

# Lab 5 – MATH 240 – Computational Statistics

Avery Johnson  
Colgate University  
Department of Mathematics  
aqjohnson@colgate.edu

## Abstract

This lab aims to determine which band contributed most to the song “Allentown” by analyzing sound and lyric features across multiple datasets. We compiled music analysis extracted using Essentia, along with Essentia model data, and integrated it with LIWC-based lyric analysis. Our methodology involved feature selection to identify distinguishing characteristics among three bands and summary visualizations to support our conclusions. Results indicate that Manchester Orchestra had the largest impact on sound, while the Front Bottoms had the largest lyrical impact.

**Keywords:** libraries; character objects; for() loops; vectors/lists; tidyverse; and summarizing data

## 1 Introduction

In 2018, the Front Bottoms, Manchester Orchestra, and All Get Out collaborated on a track called “Allentown.” To determine which band contributed most to the song, a data-driven analysis was conducted on their previous tracks. This task utilized Essentia (Bogdanov et al., 2013), an open-source tool for music analysis, to extract audio features from 181 tracks. Given the large dataset, we developed a batch file to automate the command-line extraction process, streamlining data collection. We then cleaned and integrated the three key data sets: Essentia Extractor data (Bogdanov et al., 2013), which captures detailed spectrogram-based sound analysis, Essentia Model data (Alonso-Jiménez et al., 2020), and LIWC data (Boyd et al., 2022), which analyzes lyrical content. Building on this foundation, Lab 5 focused on feature selection and analysis to summarize the data, and therefore conclude which band had the largest impact on “Allentown.”

## 2 Methods

### 2.1 Task 1: Automating Data Extraction

To automate the execution of Essentia for each audio track, we created a batch file that systematically processed all .WAV files. Using the `stringr` package in R (Wickham, 2023), we extracted subdirectories for each album, identified and counted the .WAV files, and constructed the command line calls to execute the Essentia program for each track using a `for` loop.

### 2.2 Task 2: Compiling and Cleaning Data

To analyze musical contributions, we compiled extracted data into a consolidated data frame. This involved processing JSON files, integrating Essentia extractor and model data, and incorporating LIWC lyric analysis. After extracting relevant audio features, we loaded and cleaned the data. This involved using the `rowMeans()` to average feature values across different extractors. The data was then merged into a single data frame using the `merge()` function. Finally, using `write.csv`, the final dataset was stored in training and testing sets for further analysis.

### 2.3 Task 3: Summarizing and Analyzing Feature Ranges

Building on previous labs, we examined feature distributions across the dataset to determine which characteristics best differentiated the bands. To identify distinguishing features, we analyzed the range of values for each feature and classified them as out of range, outlying, or within range based on the min, max, lower fence, and upper fence thresholds. Specifically, features where two bands were out of range while one band remained within range were considered strong indicators of influence. From an initial selection of 20 features, we identified four key lyrical features (from LIWC) and 4 sound features (from Essentia) that were most relevant. We compiled these features into summary tables using `xtable` and generated violin and box plots using `ggplot2` to visually assess differences across bands. The red dashed lines in the plots represent the corresponding feature values for “Allentown,” allowing us to determine which band’s characteristics aligned most closely with the song.

## 3 Results

### 3.1 Task 1 Results

The R script successfully identified album subdirectories and filtered .WAV files. It then generated batch commands for each track, which were saved in a text file named `batfile.txt`. These commands generated the Essentia program for each track, saving the corresponding output as JSON files. The process was automated, enabling batch processing for all audio tracks.

### 3.2 Task 2 Results

Task 2 involved first processing the JSON output for a single track. The artist, album, and track name were obtained from the file name, and relevant audio features were successfully extracted. An example of extracting one of these features from this song is shown below.

Feature	Value
Artist	The Front Bottoms
Album	Talon Of The Hawk
Track	Au Revoir (Adios)
Avg. Loudness	0.5450

Table 1: Extracted Audio Feature for “Au Revoir (Adios)”

This process was scaled to the full dataset by extracting, cleaning, and merging the Essentia Extractor, Essentia Model, and LIWC data into a single dataframe called `merged_df`, with 181 tracks and 140 features. Training data `trainingdata.csv` excluded “Allentown,” while testing data `testingdata.csv` contained only “Allentown.” These successfully written CSV files will help determine which band contributed most to the song.

### 3.3 Task 3 Results

Our analysis for task 3 involved identifying features where only one band remained within range. We specifically focused on four lyrical and four audio-based features. This analysis for the sound data can be seen in Table 2 below.

artist	description	feature
All Get Out	Out of Range	spectral_rolloff
Manchester Orchestra	Within Range	spectral_rolloff
The Front Bottoms	Out of Range	spectral_rolloff
All Get Out	Outlying	dissonance
Manchester Orchestra	Within Range	dissonance
The Front Bottoms	Out of Range	dissonance
All Get Out	Outlying	average_loudness
Manchester Orchestra	Within Range	average_loudness
The Front Bottoms	Outlying	average_loudness
All Get Out	Outlying	chords_strength
Manchester Orchestra	Within Range	chords_strength
The Front Bottoms	Out of Range	chords_strength

Table 2: Summary of Features Identifying Influencing Band

The full set of feature comparisons for both sound and lyrical features, along with the min, max, and fence values, are in the Appendix. Violin plots were created to show the range of values for each band, with red dashed lines representing the feature values for “Allentown.” We analyzed four lyrical features (`positivewords`, `OtherP`, `Perception`, and `conj`) and four sound features (`spectral_rolloff`, `average_loudness`, `chords_strength`, and `dissonance`). Figure 1 shows the violin plot for `average_loudness`, where “Allentown” aligns most closely with the Manchester Orchestra, suggesting they contributed most to the track’s loudness. Similar analyses for the lyrical and other sound features can be found in Figures 2 and 3 in the Appendix. The results indicate that the Front

Bottoms most influenced the lyrics, while the Manchester Orchestra had the greatest impact on the sound.

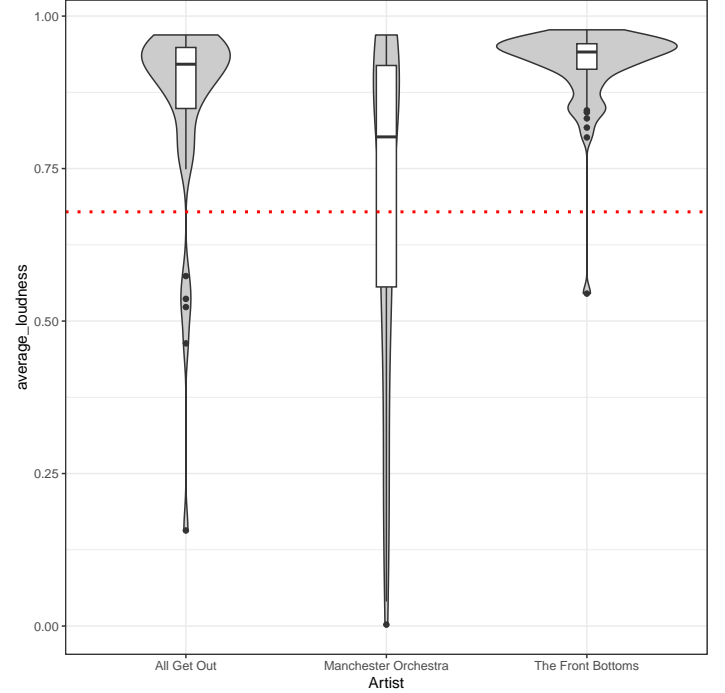


Figure 1: Average Loudness by Band

## 4 Discussion

This analysis identified key musical and lyrical features suggesting which band contributed most to “Allentown.” Lab 2 demonstrated the process of building a batch file and processing a single JSON output, while Lab 3 showed the scalability of this methodology by analyzing the full dataset. Statistical analysis revealed that Manchester Orchestra most influenced the sound of “Allentown,” with features consistently in range, while the Front Bottoms impacted the lyrics, as their features aligned with “Allentown” and others did not. The compiled summary and visual representations provide a quantitative basis for these conclusions. However, the way the data were summarized is specific to our approach. Alternative methods for feature selections could potentially yield different results. This highlights the importance of methodological choices and suggests that future studies could explore other analytical techniques for comparison.

## References

- Alonso-Jiménez, P., Bogdanov, D., Pons, J., and Serra, X. (2020). Tensorflow audio models in essentia. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 266–270. IEEE.
- Bogdanov, D., Wack, N., Gómez Gutiérrez, E., Gulati, S., Boyer, H., Mayor, O., Roma Trepal, G., Salamon, J., Zapata González, J. R., and Serra, X. (2013). Essentia: An audio analysis library for music information retrieval. In Britto, A., Gouyon, F., and Dixon, S., editors, *14th Conference of the International Society for Music Information Retrieval (ISMIR)*, pages 493–498, Curitiba, Brazil. International Society for Music Information Retrieval (ISMIR).
- Boyd, R. L., Ashokkumar, A., Seraj, S. M., and Pennebaker, J. W. (2022). *The development and psychometric properties of LIWC-22*. Austin, TX, p. 10.
- Wickham, H. (2023). *stringr: Simple, Consistent Wrappers for Common String Operations*. R package version 1.5.1.

## 5 Appendix

	artist	min	max	LF	UF	description	feature
1	All Get Out	935.91	2520.04	701.91	2767.30	Out of Range	spectral_rolloff
2	Manchester Orchestra	518.87	2566.67	151.27	2083.17	Within Range	spectral_rolloff
3	The Front Bottoms	927.04	3190.29	740.58	2421.46	Out of Range	spectral_rolloff
4	All Get Out	0.40	0.48	0.44	0.50	Outlying	dissonance
5	Manchester Orchestra	0.37	0.48	0.36	0.53	Within Range	dissonance
6	The Front Bottoms	0.43	0.48	0.44	0.49	Out of Range	dissonance
7	All Get Out	0.16	0.97	0.70	1.10	Outlying	average_loudness
8	Manchester Orchestra	0.00	0.97	0.01	1.46	Within Range	average_loudness
9	The Front Bottoms	0.55	0.98	0.85	1.02	Outlying	average_loudness
10	All Get Out	0.47	0.59	0.47	0.58	Outlying	chords_strength
11	Manchester Orchestra	0.48	0.62	0.45	0.63	Within Range	chords_strength
12	The Front Bottoms	0.48	0.57	0.46	0.58	Out of Range	chords_strength
13	All Get Out	0.91	10.68	-0.51	9.76	Out of Range	conj
14	Manchester Orchestra	0.00	14.43	0.74	10.98	Outlying	conj
15	The Front Bottoms	0.00	12.31	-1.17	13.90	Within Range	conj
16	All Get Out	4.67	20.89	4.14	18.91	Out of Range	Perception
17	Manchester Orchestra	0.00	28.37	-1.06	23.87	Within Range	Perception
18	The Front Bottoms	4.27	22.56	3.74	19.66	Out of Range	Perception
19	All Get Out	0.00	14.42	-1.50	2.50	Outlying	OtherP
20	Manchester Orchestra	0.00	7.69	-1.64	2.73	Outlying	OtherP
21	The Front Bottoms	0.00	12.50	-3.48	7.20	Within Range	OtherP
22	All Get Out	2.00	58.00	-1.50	18.50	Outlying	positivewords
23	Manchester Orchestra	0.00	27.00	-3.00	13.00	Outlying	positivewords
24	The Front Bottoms	1.00	34.00	-14.50	37.50	Within Range	positivewords

Table 3: Full Summary of Key Features Identifying the Influencing Band

## Lyrical Features

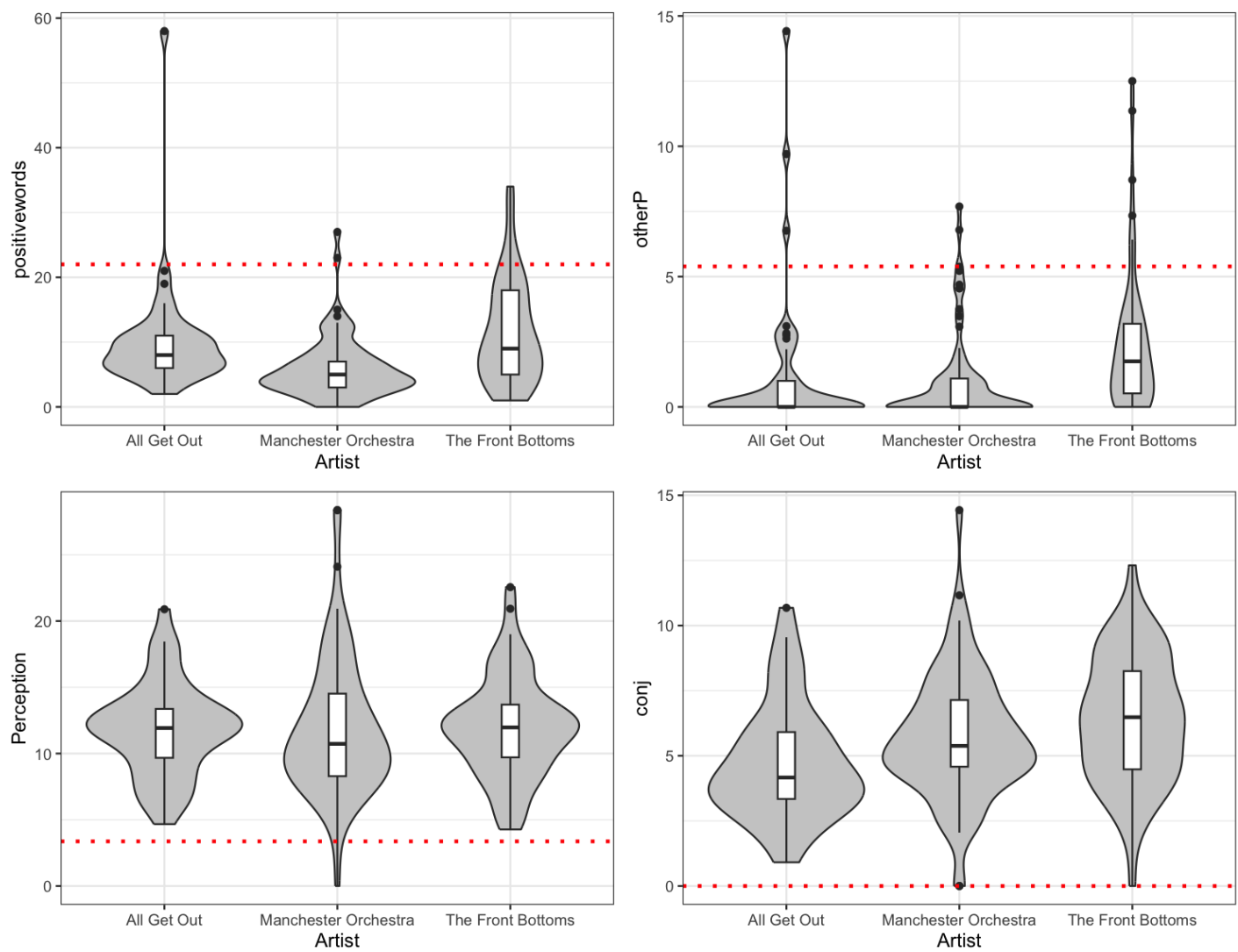


Figure 2: Lyrical Features

### Sound Features

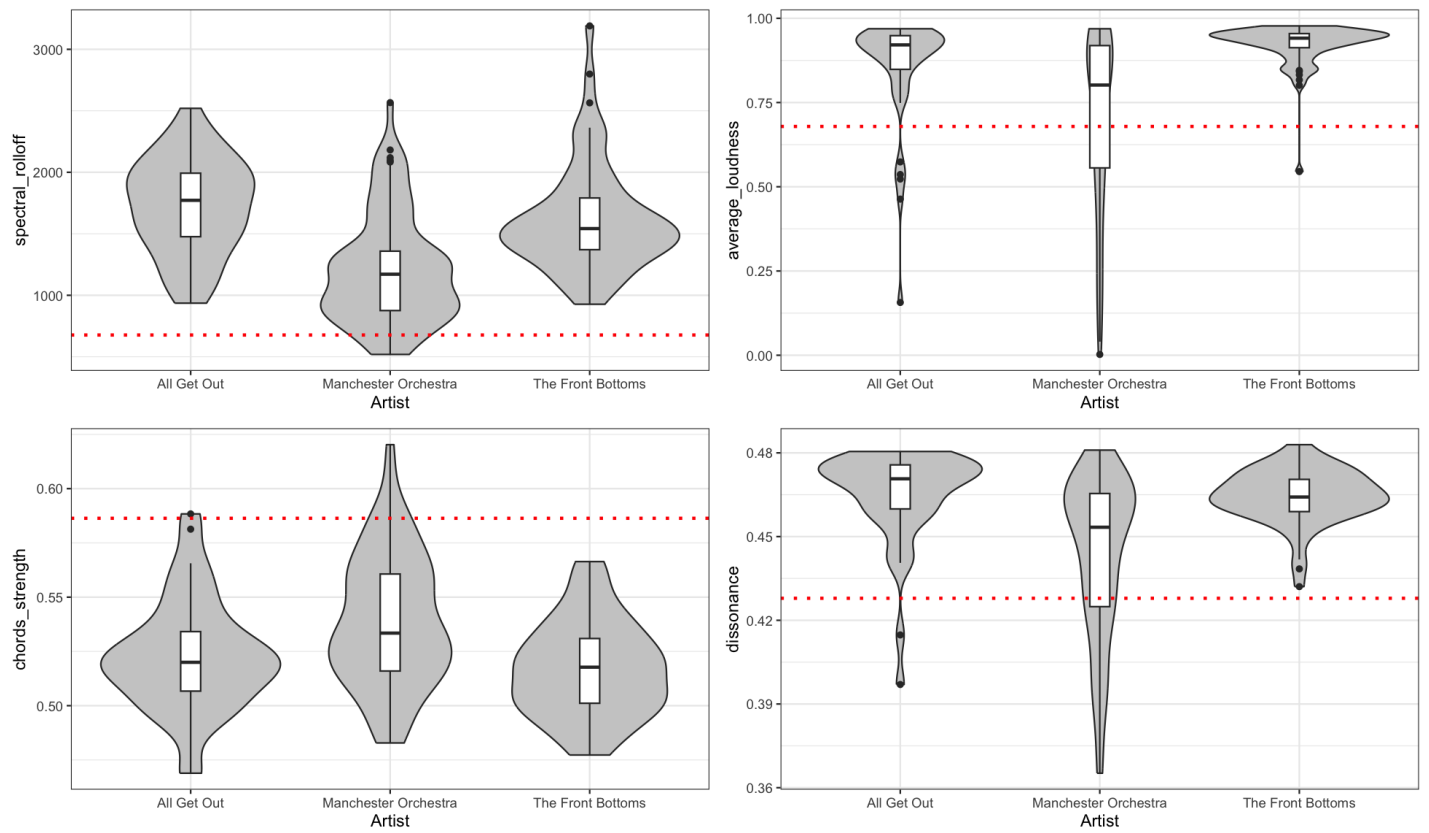


Figure 3: Sound Features