', 51), ('#TeamFollowBack', 51)] # Now hashtags for cher: text = " ".join(twitter\_df[twitter\_df['following']=="cher"]['description'].tolist()).split() hashtags = [x for x in text if x[0] == '#']counter=collections.Counter(hashtags) print("Ten most common hashtags by cher") print(counter.most\_common(10)) # Section Complete Ten most common hashtags by cher [('#BLM', 7920), ('#Resist', 5007), ('#BlackLivesMatter', 4226), ('#resist', 3134), ('#FBR', 2768), ('#black livesmatter', 2464), ('#TheResistance', 2459), ('#1', 2226), ('#', 1965), ('#Resistance', 1518)] Song Titles What are the five most common words in song titles by artist? The song titles should be on the first line of the lyrics pages, so if you have kept the raw file contents around, you will not need to re-read the data. # Your code here # After removing stopwords for robyn: text = " ".join(lyric\_df[lyric\_df['artist'] == "robyn"]['title\_token'].tolist()).split() counter=collections.Counter(text) print("Five most common words by robyn") print(counter.most\_common(10)) Five most common words by robyn [('love', 5), ('thing', 3), ('girl', 3), ('u', 3), ('like', 2), ('baby', 2), ('music', 2), ('woman', 2), ('g irlfriend', 2), ('get', 2)] In [27]: # After removing stopwords for cher: text = " ".join(lyric\_df[lyric\_df['artist'] == "cher"]['title\_token'].tolist()).split() counter=collections.Counter(text) print("Five most common words by cher") print(counter.most\_common(10)) # Section Complete Five most common words by cher [('love', 36), ('man', 15), ('song', 11), ('come', 7), ('one', 7), ('believe', 6), ('heart', 6), ('time', 6), ('go', 6), ('know', 6)] Song Lengths For each artist, a histogram of song lengths (in terms of number of tokens). If you put the song lengths in a data frame with an artist column, matplotlib will make the plotting quite easy. An example is given to help you out. num replicates = 1000 df = pd.DataFrame({ "artist" : ['Artist 1'] \* num replicates + ['Artist 2']\*num replicates, "length": np.concatenate((np.random.poisson(125,num replicates)),np.random.poisson(150,num replicates))) }) df.groupby('artist')['length'].plot(kind="hist",density=True,alpha=0.5,legend=True) Out[28]: artist AxesSubplot (0.125, 0.125; 0.775x0.755) Artist 1 AxesSubplot(0.125,0.125;0.775x0.755) Name: length, dtype: object 0.040 Artist 1 Artist 2 0.035 0.030 0.025 0.020 0.015 0.015 0.010 0.005 0.000 120 140 160 180 200 Since the lyrics may be stored with carriage returns or tabs, it may be useful to have a function that can collapse whitespace, using regular expressions, and be used for splitting. Q: What does the regular expression '\s+' match on? A: It will match 1 or more repitions of whitespace with the above regular expression. collapse\_whitespace = re.compile(r'\s+') def tokenize\_lyrics(lyric) : """strip and split on whitespace""" return([item.lower() for item in collapse\_whitespace.split(lyric)]) # Your lyric length comparison chart here. lyric df['lyric split length'] = lyric df['lyric'].apply(lambda x: len(tokenize lyrics(x))) lyric df.head() df = pd.DataFrame({ "artist" : lyric df['artist'], "length" : lyric df['lyric split length']

df.groupby('artist')['length'].plot(kind="hist",density=True,alpha=0.5,legend=True)

cher

robyn

600

AxesSubplot(0.125,0.125;0.775x0.755) AxesSubplot(0.125,0.125;0.775x0.755)

Name: length, dtype: object

100

200

300

400

Out[30]: artist

0.006

0.005

0.004

0.003

0.002

0.001

0.000

ADS 509 Assignment 2.1: Tokenization, Normalization,

This notebook holds Assignment 2.1 for Module 2 in ADS 509, Applied Text Mining. Work through this notebook, writing code and

In the previous assignment you put together Twitter data and lyrics data on two artists. In this assignment we explore some of the textual features of those data sets. If, for some reason, you did not complete that previous assignment, data to use for this assignment

This assignment asks you to write a short function to calculate some descriptive statistics on a piece of text. Then you are asked to find

**Descriptive Statistics** 

can be found in the assignment materials section of Blackboard.

some interesting and unique statistics on your corpora.

**General Assignment Instructions** 

answering questions where required.