

COMP3517 Written Report

Abstract

Music, similar to other artistic works, offers a lens into the lives of the artists that produce these poetic experiences. Music can also reveal important details about a country's history and the experience of everyday people in this country. This report details an attempt at using various off the shelf NLP tools to analyse the change in language in popular American musical works and compare that to the greater context of the American historical experience between the years of 1890 to 2019. Due to the nature of the artists that tend to populate the mainstream music charts and in attempt to offer a diverse perspective, I also include popular African American artists and rap artists that tell unique stories that are less glamorised than mainstream pop songs. These artists were chosen based on their cultural effect on hip hop, soul, R&B or funk and their involvement in African American music culture.

1 Introduction

When we think of a cultural phenomenon the we think of celebrities and popular theatrical works. Movies like Titanic and celebrities like Leonardo DiCaprio are very influential on modern day pop culture. This is no different in the realm of music. Nowadays, the most popular and well known musical artists globally are (for the most part) American. Part of this is because English is the most widely spoken language but also because American culture has managed to spread to most parts of the world and so have the works of American artists and celebrities. These artists use their music to tell stories that, accompanied by enticing instrumentals, are able entertain millions and create massive followings. It does beg the question has popular music stayed the same lyrically? Talking lyrically or grammatically one may assume that the lyrics have changed just as much as how we speak has changed. However, written word and spoken word for the most part has not changed as much as you'd think. Meaning behind

some words and the format of sentences may have changed. However, the core elements have remained the same as I hope to show is the same for music. The following sections detail the results obtained when analysing the changes in the lyrics using various off the shelf NLP tools. It is worth mentioning that a larger dataset could have been compiled for both lists used in the analysis of music lyrics in hind sight.

2 Background

2.1 NLP

NLP or "Natural language processing" is a field that began as a number of different standalone fields. The list of fields include, linguistics, machine learning, entity recognition, etc. The aim of NLP is to make use of different approaches in these different fields to teach a machine to extract meaning from text. NLP has come a long way from its inception.. Below is Fig. 1 which shows an image taken from Khurana et al. That shows the recent developments in the realm of NLP

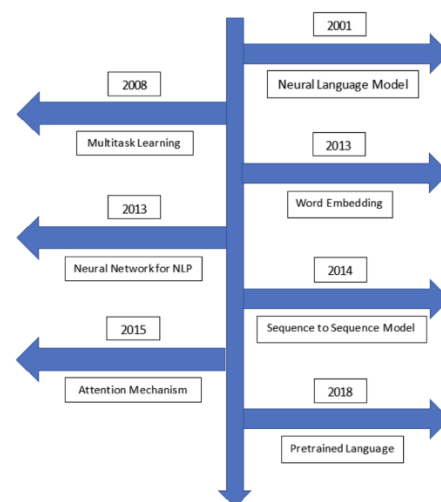


Figure. 1: Showing NLP developments from 2001 to 2018

These developments have made it easier and quicker to make use of NLP techniques such as the ones discussed in this project to automate extraction of meaning from bodies of text.

2.2 American history

America as one of the most influential countries in the world and has a very complicated history. From the importing of slaves from Africa brought by Dutch ships in the 1600s, to the abolishment of slavery and later introduction of Jim crow and segregation, to a country that prides itself on its diversity of ethnicities, it is fair to say that America has a long and complicated history. Despite this or because of this (depending which way you want to look at it). American music tends to be quite diverse. It has aspects of many different cultures and perspectives. If one is looking for the African-American perspective then you would look to soul, funk or hip-hop. If you're looking for the general/societal perspective usually you could look to the mainstream music which deals with many different themes (usually specific to the artist). If you're looking for a more Texan perspective then you might look at country music. Although the way I put it almost makes it seem like "black and white", in reality America is more than that. However, due to how long African-Americans have been in America and due to the way they were brought to America, their stories told and perspective which is reflected in their music tends to vary far more from what is portrayed in mainstream American music. Where popular mainstream music often talks about relatable concepts common to nearly all people (like hope or love), hip hop often describes elements of struggle, strife and a journey to success. Of course this also could be because I just do not know about Mexican/Latina/Spanish music as well as I do Hip Hop therefore I have a bias towards African American music/culture.

3 Methodology

In this section we discuss what methods were employed in order to be model the changes in language.

3.1 POS Tagging

POS (Part-Of-Speech) tagging is a useful tool for investigating the types of speech that are present in a body of text. In this context I am using it to try

identify changes in grammatical tendencies in the different popular musical works. For POS tagging the function "process_music" shown in Figure 2 takes a number of arguments and returns the POS distribution for each of the decades or artists depending on arguments declared (for details about arguments passed to functions for intended output refer to the README.txt). In some cases it is easier to use this method for grammatical analysis. As this method allows observations to be made and assumptions be drawn from specific distributions by looking at the more commonly occurring POS.

3.2 Topic Modelling

In order to topic model I make use of gensim's LdaModel which models the probabilities of specific words appearing for the latent topics present in the dataset. These are topics that are not directly observable and are obtained by the model through analysis of the semantic structures present in the text. Topic modelling using the Gensim library is yet another way of observing what words are commonly used by artists in their music. It is useful for identifying relationships and emotions that humans may not usually be able to detect.

3.3 Sentiment Analysis

Sentiment analysis is a highly dependable technique that allows us to extract the general sentiment expressed within the lyrics. Using this method it will be possible to check if certain historical event should have correlated with an increase or decrease in specific sentiments. Sentiment analysis tends to be more dependable because artists rarely use words that express one sentiment to express another sentiment. It is also because the dimensionality of what the machine is trying to predict is significantly reduced when compared to the other techniques. I make use of the vaderSentiment module which has the capability to measure sentiment present in bodies of text. It is worth noting that, although sentiment analysis is not flawless. It works well on general text however, in cases where literary devices like metaphors are often used and the representations of meaning are hidden in the words there is a possibility of incorrect scores being given. There is also a possibility of bias towards certain words, as later discussed, interfering with the scores.

4 Data

4.1 Collecting data

In order to collect the data for this investigation it was necessary to first decide on the what data would be appropriate to include. Since there wasn't an already existing corpus of data that would fit this specific task creating an appropriate corpus was necessary. Since there were only 130 songs that had to be included on the side of songs by decade (technically 128 excluding the non English ones) this was done manually as it would take less time to do it this way than create a script to automate the process. Furthermore, implementing web scraping would have taken time to learn and so it would have taken longer. The first corpus consisted of 10 songs per decade from the decades 2010s to 1890s. The list for the top 10 list for each decade was taken from "DavesMusicDatabase.com". Whether these songs were really the best 10 songs of each decade wasn't all that important. The most important detail was that they were songs that were quite popular post release as this lets us draw a relation between popularity and grammatical patterns. It is worth noting that the second dataset which was the hip hop dataset, was a set of 24 songs chosen to add a unique perspective to the investigation. The songs were chosen in pairs where I attempted to use songs released further apart from each other where possible. There were a total of 12 artists chosen and 2 songs were chosen for each of them. The lyrics were collected in the same fashion as the first dataset. The artists were chosen based on who fit the following criteria. The artist had to be either an influential African American Artist or an influential hip hop artist. They also had to have an appropriate understanding of the African American perspective of life in America. Now in the case of Eminem, he was chosen because of who he grew up with and how this lead him to become a hip hop artist. He is also a cornerstone of hip hop and one of the greatest hip hop artists of all time. To snub him from the list due to his ethnicity seemed inappropriate.

4.2 Processing data

In order to use the NLP techniques that would allow us to extract meaning from the lyrics it is important to pre-process the data so that characters and/or words that don't contribute any

meaning do not potentially confuse or interfere with the algorithm. In Figure 2 there are a number of different defined functions. Firstly the "remove_newlines" function removes any newline characters in the text and replaces it with a space between the last character of the previous line and the following line. This was done because the platform that had these lyrics stored had them

```
##### START OF FUNCTION CODE BLOCK #####
# Function for removing any newline characters from corpus
> def remove_newlines(strings):
# Function for removing all escape characters from corpus
> def remove_escape_characters(strings):...
> def tokenize_and_remove_stop_words(lyrics):...
# tokenize_lyrics(e)
# remove_lyrics = ['don', 'asasou', 'ta', 'na', 'wan']
# Function for removing words in a song
> def remove_words(array, words_to_remove):...
> def tag_parts_of_speech(word_array):...
> def lematize_and_stem(lyrics):...
# It's all well and good that we have counted the number of times these different parts of speech have appeared but we
# index position 0 and 2 have 1 less song therefore 2010s -> [0:8], 1990s -> [10: 27]
# Decade indices must be in ascending order. So the affected indices list cannot have a later decade come before an earlier
> def decade_pos_tagger(array, decade_length=[]):...
# Function should return the number of each type present within a 2d array where
# the first dimension is the songs present in passed decade and the second is the song lyrics of each song
> def pos_tag_decade(lyrics=[], decade_separated=True, artists_names = ["Laurn Hill", "Nas", "Tupac", "Eminem", "Biggie"])
> def topic_model(process=True, dictionary1=pre_processed_hpop_lyrics, dictionary2=pre_processed_lyrics):...
> def sentiment_analyse(n_songs_per_artist=2, decade_separated=True, decade_length=[9,10,9,10,10,10,10,10,10,10,10,10,8],)
> def process_music(pre_processed_data=pre_processed_lyrics, process=True, print_dist=False, decade_length=None, decade_
```

Figure. 2: Showing a list of defined functions

stored in centred paragraphs which lead to newline characters. There were also some weird escape characters similar to the newline character "\n" and so the function "remove_escape_characters" removes these characters from the corpus. Those two functions are the only ones used for pre_processing.

5 Results and Discussion

5.1 Sentiment Analysis

The results of the sentiment analysis showed a pattern of general positivity in the lyrics of the most popular song of each decade. This could be because many popular mainstream songs tend to

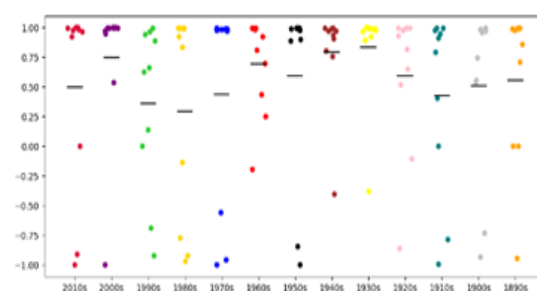


Figure 3: showing distribution of sentiments for the decade songs

deal with happier topics such as love or hope than in hip hop which tends to have songs about reflection and tells stories about struggle. For example Nas's "N.Y. State of Mind" is a song describing the city of New York what was in his head when he lived in there when he wrote the song. In the 90s it was a city riddled with high violent crime rates where the streets were seeing the effects of the "Crack epidemic". Mainstream music misses out on this side of the American story. Where rap as an art form often tells the story of individuals and how they got to where they are, mainstream often is more about experiences that everyone can relate to. Stories about love, partying or songs that are just plain catchy often dominate these platforms. There are some values missing due to weird interactions in the code and so there is a possibility that the extent to which some decades may have had negative values may have been inaccurate. However this is not really an issue since the decade most affected by this (1890s) is only missing 2 values in an era that wasn't of extreme historical importance when compared to some of the other decades. The decade with the least negative sentiment in American music was the 60s which would make sense as that was the decade of the hippies and the peace movement it would then make sense that out of all the eras that are supposed to have a general positive outlook on life, where everyone would be producing some of the happiest music would be the 60s. The hip hop sentiments were evenly spread. However there is a possibility that the classifier was biased towards the profanity expressed in the Hip Hop lyrics. The song **** in Paris had a value of -0.9993 due to the presence of a large amounts of profane language. For this reason the results for this section as it pertains to the second dataset have been omitted but can be found in the README.txt file.

5.2 Topic modelling

	Song 1	Song 2	Song 3	Song 4	Song 5	Song 6	Song 7	Song 8	Song 9	Song 10
2010s	0.9248	0.9994	0.9981	0.9997	0.9865	0.9776	-0.9991	-0.908	0.9635	0
2000s	0.9992	0.9941	0.9994	0.9981	0.5373	0.9966	0.9704	-0.9972	0.9478	0.9988
1990s	-0.6854	0.9405	-0.921	0.9971	0.8885	0.6249	0.6602	0.9661	0.1391	0
1980s	0.9953	-0.7731	0.9978	-0.1333	0.9913	0.9959	0.923	-0.922	-0.9698	0.8359
1970s	0.9989	-0.5538	0.988	0.9803	-0.998	0.9747	-0.956	0.9844	0.987	0.9893
1960s	-0.1943	0.987	0.9966	0.9985	0.8126	0.9906	0.4356	0.25	0.926	0.7003
1950s	-0.9954	0.9952	0.9959	-0.8407	0.8915	0.9947	0.901	0.9898	0.986	0.9778
1940s	0.9959	0.9392	0.9705	0.9062	0.8045	0.9697	-0.4019	0.9894	0.7555	0.9957
1930s	0.9274	0.9994	0.9897	-0.3792	0.9957	0.9878	0.9663	0.9885	0.8957	0.9795
1920s	0.9982	0.9962	-0.1033	0.8153	0.6537	-0.858	0.9975	0.9765	0.5279	0.5207
1910s	0.9937	0.9515	0.9777	0.9136	0	-0.9879	0.4082	0.7964	-0.7815	0.9944
1900s	0.9976	0.9623	-0.9303	0.5588	0.5373	0.9872	0.7469	0.9712	0.9865	-0.7311
1890s	0.8591	0.9918	0.9875	0.9972	0.9802	-0.9411	0.9962	0.7096	0	0

Topic modelling is a technique that identifies the topics present in the latent space of the corpus of text as previously described and the probabilities of words associated with those topics. Now, depending on the number of topics you want the

Word	Probability	Probability	Probability	Probability	Probability	Probability	Probability	Probability
i	0.058	0.029	0.038	0.064	0.041	0.029	0.029	0.083
and	0.018	0.007	0.005		0.023	0.014	0.007	0.014
im	0.018	0.007	0.012	0.018	0.017	0.017	0.023	0.015
got	0.014			0.012	0.008	0.011	0.009	0.017
like	0.012		0.007	0.008	0.013	0.018		0.016
know	0.012				0.016	0.008		0.009
you	0.011				0.01			
dont	0.007	0.007		0.009	0.009	0.009		0.009
ni**a	0.007				0.009	0.008		

Figure 4: Showing the probability distribution of words in the hip hop corpus

machine to attempt to model, you can pass the function a different number. However, there is a limit to the number of latent topics identified. Therefore even if you give it an argument of `n_topics = 200`, it will only model 200 if there are 200 topics present in the latent space. Otherwise it will pass back as many as it can find. For the purpose of this project I chose 20 topics for both bodies of text. I chose this because that was what I discovered to be the limit for the number of topics in the latent spaces. For the hip hop topic modelling, the words that had the highest probability can be seen in Figure 4. With subject "i", occurring with the highest probability. Since rap is an art form that tends to describe the individual's perception and experience this is to be expected. In the list of popular songs by decade, the words the distribution of occurring words was spread slightly more. However, "I" still showed the highest probability of occurring.

What the output of the topic modelling looked like. The one thing both datasets demonstrated was that "I" occurred at the highest rate between both datasets. This is likely because music most often is used to describe one's experience. It is usually used to tell stories from the individuals perspective making the story personal.

Figure 7: Sentiment analysis table for decade songs

Based on what was observed in this project it is appropriate to say that music as a whole on the grammatical level has not changed as much as one would think. It was an art form that was used to tell a story of a subject or number of subjects and seems to have remained that way. Furthermore, HipHop and mainstream music have shown great similarities in most categories. In the future, maybe when there is time to find a way to adjust for the presence of specific words due to the nature of HipHop, it would be interesting to compare the sentiments of the popular HipHop songs to the sentiments of the popular mainstream songs. As for the overall positivity or negativity of popular works, generally, mainstream music has more

positive music that goes popular than it does have negative. Probably because popular music in the mainstream tends to be catchy and upbeat. This is left up to speculation for now. In this paper we have shown the differences and similarities in different musical works from the 1980s to 2019 across a number of different genres.

9 References

[1] Top 10 songs by decade:

<https://davesmusicdatabase.blogspot.co.uk/2012/04/top-songs-by-decade-1900-present.html>

[2] Most influential HipHop artists:

<https://www.standard.co.uk/culture/music/the-most-influential-hip-hop-artists-of-all-time-a3863356.html>

[3] American history: [American History Timeline 1800-1900 \(american-history.net\)](http://american-history.net)

[4] African American history: [Black History: Facts and People | HISTORY.com - HISTORY](http://black-history-facts-and-people.com)

[5] Diksha Khurana, Aditya Koli, Kiran Khatter, Sukhdev Singh, Natural language processing: State of the art, current trends and challenges, [Natural language processing: state of the art, current trends and challenges \(springer.com\)](http://springer.com)

[6] Code for plotting a vertical scatter:

<https://stackoverflow.com/questions/44982574/how-to-plot-vertical-scatter-using-only-matplotlib>