

Stephen Hullender  
CIS 3715 001  
April 24, 2022

## **Predicting and Identifying Genre of Song Based on Lyrical Content: A Final Report**

### **Introduction**

Music is integral in our daily lives, and the composure of a song's lyrical content is an important attribute to consider when listening to music. Lyrics can influence our choice of music, and impact the song's significant meaning, as well as influence the reputability of a genre. When looking at the words of a song, it can be interesting to consider these factors and figure out what words define different types of mainstream music. What distinguishes rock music from pop, or rock from metal? How about gospel music to romance?

This project seeks to understand how certain vocabulary influences genres, such as measuring what words are common among different genres and/or artists. The goal of the project is to accumulate the lyrical content of a variety of popular songs and predict the genre of any given sample of lyrics provided. The worksheet will also look at the most common words for all genres provided, assuming that the genres meet the criteria of having a feasible number of songs. For this project, I was the sole contributor of the project, responsible for all implementation of libraries and Python code in these worksheets, and am the sole author of the report provided.

### **Approach**

The main approaches that will be followed include using K-nearest neighbors to determine how many nearest neighbors are possible, logistic regression to calculate the accuracy and F-1 scores, and K-means classifier to distribute the top ten most common words across clusters representative of genres. The purpose of using KNN is to read accuracy of the data provided and to read the best parameters (i.e. nearest neighbors). Logistic regression was initially chosen as the project was intended on researching more into other variables such as artist and popularity score, two factors that may or may not have influenced how lyrics would be predicted. Lastly, KMeans was used to arrange a cluster of words that most resemble the given features, or genres.

To begin the project, we are introduced to two datasets: the artists dataset, and the lyrics dataset. These datasets were downloaded from Kaggle as 'Song lyrics from 79 musical genres'<sup>[1]</sup>. The artists dataset consists of the following columns; "Artist", "Genres", "Songs", "Popularity", and "Link". The lyrics dataset contains columns: "Artist Link", "Song Name", "Song Link", "Lyric", and "language". The common identifier is the artist link, "Link" and "ALink" for artist data and lyrical data, respectively.

To reduce the amount of clutter from the given information, we use Exploratory Data Analysis to clean out unwanted data. First, we delete unnecessary columns such as artist name, number of songs, and popularity score of the artist. The lyrical column is deleted because the lyrics provided are only a sample and don't represent the song's entirety, but will be replaced with a full string of lyrics eventually. Second, we halve our data by eliminating non-English entries, indicated by the "language" column in the lyrics

data; if an artist has no more songs as a result, they are eliminated from the artist data. Next, we iterate through each genre item in the “Genre” column, and eliminate any entries with certain genres. Some genres are removed for the following reasons:

- not enough entries
- redundant genres that can be perfectly matched with other existing, and bigger, genres
- a reasonable assumption that the majority of songs within these genres are non-English (even if non-English entries are already cut out of the dataset)

Afterwards, we cleared out the Genre column and replaced it with two columns for the first and second (if applicable) genre. The reason why I chose two choices for each song is to analyze the results for the genres provided (i.e. “is the genre provided accurate? How about the second genre choice?”). By this point, we have cut down from almost 380,000 songs to around 184,300 songs. In all that we have processed, the end result will only contain about 10,000 songs; this is mainly due to a conflict between time and the large amounts of data being processed, as well as filtering out bad entries from the datasets.

After we are done getting the songs we want, we start to get the lyrical data we need. All web scraping and data filtering is to be done through the Python Beautiful Soup library, which pulls up the HTML code for the web page containing a song’s lyric; the link to this page is provided through the “Song Link” column of the lyrics dataset. The HTTP requests are also performed by another Python library called urllib. Another major library that will be used to clean up words is to use the Natural Language Toolkit (nltk) to remove stopwords. Stopwords are words that take up space, such as “I”, “you”, “and”, “a(n)”, etc.

One other approach in relation to cleaning text is using regular expressions to clear out most punctuation symbols, with the exception of the single quote, for words like “i’ll” or “we’ve”. With that, there is an important issue that should be brought to attention when handling the datasets: all of our lyrical content is to be extracted using web scraping tools. We will be pulling data from a lyrics website called Vagalume.com (as the author of the dataset got their data from this source). It is possible that the format of the lyrics received will vary, and will require some discretion even after cleaning up unneeded characters. This also means that it is possible that the lyrical data may not accurately predict the genres due to garbage text that can only be caught by chance.

Due to time constraints, and lack of knowledge of parallel processing, despite attempts to connect relevant libraries to expedite the process, I was left with little choice but to iterate through each entry to get the lyrical content of the song, as well as summatting the number of instances of each word based on genre (this is to compare what words are used most in a given genre). Because of my naive proposition, the maximum that I allowed is around half, or 95,000 songs, to be processed into our next step. To save the data and save time, I came up with the idea of printing it to a .csv file and running the project on two worksheets. An aforementioned, and last minute, decision was to pick a randomized 10,000 songs instead.

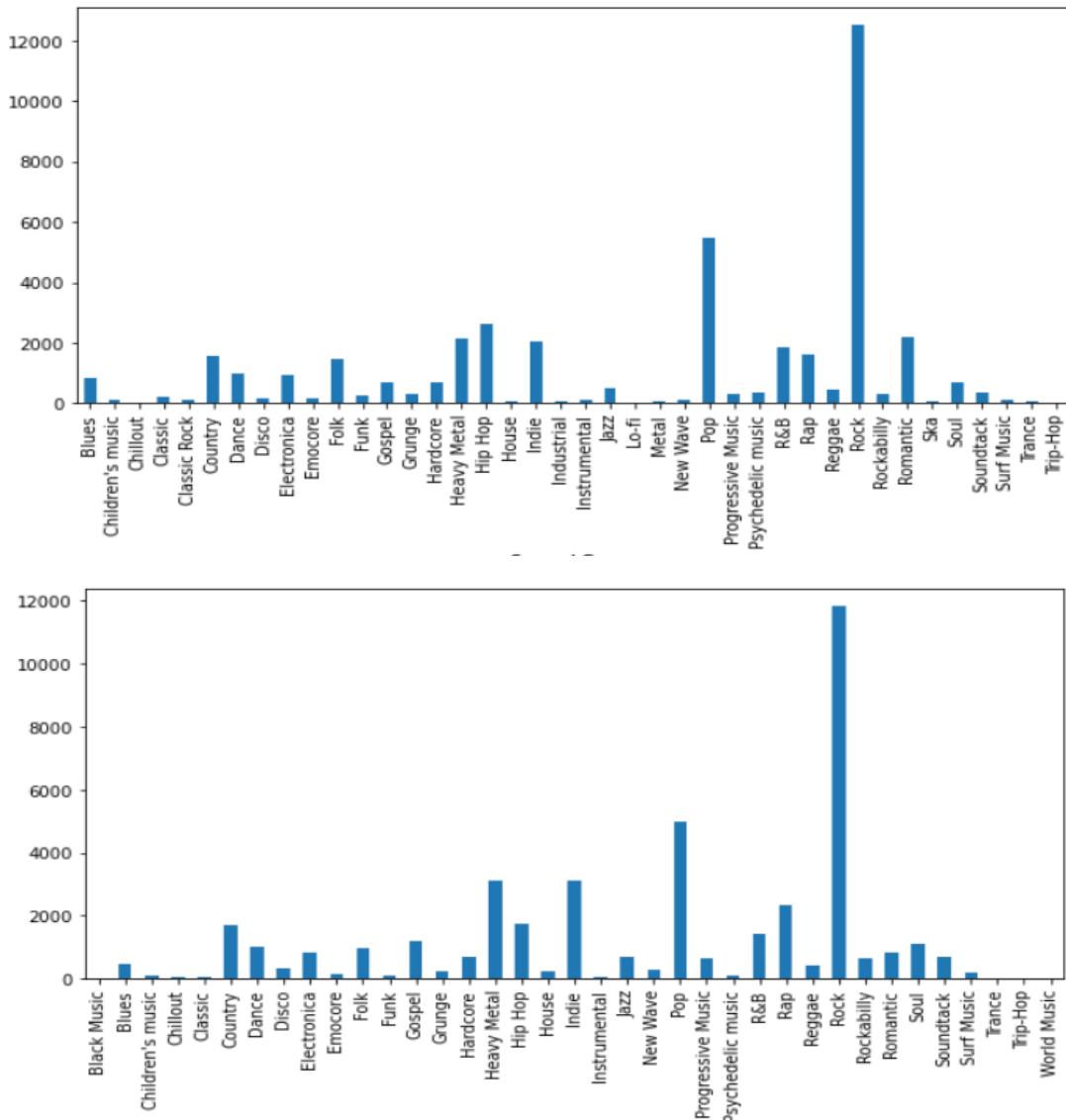
Since the project involves extracting text, one approach is to use the Bag-of-Words model, which includes the use of TF-IDF, or term frequency-inverse document frequency. This model measures how relevant a word is on the condition of its presence among lyrics from several songs. The reason why we also use logistic regression to begin making our predictions is because logistic regression is easy to

implement and fits with the data that we want to classify, such that our data contains 0's and 1's for every word entry provided per document (binary classification), as well as data involving the frequency of each word using the aforementioned Bag-of-words approach. In the case of clustering words, we decided that K-Means would be a great fit for arranging the words for each genre.

## **Results**

Based on the data given for all words captured for all genres, we have the following visualizations, giving us the top 10 words from each major genre:

**First genre choice results (top), second genre choice results (bottom)**



The results indicate an obvious skew, which is that the majority of songs belong in a select few genres, noticeably 'Rock', 'Pop', along with secondary genres such as 'Heavy Metal' and 'Indie'.

The reason why I chose to look through two sets of genres is based on the multitude of genres that the majority of the songs given provided. Because processing the data with more than one genre at once would probably complicate results, it was easier to iterate through both TF-IDF vectorization and logistic regression using each set of genres each time. Furthermore, two columns were used as an assumed factor that the majority of songs would average out to about 2 genres per artist (i.e. the distribution of genres per artist would be equal among artists that have 1, 2, or 3 genres, but no more than that).

Initially when performing TF-IDF and logistic regression, the implementation of the models were not clear, mostly from lack of knowledge on how to use certain sklearn libraries. Because of this, the results given for the first implementation show very different, and bad, results compared to the final results. For the next few paragraphs, I will be discussing the results made based on these implementations:

### Rough comparison of genre predictions vs. actual genre names

	Real	Predictions	Correct?		Real	Predictions	Correct?
<b>0</b>	Hip Hop	Pop	False	<b>0</b>	R&B	R&B	True
<b>1</b>	Rock	Romantic	False	<b>1</b>	Rock	Romantic	False
<b>2</b>	House	Pop	False	<b>2</b>	Electronica	Hip Hop	False
<b>3</b>	Indie	Romantic	False	<b>3</b>	Surf Music	Romantic	False
<b>4</b>	Gothic Rock	Folk	False	<b>4</b>	Gothic Rock	Country	False
...	...	...	...	...	...	...	...
<b>12543</b>	Heavy Metal	Heavy Metal	True	<b>12543</b>	Hard Rock	Progressive Music	False
<b>12544</b>	R&B	Country	False	<b>12544</b>	Hip Hop	Hip Hop	True
<b>12545</b>	Pop	Pop	True	<b>12545</b>	Pop	R&B	False
<b>12546</b>	Rock	Progressive Music	False	<b>12546</b>	Romantic	Progressive Music	False
<b>12547</b>	Soul	Romantic	False	<b>12547</b>	R&B	Romantic	False
-----	-	-	-	-----	-	-	-

The two results were an attempt at comparing predicted genres to actual song genres based on first and second genre choice, respectively. The results show that the majority of predictions point to an incorrect result. Remarkably, despite only making up only a small percentage of the total songs, “Rap” appears in a significant number of predictions, even if the genres do not match whatsoever.

### Comparison of real vs. predicted genres for first and second genre choice

	Real	Predictions	Correct?		Real	Predictions	Correct?
<b>count</b>	12548	12548	12548	<b>count</b>	12548	12548	12548
<b>unique</b>	42	34	2	<b>unique</b>	47	34	2
<b>top</b>	Rock	Rap	False	<b>top</b>	Rock	Rap	False
<b>freq</b>	2126	4567	10115	<b>freq</b>	2266	4250	10499

### Accuracy, F-1 (micro), recall, precision scores for first and second genre choice

```
accuracy is 0.13232023721275019,  
f1 score is 0.13232023721275019,  
recall is 0.13232023721275019,  
precision is 0.13232023721275019
```

```
accuracy is 0.09673832468495182,  
f1 score is 0.09673832468495182,  
recall is 0.09673832468495182,  
precision is 0.09673832468495182
```

One observation, from using various sample sizes for TF-IDF and logistic regression, is that the bigger the size of the processed dataset, the better the results are. Between the accuracy, recall, precision and F1 scores, all shared about 0.1 when using around 90 songs. This increased to around 0.13 when the final processing returned around 8000 songs. When looking at the results from a bird's eye view, the overall scores are lacking, possibly due to incorrect implementation of logistic regression and/or TF-IDF. The formulas for TF-IDF were created manually, before implementing the TF-IDF vectorizer from sklearn, and therefore took longer than expected.

One of the reasons why the logistic regression most likely did not work as expected is because the values that we want to compare are text-based. Unlike other lab assignments where regression models were practiced using binary values for data (0's or 1's), or even numerical data, we were dealing with text. Consequently, this worksheet was scrapped and a new implementation of TF-IDF and Logistic Regression was created.

After revamping the models, I started by adding the TF-IDF vectorizer to the solution; the vectorizer converts a list of text into assorted numerical values accordingly by using the TF-IDF equation. The Logistic Regression model differs in that the original implementation excluded GridSearchCV, and attempted to calculate all metrics, including precision and recall, which would always return the same, low, result as the accuracy and F-1 score. GridSearchCV, part of the model selection library from sklearn, is used to parametrize data using a cross-validation method; this was used in both the K-neighbors classifier and logistic regression models to consider different combinations of data when passing it through their respective models.

#### Best scores from GridSearchCV, using KNearestClassifier, for first & second genre choice

```
0.18073170731707316 0.175  
{'n_neighbors': 4} {'n_neighbors': 4}  
Wall time: 10.8 s Wall time: 11.3 s
```

#### Accuracy and F-1 (macro & micro) scores, logistic regression, first & second genre choice

```
{'C': 3}  
0.35555555555555555557 0.11614313256404439 0.35555555555555555557  
Wall time: 39min 44s  
  
{'C': 2}  
0.3522222222222222 0.0679471926576408 0.3522222222222222  
Wall time: 50min 6s
```

The accuracy scores show that there are some possible hindrances with getting the best data with the lyrics provided. Possible errors may be the result of bad web scraping or words that have made it past filtering for stopwords and/or punctuation. Overall, the accuracy and F-1 scores improved much more compared to the first implementation of the TF-IDF and logistic regression models. As predicted, the results also suggest that the micro F-1 score works better than the macro F-1 score; this was originally chosen under the condition that there were multiple variables that would determine the outcome of the predictions made (i.e. measuring more than just binary data).

Finally, we were able to create clusters to fit in the most common words per genre. All clusters were generated using the K-means classifier, and retrieving the cluster centers provided. For every cluster we obtained, we iterated through each center and got the most common words using argsort(), which sorts out the most frequent values. Here are the results provided:

#### **Results for each cluster, first genre choices**

Cluster 0	ve	Cluster 4	^	Cluster 8	roll	Cluster 12	make
away	life		yeah	oh	girl	Cluster 16	gonna
far	eyes		know	know	know		know
don	world		got	like	like		don
fly	heart		love	got	got		ve
day	love		day	don	don		way
just	day		lost	baby	love		love
ll	inside		know	ll	man		live
run	know		pain	ve	just		just
sail	like		like	like	baby		ain
love	world		world	world	world		like
Cluster 1	little	Cluster 5	home	ve	Cluster 13	nigga	Cluster 17
da	coming		got	shit	Cluster 17	ll	
bit	come		know	niggas		love	
just	just		just	like		just	
love	gone		love	fuck		know	
oh	know		don	ain		come	
baby	ve		life	bitch		hold	
tell	don		ll	niggaz		make	
don	like		way	got		day	
like	ll		time	em		way	
Cluster 2	time	Cluster 6	christmas	Cluster 10	Cluster 14	right	Cluster 18
say	don		merry	thing	Cluster 18	lyrics	monroe
goodbye	know		year	just		leitch	donovan
way	just		love	know		mercyme	sentenced
tell	say		santa	time		baez	
just	need		time	cause		joan	
love	love		tree	ll		ashanti	
don	like		blue	doesn		boney	
day	cause		mistletoe	don			
long	tell		belts	heart			
Cluster 3	love	Cluster 7	oh	Cluster 11	Cluster 15	believe	Cluster 19
need	yeah		got	ll	Cluster 15	night	night
know	love		ain	11		holy	holy
heart	don		like	just		sleep	sleep
just	baby		rock	love		day	day
don	know		man	don		heavenly	heavenly
time	just		just	know		light	light
11	make		don	dreams		peace	peace
say	come		know	miracles		long	long
ve	ll		uh	heart		silent	silent
Cluster 4	sight	Cluster 8	feel	Cluster 12	Cluster 16	make	sight
Cluster 20	sight	Cluster 24	let	Cluster 28			
want	let		na				
don	love		oh				
know	don		know				
just	11		ve				
love	heart		better				
need	just		ba				
11	come		make				
like	know		11				
say	oh		got				
come	time		look				
Cluster 21	come	Cluster 25	Cluster 29				
la	like		let				
love	feel		na				
oh	just		oh				
don	know		know				
know	don		ve				
like	love		better				
baby	cause		ba				
come	think		make				
fa	going		11				
somebody	make		got				
Cluster 22	somebody	Cluster 26	Cluster 30				
baby	just		ya				
love	don		shake				
don	like		like				
know	know		wanna				
just	oh		wanna				
oh	man		got				
got	11		know				
want	love		love				
come	ain		ain	baby			
11	way		don	hey			
Cluster 23	way	Cluster 27	Cluster 31	Cluster 32			
ah	tonight		wanna	like			
oh	alright		ooh	oh			
ha	love		don	love			
know	11		know	don			
love	gonna		love	got			
want	right		just	just			
got	cause		like	way			
yeah	make		oh	make			
let	know		want	know			
feel	let		baby				
Cluster 24	let						

Wall time: 1min 25s

## Results for each cluster, second genre choices

Cluster 0	Cluster 4	Cluster 8	Cluster 12	Cluster 16
ll love just know away need alright come way make	day away just long love night today 11 come don	girl know like got oh just love 11 don baby 11	lost carry wanted ve life know love oh just long	god come lord sun won free love life light like
Cluster 1	Cluster 5	Cluster 9	Cluster 13	Cluster 17
hey don like love got oh know just make girl	sing sweet song thing love dream oh just 11 come	la love da oh come don baby fa tomorrow	home gone coming ve away long just know love	time say goodbye just bye don remember love need way
Cluster 2	Cluster 6	Cluster 10	Cluster 14	Cluster 18
yeah oh got know love don baby 11 like just	talk don step like say know just need way love	na oh know ve ba 11 grind better make just	dance let music wanna just baby like night come ya	believe just don love 11 make know dreams ve life
Cluster 3	Cluster 7	Cluster 11	Cluster 15	Cluster 19
like man just don know love 11 got oh little	ooh oh yeah love like don know feeling just baby	oh love yeah don baby know just got like make	life blood death away die pain black dead fight hell	don know say love just 11 like away ve need
Cluster 4	Cluster 20	Cluster 24	Cluster 28	Cluster 20
	wanna don know love just like baby feel want oh	heart break love don know just oh apart broken let just	ve got know don just love life like 11 time	heart break love don just life like 11 time
Cluster 21	Cluster 25	Cluster 29		
shake boom baby like ya gonna let got hot mellow	don just like chorus know think tell got cause way	love need just know don heart time 11 say like		
Cluster 22	Cluster 26	Cluster 30		
gonna know ain don make love live 11 just tonight	nigga shit like niggas got ain fuck bitch ya em	baby love know don got just oh come yeah 11		
Cluster 23	Cluster 27	Cluster 31		praise
christmas merry year santa love tree time 11 mistletoe snow	let love don just come know won away time	jesus holy heavenly night lord christ silent peace sleep praise		Cluster 32
Cluster 24				feel night like just right kiss love 11 know tonight

Wall time: 1min 25s

The results are not the best, and despite efforts to contain stopwords, as well as onomatopoeic sounds like 'oh' and 'lalala', some words still show up. Due to this being a last-minute approach, I was

unable to identify which genre is which based on the clusters given, although a few were identified using context clues with the words given. For example, it is very likely that Cluster 10 for the first set of genres, and Cluster 23 for the second set, contain holiday or Christmas music; this is based on words such as “christmas”, “merry”, “santa”, etc. Likewise, a genre like “Gospel” is identified from Cluster 19 (first genre choice) and Cluster 31 (second cluster choice) with words such as “jesus”, “holy”, “heavenly”, etc.

Finally, for reference, when extracting lyrical data, I also made the step of storing all words that ever existed (within the constraints of the lyrics dataset) into one, and very heavy, matrix array. Here is a sample of the data provided:

### All word count of words to ever exist from all genres and all songs

1	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD		
2	know	36604	Rock	Indie	Pop	Romantic	Folk	Grunge	Jazz	Classic	New Wa	Rap	47	4340	2779	1928	679	1812	9961	3162	6367	2014	887	566	772	492	420	555	3192	800		
3	love	34587		25871	7260	2760	803	1573	14783	47	539	968	9	3143	4072	7885	6252	2278	1236	1438	595	397	871	438	959	3884	875					
4	like	29495		25256	5073	2316	563	524	392	458	8594	2045	61	5721	1382	1648	1491	503	689	472	228	593	346	515	379	496	270	2333	470			
5	oh	26644	4883	22753	6018	535	560	205	474	6852	14	2448	207	794	618	1666	7004	3246	6659	584	650	505	379	346	379	496	270	2333	470			
6	got	23047	15393	3685	1620	551	591	102	299	13387	14732	29	2441	1912	1137	495	1016	4302	1850	2767	386	369	403	384	362	266	321	1412	812			
7	time	23001	3143	12420	3719	1759	366	533	191	460	4547	5517	58	4077	1881	1398	4034	1504	434	727	3919	1501	355	398	212	176	278	1291	578			
8	never	22618	3104	12424	3874	1575	466	565	248	355	4516	5379	50	4034	1599	1504	434	727	3919	1501	355	398	212	176	278	1291	578					
9	one	22317	3104	13160	3888	1717	415	627	194	364	5631	6848	70	4060	1901	1286	463	934	4731	1833	3224	1258	502	402	313	299	251	404	1676	738		
10	go	21633	3436	14655	3935	1835	292	6321	8385	33	2613	1708	1204	507	1041	2511	1804	4123	999	611	394	443	314	253	405	2293	651					
11	need	21544	3243	1622	3501	1351	442	997	106	299	6321	9448	30	2895	1523	1207	477	1229	762	4549	714	818	244	252	464	2699	598					
12	get	20916	2343	1622	3501	1351	442	997	106	299	6321	9448	30	2895	1523	1207	477	1229	762	4549	714	818	244	252	464	2699	598					
13	can't	19460	2122	3485	1622	5001	478	102	305	391	6098	46	3116	1473	1333	379	940	4637	1670	3068	479	369	317	246	207	308	1608	465				
14	want	18022	2776	1904	3031	1084	335	352	80	350	4325	5924	29	2009	906	1011	282	739	4157	1471	3448	777	556	210	341	197	357	1593	374			
15	come	17345	2649	10623	2958	1481	376	537	225	282	4364	5200	55	2817	1418	690	480	988	3801	1532	2381	1134	583	216	237	319	252	349	1419	510		
16	take	17235	2254	9919	3163	1583	459	151	45	181	318	4263	5357	36	2987	1197	1147	488	1785	2425	891	332	249	243	233	148	224	1287	533			
17	yeah	17115	2415	14358	3163	1583	459	151	45	181	318	4263	5357	36	2987	1197	1147	488	1785	2425	891	332	249	243	233	148	224	1287	533			
18	sway	16708	2123	10041	2471	1581	355	409	187	246	3358	4238	25	2547	1207	928	1763	1207	1763	2425	891	332	249	243	233	148	224	1287	533			
19	say	16347	2289	1780	2254	1245	377	361	127	253	1203	1895	37	303	1111	996	523	1505	1893	1549	883	754	182	239	160	145	151	937	554			
20	baby	15778	2386	11303	3307	1195	319	511	188	281	4625	5883	29	1721	1235	850	352	1141	1611	4612	2032	2597	732	574	260	238	140	264	1251	461		
21	ill	15616	2435	9558	2986	1278	336	587	169	249	2700	3349	26	2199	1459	960	539	716	2959	1209	2303	307	679	268	285	318	256	318	1194	378		
22	bad	15456	2295	8577	2356	1301	370	254	83	184	6155	6954	53	2209	1531	926	429	719	3161	948	1658	621	218	278	234	171	137	174	1255	366		
23	gonna	15156	1972	10039	2527	1041	224	381	107	195	257	3227	15	1821	1417	606	513	159	3013	1369	2338	549	409	136	209	259	286	175	996	274		
24	let	15020	2124	11776	3187	1109	292	397	165	224	4559	5906	57	2340	1229	916	312	544	4466	1638	3103	1256	321	323	275	238	175	342	1682	393		
25	life	14651	2144	1444	2471	1581	355	409	187	246	3358	4238	25	2547	1207	928	1763	1207	1763	2425	891	332	249	243	233	148	224	1287	533			
26	feel	14499	2301	9659	2468	753	374	251	98	240	2700	3691	51	2386	906	820	213	509	3269	1639	2693	541	278	281	250	190	16	477	1663	409		
27	make	14176	1807	10433	2732	1849	256	367	12	237	234	6368	1864	959	668	261	511	4683	1712	2802	695	608	257	260	204	175	266	150	402			
28	well	14067	2101	5710	1892	1783	362	447	154	229	1776	11	1591	1691	603	751	1211	1425	921	1038	492	222	233	191	319	206	103	566	426			
29	baby	13870	1955	13008	4473	1023	249	493	56	153	4473	7455	14	1022	1634	279	931	2164	8958	3725	4153	329	1263	69	174	313	265	458	1598	140		
30	cause	12720	1888	16847	3181	970	228	326	87	157	6699	7676	23	1494	1240	836	261	548	4200	1253	2941	732	335	178	255	148	89	230	1340	150		
31	rain	12670	1879	13037	2685	756	202	384	169	226	1041	2247	71	1479	937	601	1587	1146	708	274	538	1253	461	384	156	177	130	261	1071	321		
32	need	12632	1805	6727	2468	756	202	308	77	188	3243	449	22	1520	2700	3691	1255	2727	741	324	370	219	233	147	231	1171	325					
33	could	12632	2057	7445	2287	1132	195	281	89	219	2506	2892	12	1255	1117	549	203	410	2202	824	1757	655	320	203	202	132	93	227	82	436		
34	day	11674	1717	5683	1865	191	176	435	216	200	2198	2566	1959	12	1589	1691	598	335	603	1834	936	1180	789	292	110	172	251	138	92	693	475	
35	right	11593	1785	7195	2325	1104	167	423	204	213	1720	2495	26	2138	1148	411	374	651	2316	1062	1947	417	259	164	182	240	347	167	376			
36	wanna	11207	1752	10733	2479	534	158	149	202	202	1204	1680	1658	5	1783	1353	457	411	429	2291	1280	1828	384	251	176	197	168	203	1065	373		
37	heart	11125	1815	8178	3151	1128	147	487	197	202	1204	1680	1658	5	1783	1353	457	411	429	2291	1280	1828	384	251	176	197	168	203	1065	373		
38	hear	11042	1845	8147	2040	2262	231	130	131	134	1469	4823	205	1206	1658	5	1783	1353	457	411	429	2291	1280	1828	384	251	176	197	168	203	1065	373
39	give	9757	1223	6380	3841	557	158	380	66	150	3057	3788	44	1436	443	549	202	407	2686	958	1608	534	357	329	130	134	36	165	1016	291		
40	give	9757	1223	6380	3841	557	158	380	66	150	3057	3788	44	1436	443	549	202	407	2686	958	1608	534	357	329	130	134	36	165	1016	291		
41	evens	9630	1472	4892	1594	842	176	245	103	180	1089	1308	39	2052	689	592	140	302	180	548	1100	548	136	191	124	129	111	113	612	391		

The majority of the words provided throughout the worksheet comes from genres that make up a significant share of the dataset, like Rock and Pop music, but there were a couple outliers that stood out, such as Rap. Although this data was not used for the final product, and was realized to be futile despite taking a lot of time and RAM to render, I found it fascinating given the vast array of words available, as well as the frequency in which the words are used among all given genres. Some of the words provided can also be recognized in the clusters generated too.

## **Conclusion**

In conclusion, the results indicate that the song's genre can be difficult to predict based on the lyrical content of the song. By using approaches such as the Bag-of-words model (with TF-IDF), logistic regression, and the K-means classifier, we were able to make a reasonable estimate on how well a song can be matched to its genre based on the types of vocabulary that is used in the composition. The results indicate that the accuracy scores were not as great, but were improved after redesigning the logistic regression and TF-IDF models. Clusters made from the K-means classifier were shown to have given accurate results based on the features of the data provided.

To summarize the biggest problems with this project, they either involved lack of preparation and knowledge using text analysis libraries (TF-IDF), waiting a long time to gather the data, and variations of certain text messing up the results despite attempts to clean the data while extracting the lyrics. A lack of knowledge on text in relation to the logistic regression model also made it challenging to create a functional model for our results, whereas numerical data was handled much more easily in previous works. Handling large datasets was also a challenge given that previous projects were handled with smaller, and much more manageable, data.

## **Acknowledgements**

The idea of this project came about after skimming through numerous articles that explained how music impacts people's mood. One journal written by Christenson et al. argues that lyrics are important because they reflect the consensus of the audience they are targeting. The lyrical content matters most when the listener can resonate with the connotation of the song. There is also some notion of certain genres being assigned a certain image due to stereotypes involving the content of the song's lyrics, such as a controversial case with rock or metal music being correlated with violence. More information about this journal is linked in *References* as an external link.

Although I was inspired to create a project about the linguistics of music, I initially struggled to pick a topic that could relate to the methods of data science. The reason why I chose to analyze the genre of the song is because it is the best factor that can be predicted based on the lyrical content of the song. When brainstorming, there was the consideration of predicting the artist based on the song, but that proved to be more difficult in practice because of the wider arrangement of artists compared to genres.

The flow of the project was influenced by a similar project done by Alda Sianipar, a graduate data science student from UCLA, who had performed similar approaches in creating the project<sup>[2]</sup>. The only major difference between the inspired project and the project undergone here is that lemmatization was removed; lemmatization is the process of grouping words based on a common root word, like, for instance, "mobile", "mobility", and "mobilize". The reason for skipping over this was because of unfamiliarity with lemmatization.

Assistance from this project was also possible with other articles and sources like GitHub and *Towards Data Science*. One of the biggest motivations for using these sources came from remodeling my approaches when trying to display the best data possible using TF-IDF and logistic regression. The rough

draft implementation of the logistic regression and TF-IDF model was inspired by a GitHub-uploaded IPYNB worksheet from data scientist Idil Ismiguzel, who also wrote an article in relation to how and why she chose logistic regression for text classification<sup>[3][4]</sup>. The idea of clustering also came about when revisiting the sources provided for Labs 8 and 9, with an article that explains how K-Means clustering works and is demonstrated<sup>[5]</sup>.

## References

- [1] Song lyrics from 79 musical genres, *Kaggle*,  
<https://www.kaggle.com/datasets/neisse/scrapped-lyrics-from-6-genres?select=lyrics-data.csv>
- [2] Sianipar, Alda. “Predicting a Song’s Genre Using Natural Language Processing”, *Better Programming*, 24 July 2019,  
<https://betterprogramming.pub/predicting-a-songs-genre-using-natural-language-processing-7b354ed5bd80>
- [3] Ismiguzel, Idil. “Applying Text Classification Using Logistic Regression”, *Analytics Vidhya*, 7 May 2020,  
<https://medium.com/analytics-vidhya/applying-text-classification-using-logistic-regression-a-comparison-between-bow-and-tf-idf-1f1ed1b83640>
- [4] Text Classification on Amazon Find Food Reviews Dataset, *GitHub*,  
<https://github.com/Idilismiguzel/NLP-with-Python/blob/master/Text-Classification.ipynb>
- [5] Foley, Daniel. “K-Means Clustering: Making Sense of Text Data using Unsupervised Learning”, *Towards Data Science*, 8 Feb 2019,  
<https://towardsdatascience.com/k-means-clustering-8e1e64c1561c>
- **EXTERNAL LINK** Christenson et al., “What has America been singing about? Trends in themes in the U.S. top-40 songs: 1960-2010”, *Sage Journals*, 23 Jan 2018,  
<https://journals.sagepub.com/doi/full/10.1177/0305735617748205>