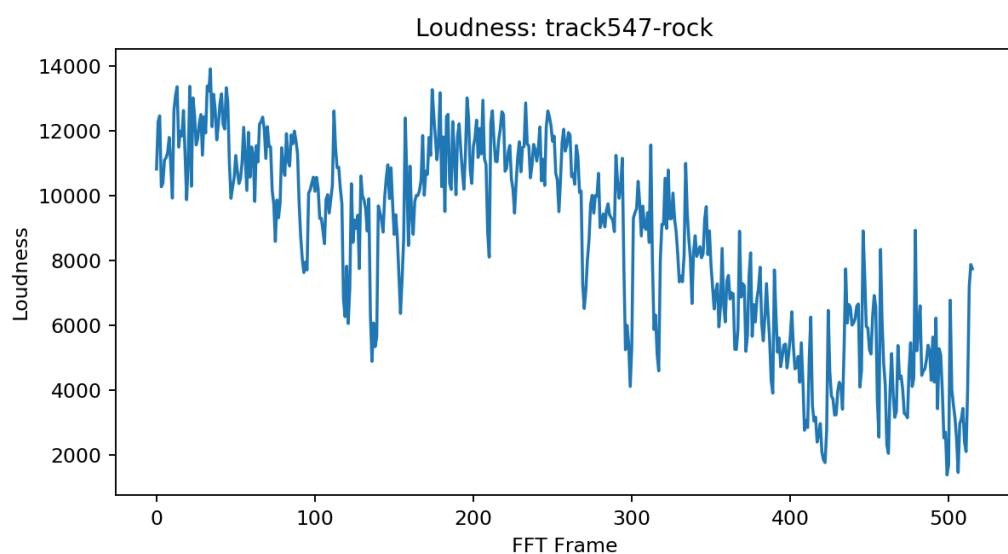


ECEN 4532 - Lab 1: Sound Processing in Python

Andrew Teta

January 28, 2019



Contents

Introduction	3
Format	3
0.1 Background	3
1 Time-domain Analysis	3
1.1 Loudness	3
1.2 Zero Crossing Rate	7
2 Spectral (Frequency) Analysis	11
2.1 Spectrogram	13
2.2 Spectral Centroid and Spread	16
2.3 Spectral Flatness	20
2.4 Spectral Flux	22
3 Psychoacoustics	24
3.1 The mel/Bark Scale	24
3.2 The Cochlear Filterbank	24
4 Conclusion	26
5 Appendix	27
5.1 Appendix A: Main	27
5.2 Appendix B: Sample Extraction	29
5.3 Appendix C: Time-domain Analysis	29
5.3.1 Loudness	29
5.3.2 ZCR	29
5.4 Appendix D: Frequency-domain Analysis	30
5.4.1 Spectrogram	30
5.4.2 Spectral Centroid and Spread	31
5.4.3 Spectral Flatness	32
5.4.4 Spectral Flux	32
5.5 Psychoacoustics	33
5.5.1 Mel Filter Banks and Coefficients	33
References	34

Introduction

In this discussion, we will be exploring some techniques recently developed in the audio industry to organize, search, and classify large music collections. We are developing some statistical methods, along with frequency analysis to automatically characterize songs and genres.

0.1 Background

Fundamentally, we begin with a '.wav' file. This is one variety of many signal encodings, generated by sampling a physical sound wave using a microphone. The microphone detects pressure differences within a medium (sound) and converts them into electrical voltages. As the voltages are converted into a digital representation, they are quantized at a sampling frequency (fs). Furthermore, the Nyquist sampling theorem dictates that for a signal to be represented accurately in quantized form, it must be sampled at twice the rate of the highest frequency component. For example, a dolphin can only hear sound in the range $7kHz - 120kHz$, so if we want to record the sound that a dolphin would be able to hear, we have to sample at a rate $fs = 240kHz$ or higher.

In this lab, we consider a small collection of music, organized by genre, in '.wav' format. Each song is sampled at $fs = 11,025Hz$. Narrowing our data set even further, we will extract a 24 second sample from the middle of each song. We will then implement a short-time Fourier transform (STFT), which divides the sample into $N = 512$ samples, or $46ms$ and call each 46 ms interval a 'frame'. The STFT is a good way to obtain frequency spectrum data, while maintaining a level of time-domain relevance. For the purposes of this lab, the frames will be non-overlapping.

1 Time-domain Analysis

We begin by extracting a $24s$ sample from the middle of each song using the python function `scipy.io.wavfile.read`. See section 5.3 for Python implementation.

1.1 Loudness

To get a sense of 'loudness', we will compute the standard deviation over a 'frame' of size $N = 512$ defined as:

$$\sigma(n) = \sqrt{\frac{1}{N} \sum_{m=0}^{N-1} [x(nN + m) - E[x_n]]^2} \quad \text{with} \quad E[x_n] = \frac{1}{N} \sum_{m=0}^{N-1} x(nN + m) \quad (1)$$

See section 5.3.1 for Python implementation.

Results The output of the loudness calculation for one song of each genre is shown in the figures below.

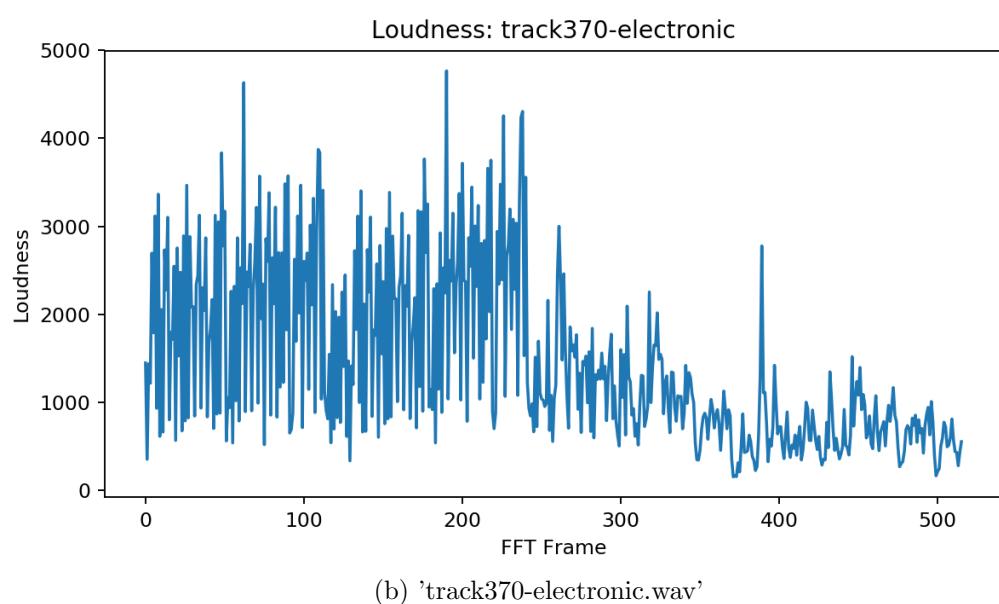
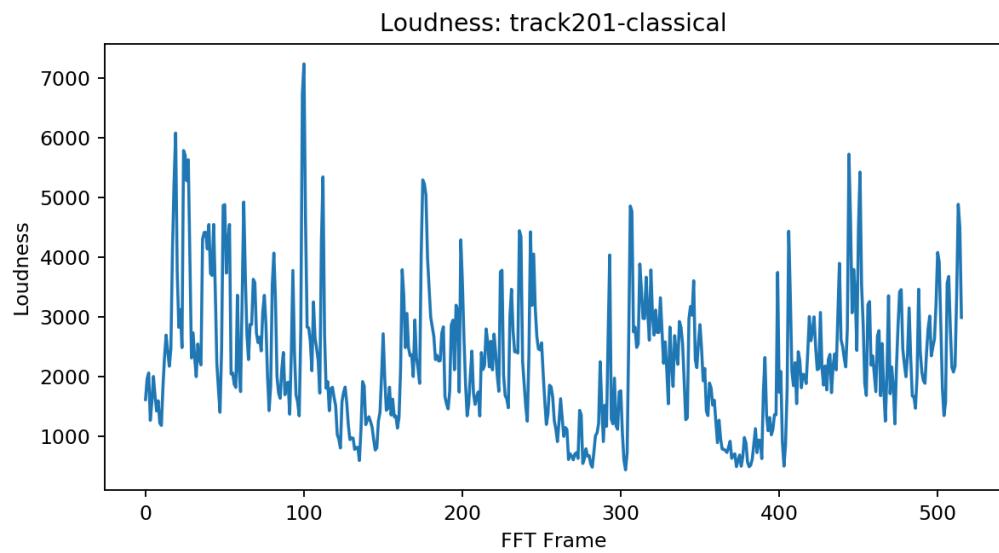


Figure 1: Loudness vs. frame

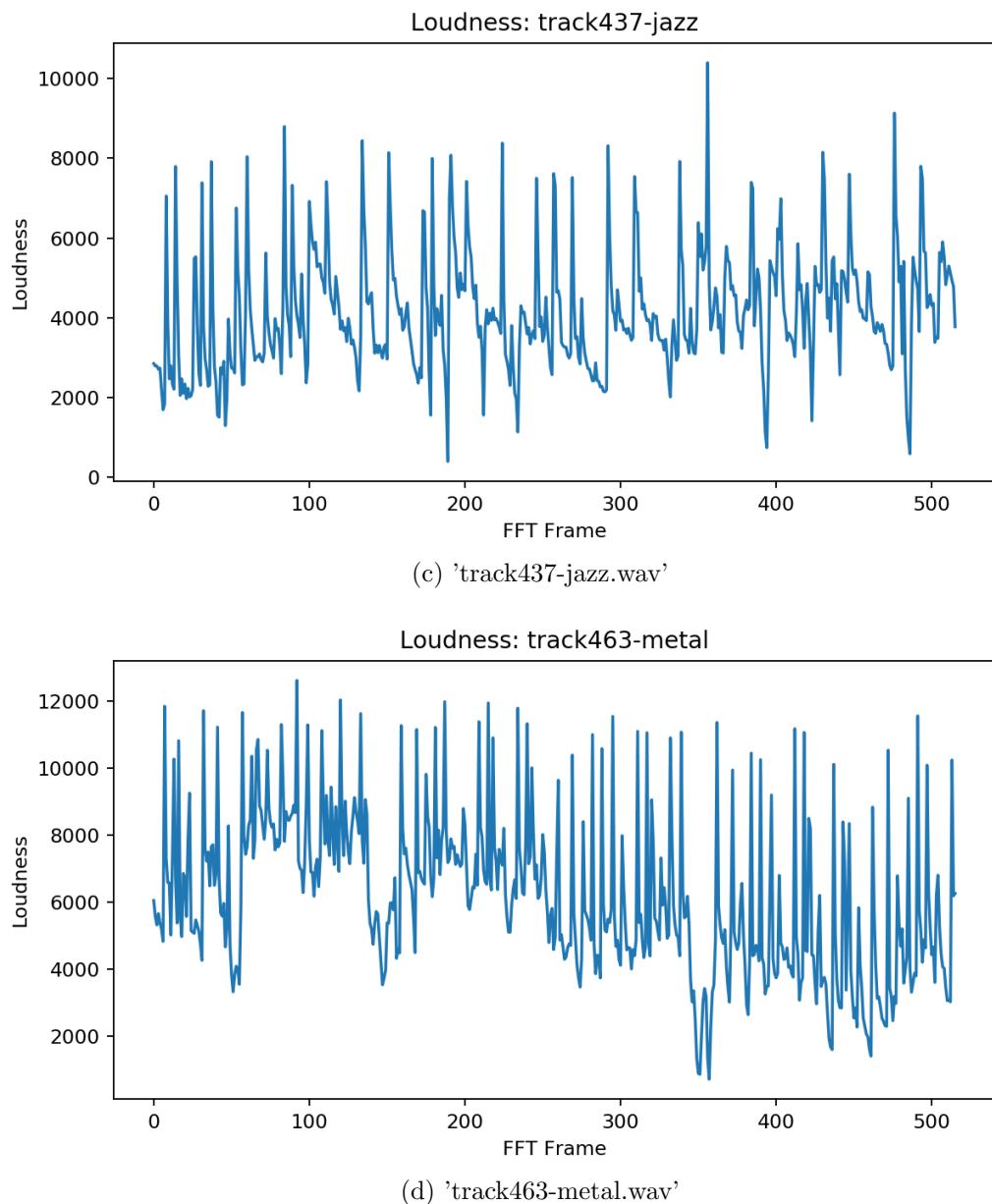
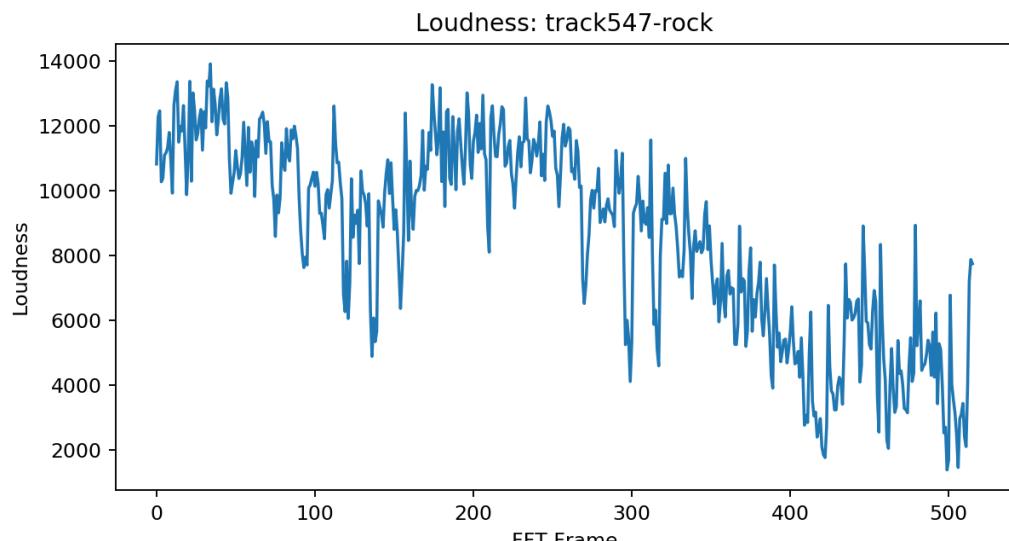
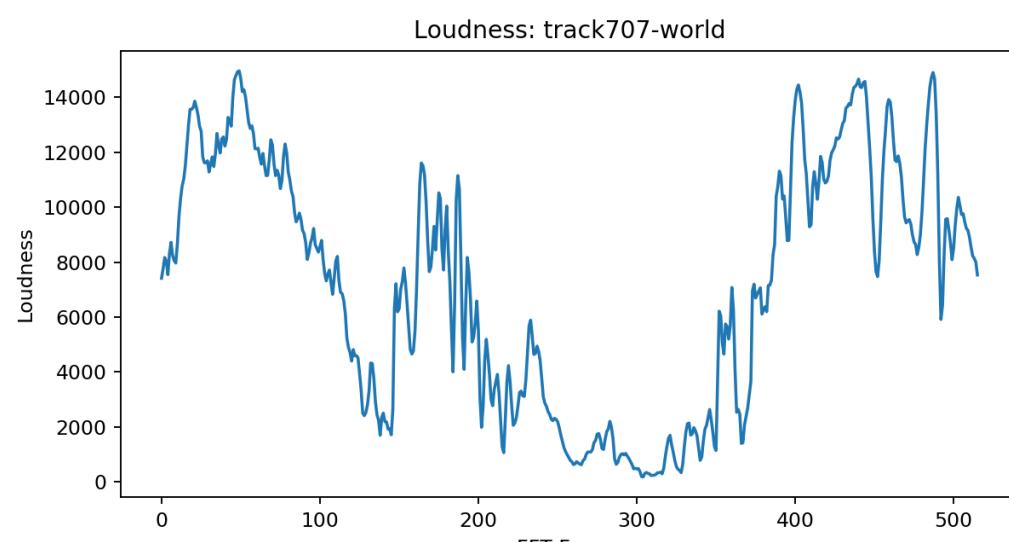


Figure 1: Loudness vs. frame



(e) 'track547-rock.wav'



(f) 'track707-world.wav'

Figure 1: Loudness vs. frame

Comments While each of the plots in figure 1 has its own characteristics, loudness is not the best tool for characterizing genres. Figure 1b shows the uniformity of the track over time and contrasts with figure 1f as well as figure 1a, however it would be hard to distinguish between figure 1b and figure 1d. However, these plots do give a good sense of the progression of each song during the sample considered.

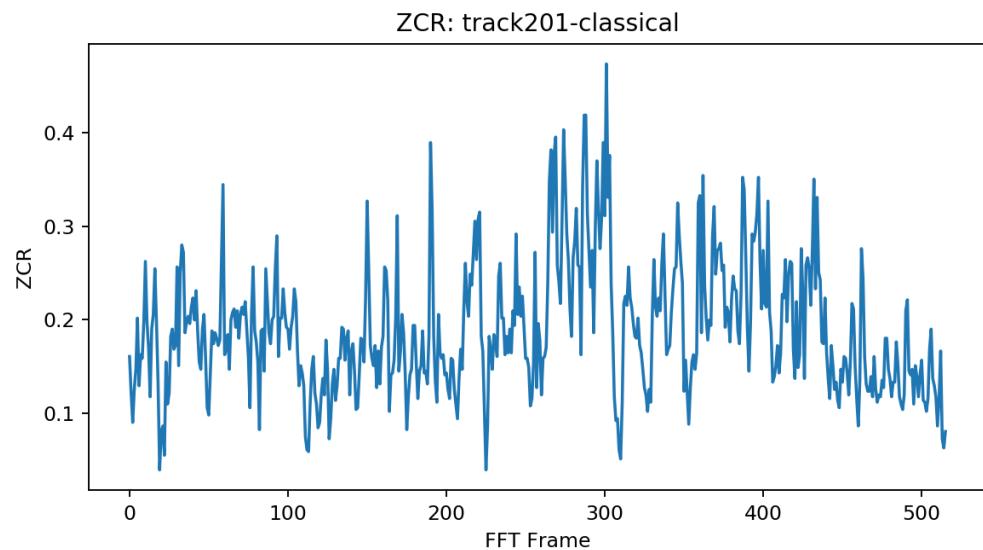
1.2 Zero Crossing Rate

The zero crossing rate (ZCR) is the average number of times the audio signal crosses the zero amplitude line per unit time. It is related to pitch height and correlated to the noise in the signal. For this lab, ZCR is defined as:

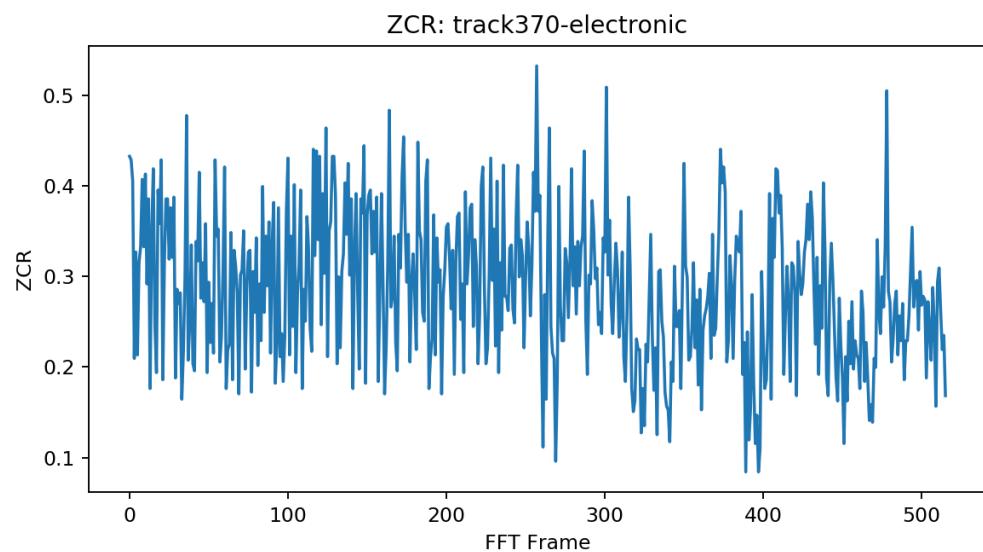
$$ZCR(n) = \frac{1}{N-1} \sum_{m=1}^{N-1} \frac{1}{2} |sgn(x(nN+m)) - sgn(x(nN+m-1))|. \quad (2)$$

See section 5.3.2 for Python implementation.

Results Again, we display the output of computed ZCR for one track of each genre in the figures below.

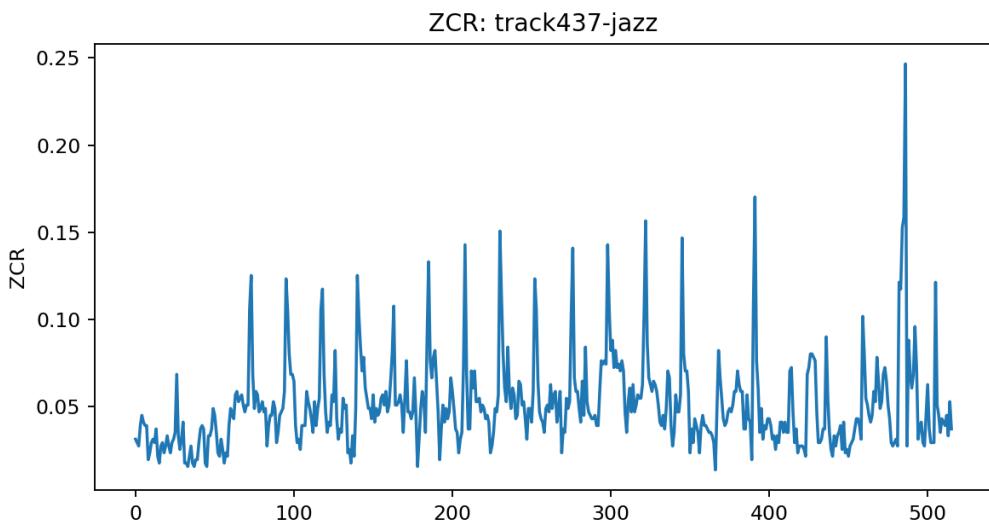


(a) 'track201-classical.wav'

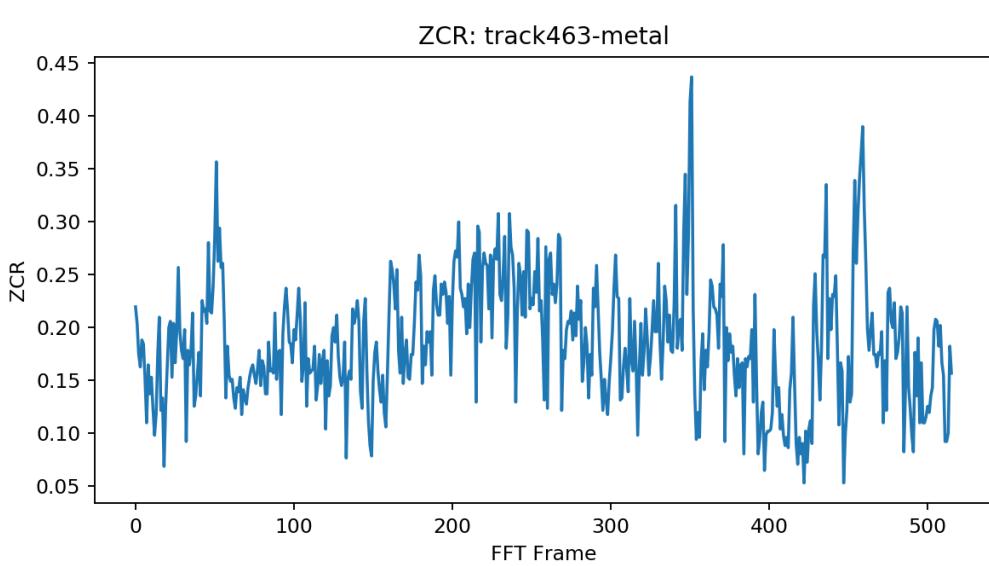


(b) 'track370-electronic.wav'

Figure 2: ZCR vs. frame



(c) 'track437-jazz.wav'



(d) 'track463-metal.wav'

Figure 2: ZCR vs. frame

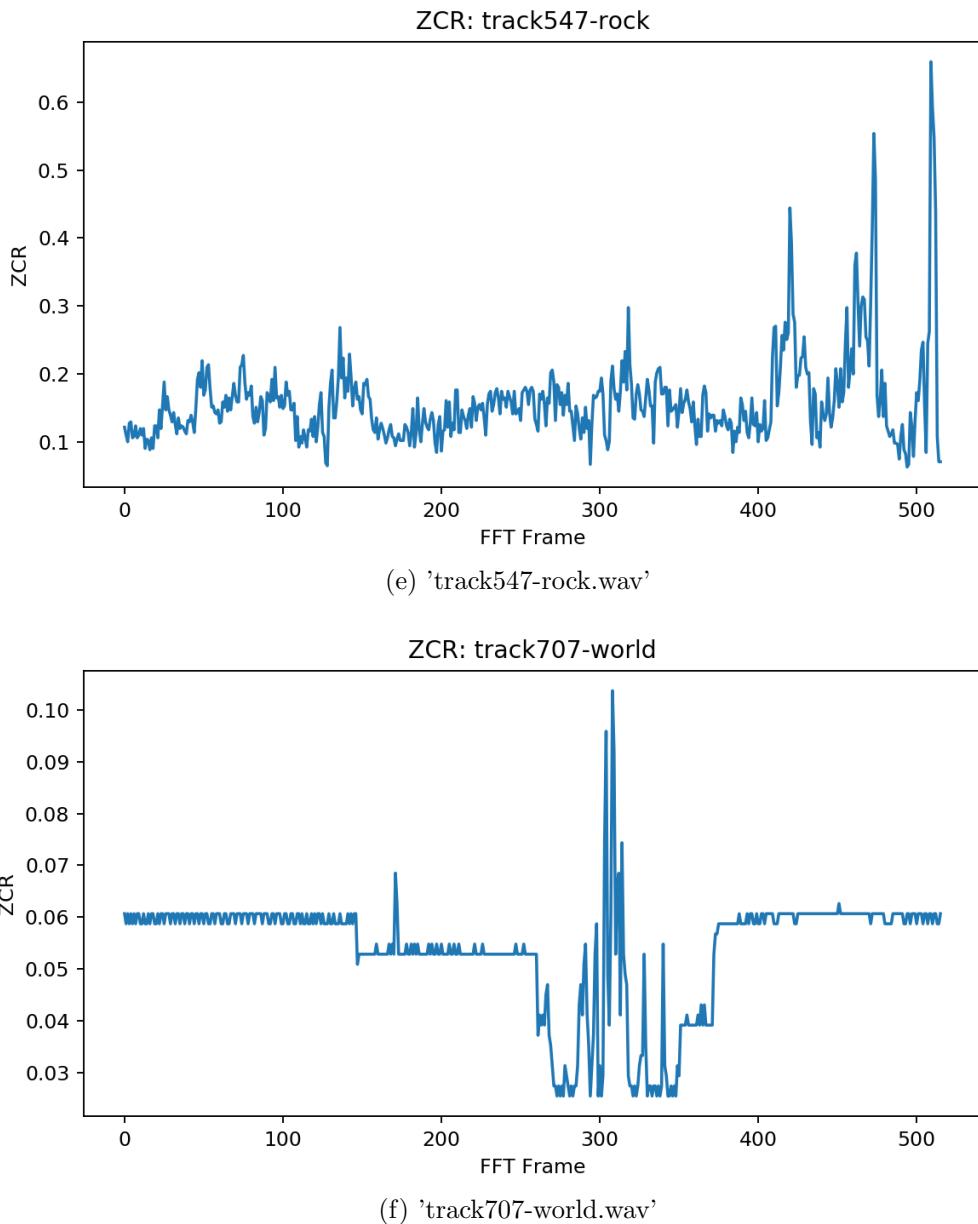


Figure 2: ZCR vs. frame

Comments From the plots in figure 2 we can definitely see some variability between genres. For example, the world track in 2f is easily distinguishable, however most of the other tracks are similar enough that this is also not likely a good measurement for classifying genre. This can be supported by the similarity between jazz (2c) and rock (2e).

2 Spectral (Frequency) Analysis

Sound, by nature is made up of many different frequencies or notes of different pitch. Thus it is natural to analyze audio in the spectral (frequency) domain.

Method As discussed in the introduction, we will accomplish our spectral analysis using a STFT, which decomposes short-time frames of a signal into discrete frequency bins, maintaining some time-domain relevance within the spectral domain. Then, we will perform some statistical analysis on the signal in an attempt to characterize some songs.

Rather than directly splitting each of our $N = 512$ frames, we will convolve (a fancy term for multiplication in the frequency domain) a taper window, w , with the signal and take the Fourier transform (FT). Since the signal is real, and we are concerned only with the magnitude, we will compute the FT of only half of the frequency spectrum to arrive at $N/2 = 256$ frequency bins for each of our $N = 512$ frames. Finally, for the purposes of this analysis, we are interested in the power spectrum of the original song, so we take the magnitude squared (i.e. $|X_n(k)|^2$, as a function of the frame index n and frequency index k). See section 5.4 for Python implementation.

DTFT Window Derivation To give some motivation for the taper window, w , we will derive the theoretical discrete-time Fourier transform (DTFT) for a given signal $x[n]$ where,

$$x[n] = 1, \quad -N/2 \leq n \leq N/2 \quad (3)$$

$$x[n] = 0, \quad \text{else} \quad (4)$$

Recall the definition of the DTFT.

$$X(\omega) = \sum_{n=-\infty}^{\infty} x[n]e^{-j\omega n}. \quad (5)$$

Since $x[n] \neq 0$ over only a subset of values in the DTFT range, equation (5) becomes,

$$\begin{aligned} X(\omega) &= \sum_{n=-N/2}^{N/2} x[n]e^{-j\omega n} \\ &= \sum_{n=-N/2}^{N/2} (e^{-j\omega})^n && \text{which we can rewrite as (a geometric series)} \\ &= e^{j\omega N/2} \frac{1 - (e^{-j\omega})^{N+1}}{1 - e^{-j\omega}} && \text{factoring,} \\ &= e^{j\omega N/2} \frac{e^{-j\omega(N+1)/2} (e^{j\omega(N+1)/2} - e^{-j\omega(N+1)/2})}{e^{-j\omega/2} (e^{j\omega/2} - e^{-j\omega/2})} \\ &= \frac{e^{j\omega(N+1)/2} - e^{-j\omega(N+1)/2}}{e^{j\omega/2} - e^{-j\omega/2}} && \text{and by Euler's formula,} \\ X(\omega) &= \frac{\sin(\frac{\omega(N+1)}{2})}{\sin(\frac{\omega}{2})} \end{aligned} \quad (6)$$

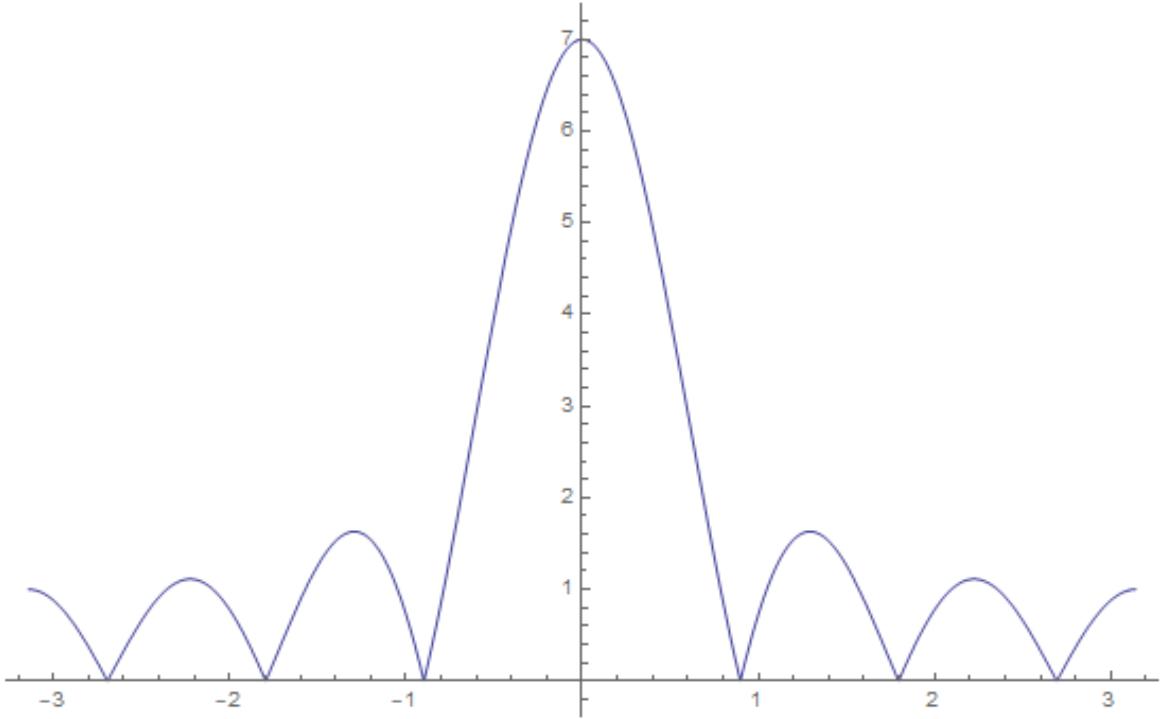


Figure 3: Plot of (6) with $\omega = [-\pi, \pi]$ ($N = 6$).

Equation (6), the magnitude of which is plotted in figure 3, is easily recognized as the *sinc* function. Notice the relatively large amplitude side-lobes. If we consider constant total energy of this function when used as a window, the narrow main-lobe will extract each frequency very selectively, with the side-lobes producing frequency leakage. In contrast, if we aim to extract frequencies with minimal leakage, we could choose a window with a wider main-lobe. In that case, the energy would be more concentrated at the central frequency, albeit less selectively. In this lab, we will focus on some windows that are relatively good approximations to a *sinc* window, but with no side-lobes, so the frequency leakage is essentially eliminated. The tradeoff is that we produce a low-passed version of the original signal at the output of the window operation.

Windowing a Signal The process of multiplying a signal $x[n]$ by a window $w[n]$ is relatively straightforward. We consider the operation in the time-domain and derive the frequency-domain equivalence.

$$y[n] = x[n]w[n] \quad (7)$$

Recall that multiplication in time is a convolution in frequency, such that

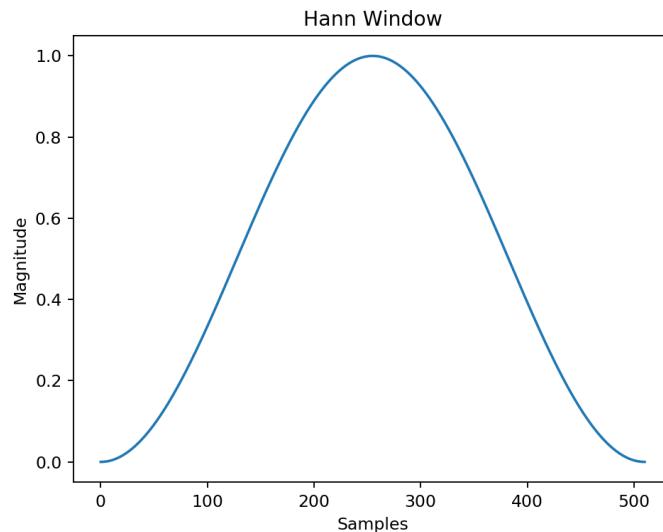
$$\begin{aligned} Y(\omega) &= X(\omega) * W(\omega) \\ &= \sum_{k=-\infty}^{+\infty} W(\omega)X(\omega - k) \end{aligned} \quad (8)$$

$$= \sum_{k=-\pi}^{\pi} \frac{\sin(\frac{\omega(N+1)}{2})}{\sin(\frac{\omega}{2})} X(\omega - k) \quad (9)$$

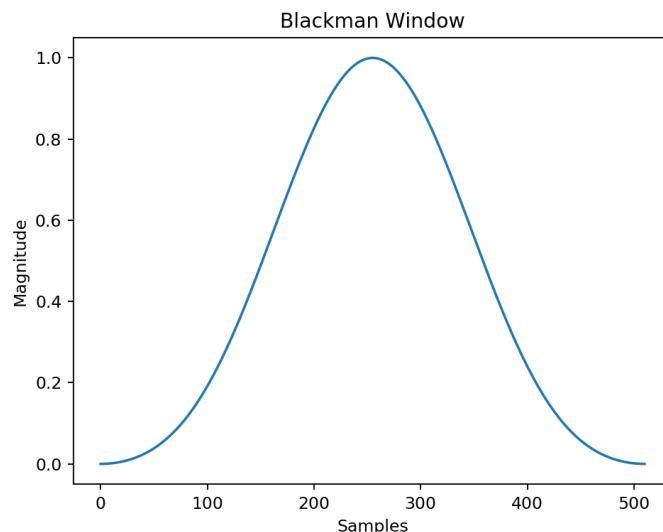
where (8) is the definition of convolution[1] and (9) is the output of the window operation in frequency. We will use this over a sequence of 'frames' of a DTFT algorithm in Python to generate a spectrogram of a signal.

2.1 Spectrogram

As promised, we will be implementing a STFT. In Python, we will use the `scipy.signal.spectrogram` function to automate this. However, we first need a window. For this, we will use the library function, `scipy.signal.get_window` which allows us to quickly generate taper windows of the same size as our STFT frames. There are many choices of windows, but we will look at the 'Hann' and 'Blackman' windows. See figure 4.



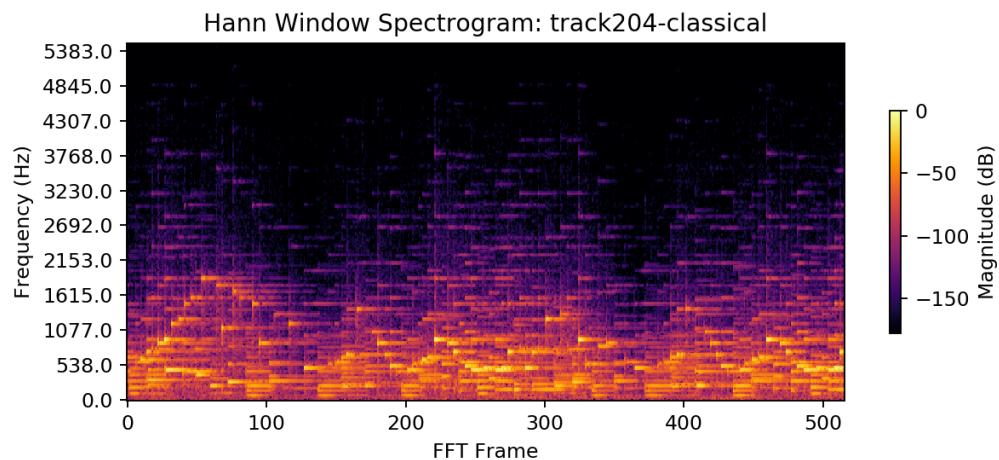
(a) Hann window



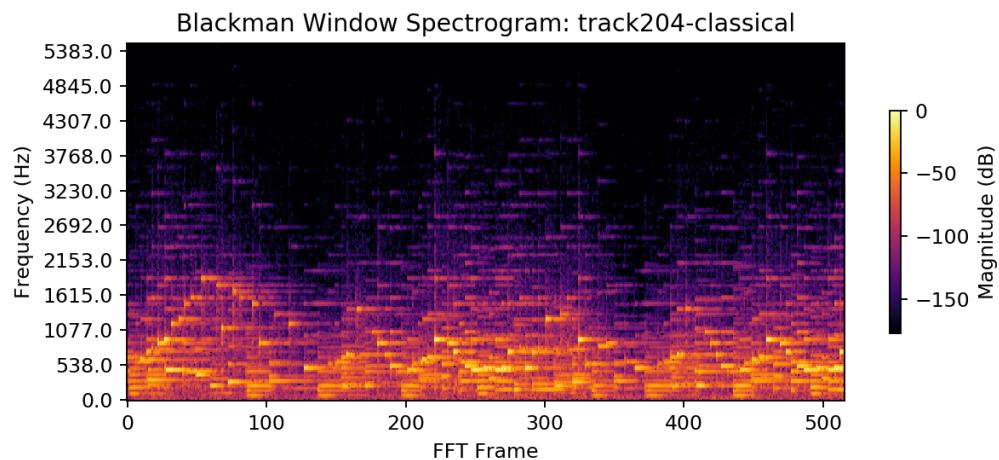
(b) Blackman window

Figure 4: Windows

As you can see, the Hann window (4a) tapers a little less steeply than the Blackman window (4b).

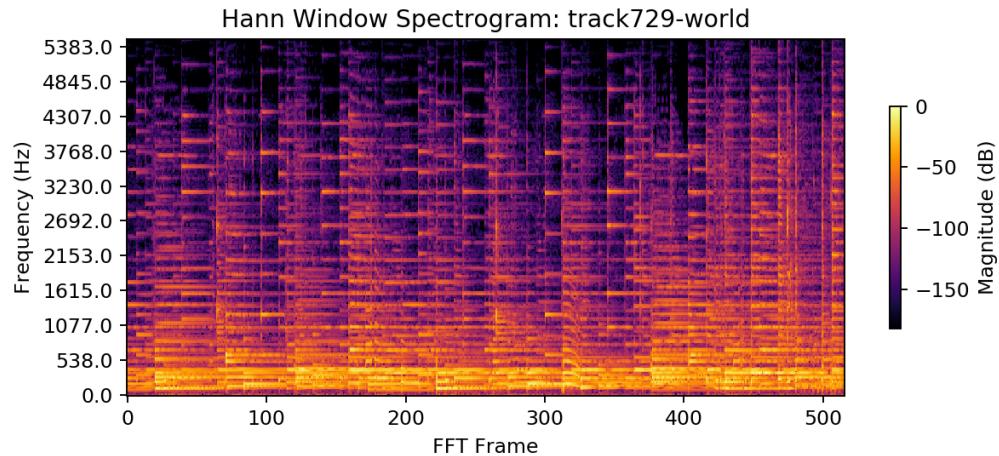


(a) Spectrogram of 'track204-classical' using a Hann window

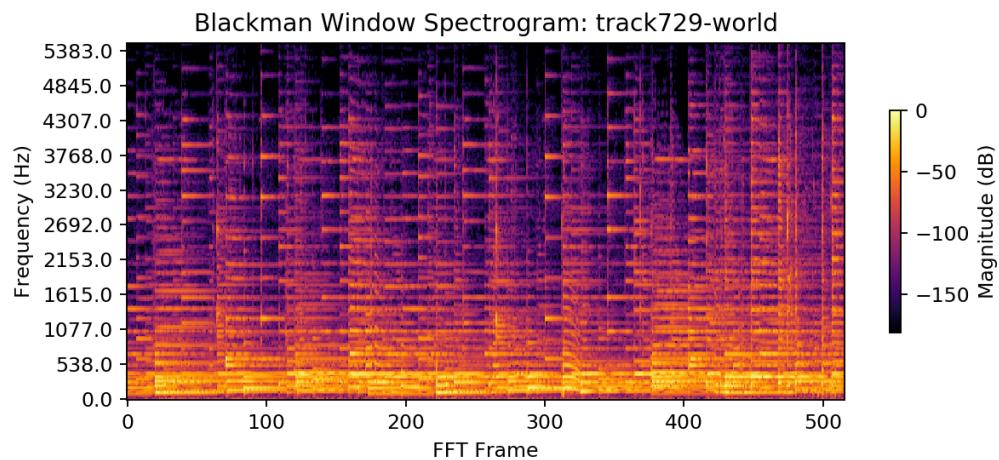


(b) Spectrogram of 'track204-classical' using a Blackman window

Figure 5: Hann vs. Blackman windows for a classical music track



(a) Spectrogram of 'track729-world' using a Hann window



(b) Spectrogram of 'track729-world' using a Blackman window

Figure 6: Hann vs. Blackman windows for a classical music track

See section 5.4.1 for implementation in Python.

Comments If you look closely at figures 7a and 7b, you can see that there is more spectral bleeding in the Blackman window, however the magnitude of center frequencies is relatively higher. We will proceed with analysis using only the Blackman window.

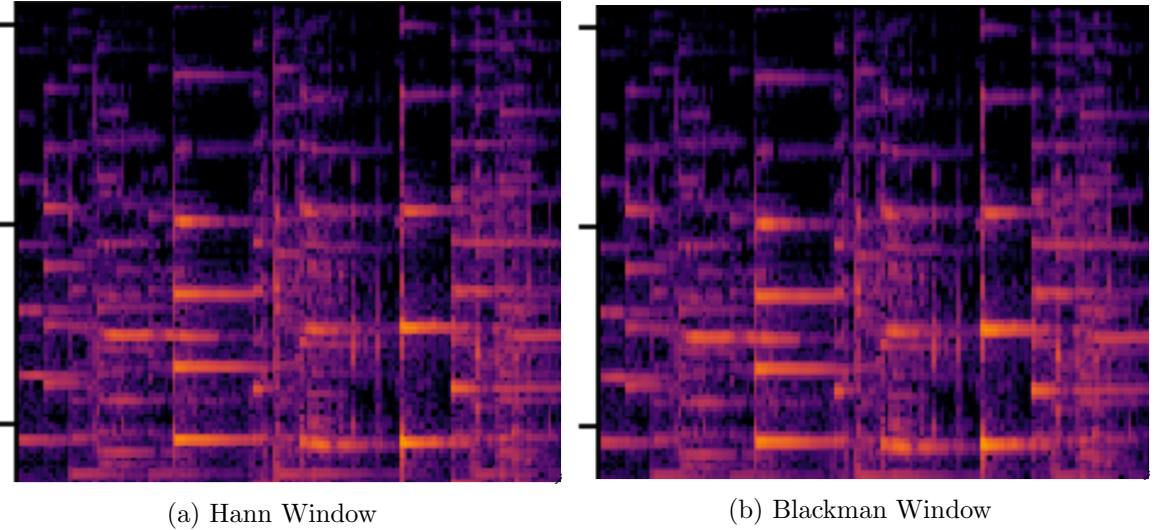


Figure 7: A closer look at Hann and Blackman windows.

2.2 Spectral Centroid and Spread

We wish to find the 'center of mass' of each frame in a spectrogram and will call this the spectral centroid. The centroid can be used to quantify sound sharpness or brightness. Additionally, we would like to find the spectral spread, or width of the spectrum around the centroid. Thus, we can then compare tone-like and noise-like sounds. For these concepts, we will treat the normalized magnitude of a spectral coefficient as if it were a 'probability' of that particular frequency. Then, for frame n , we have the 'probability' of frequency k

$$P_n(k) = \frac{|X_n(k)|}{\sum_{l=0}^K |X_n(l)|} \quad (10)$$

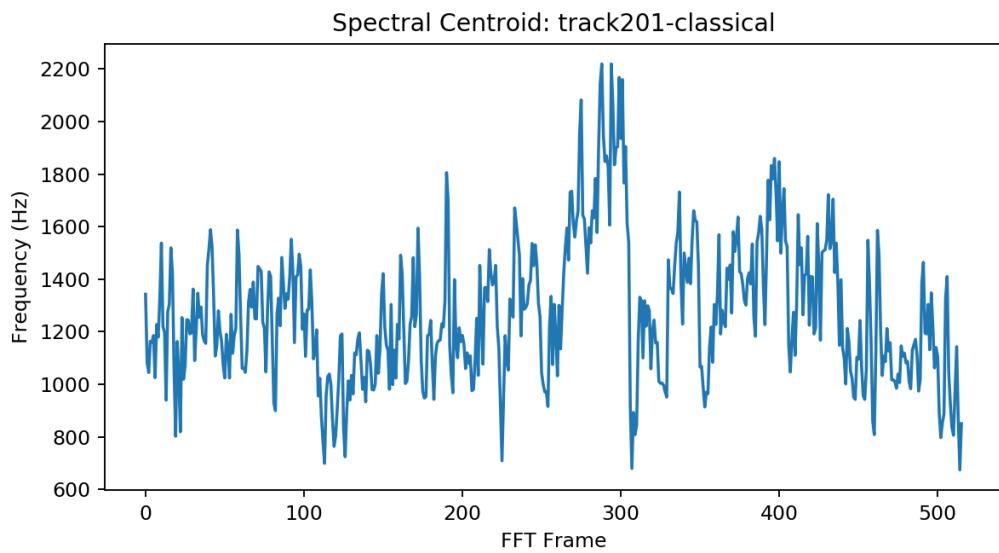
Which we can then use to define the spectral centroid as

$$\mu_n = \sum_{k=0}^K k P_n(k) \quad (11)$$

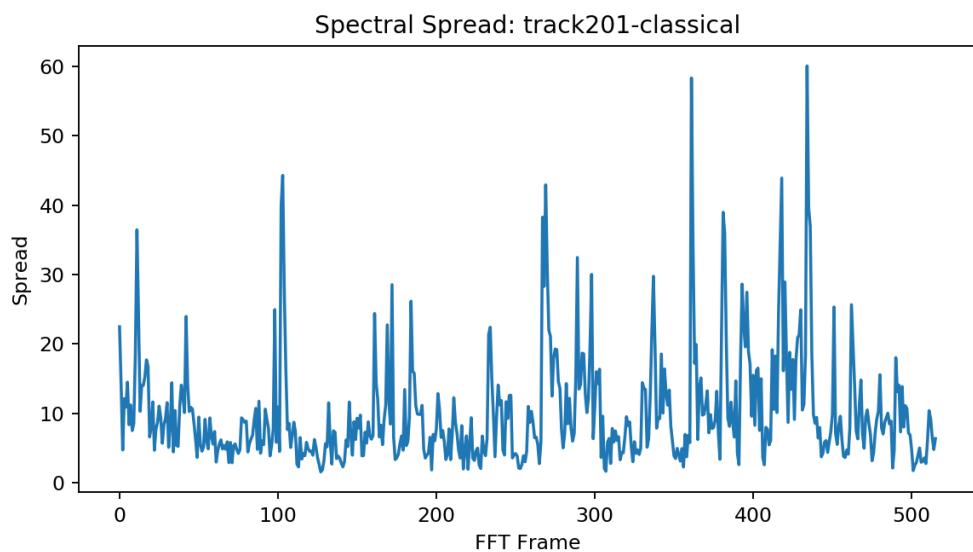
And the spectral spread for frame n is the standard deviation given by

$$\sigma_n = \sqrt{\sum_{k=0}^K [k - \mu_n]^2 P_n(k)} \quad (12)$$

The spectral centroid and spread were calculated for a classical, jazz, and metal song and a time-domain plot is shown in figure 8. See section 5.4.2 for Python implementation.

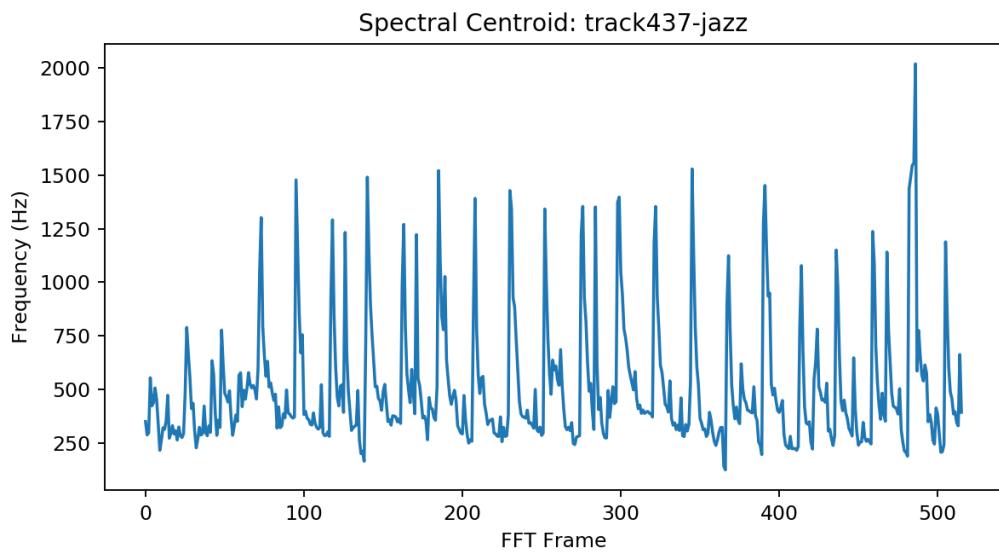


(a) Spectral centroid of 'track201-classical' as a function of frame number.

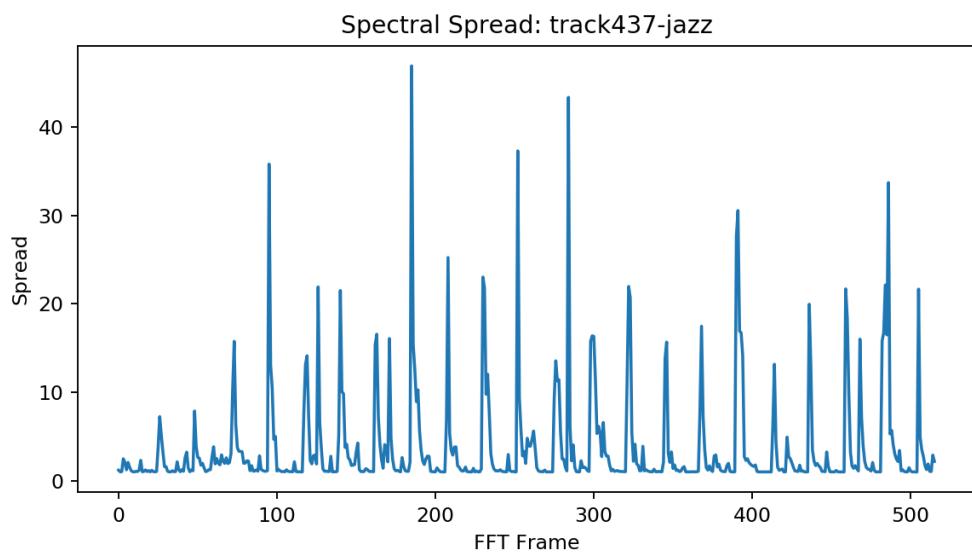


(b) Spectral spread of 'track201-classical' as a function of frame number.

Figure 8: Centroid and spread

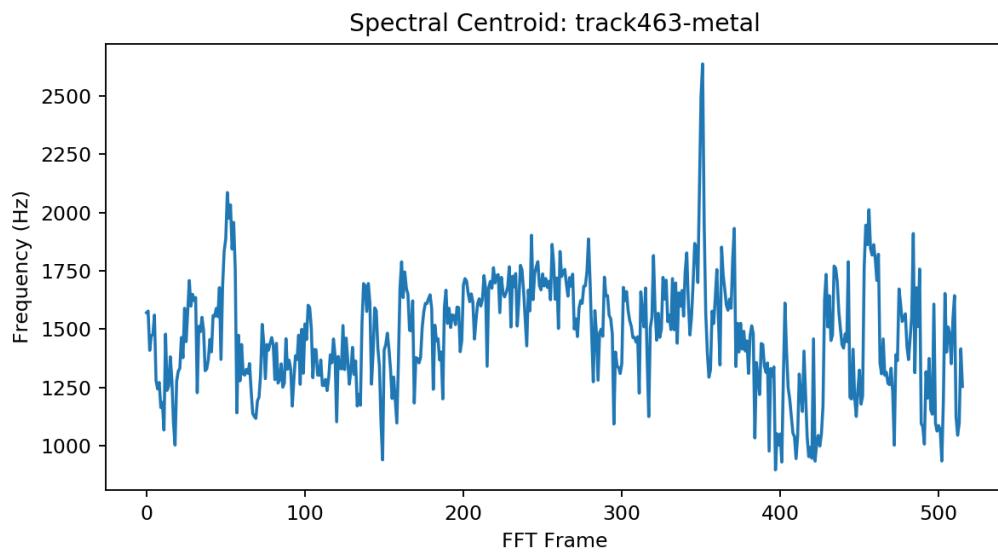


(c) Spectral centroid of 'track437-jazz' as a function of frame number.

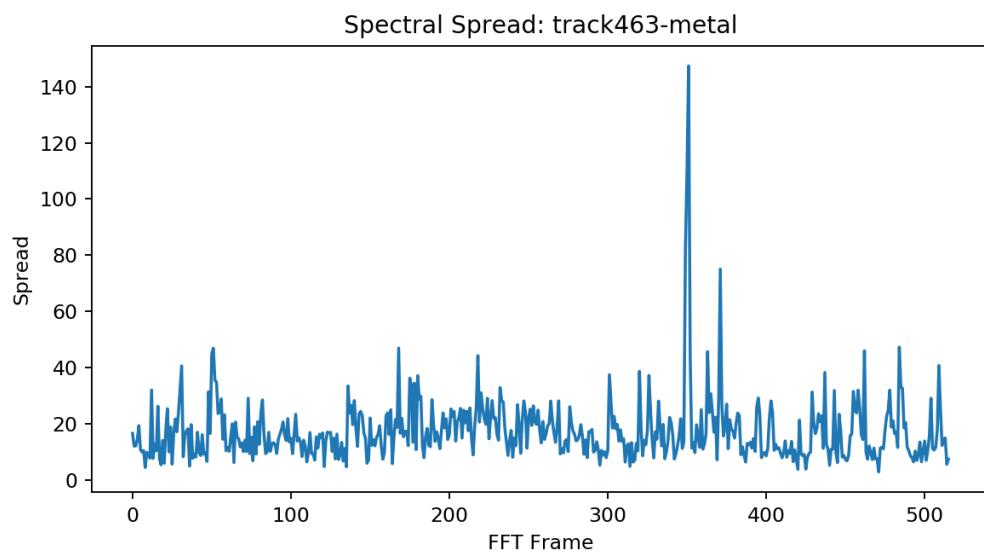


(d) Spectral spread of 'track437-jazz' as a function of frame number.

Figure 8: Centroid and spread



(e) Spectral centroid of 'track463-metal' as a function of frame number.



(f) Spectral spread of 'track463-metal' as a function of frame number.

Figure 8: Centroid and spread

Comments The centroids of this classical track, shown in figure 8a the most statistically significant frequency changes from frame to frame, sometimes drastically. This is most likely due to the violin’s sharp transitions between tones. Then in the spread (figure 8b), we can see that there is relatively high spread for many frames, indicating the high tones in the song and their inherent noise (potentially from under-sampling).

Then, we see in the analysis of a jazz song (figures 8c and 8d) There is a contrast between the percussion and electric keyboard, which is apparent in the consistently alternating peaks and valleys. The spread shows that many of the peaks are also high-frequency and noisy.

In the metal track, we see a relatively uniform centroid (figure 8e) and spread (figure 8f) across frames. Notice also the different scales on the vertical axis of these plots, indicating that the metal track is fairly noisy compared to the other two songs, although it appears to have smaller magnitude on first glance. We see a large spike around frame 350 corresponding to a high-pitch yell from the vocalist. In the centroid plot, this can be seen as a spike in frequency with an associated spread due to the timbre of his voice and noise in the signal.

Comparing these between genres, notice that the classical and jazz have some similarities, which would make sense based on the similarity of genres, however we can see a contrast in the centroids, due to the electronic composition of the jazz track as opposed to the natural, dissonant composition of the classical. Furthermore, the small spread shown in figure 8d indicates that notes were generated electronically. The metal track stands out in a couple regards. Figure 8e is vastly different and shows the dense composition in the track, with the centroid in each frame nearing the middle of the frequency spectrum and indicating that the range of frequencies is large at any given point in the song.

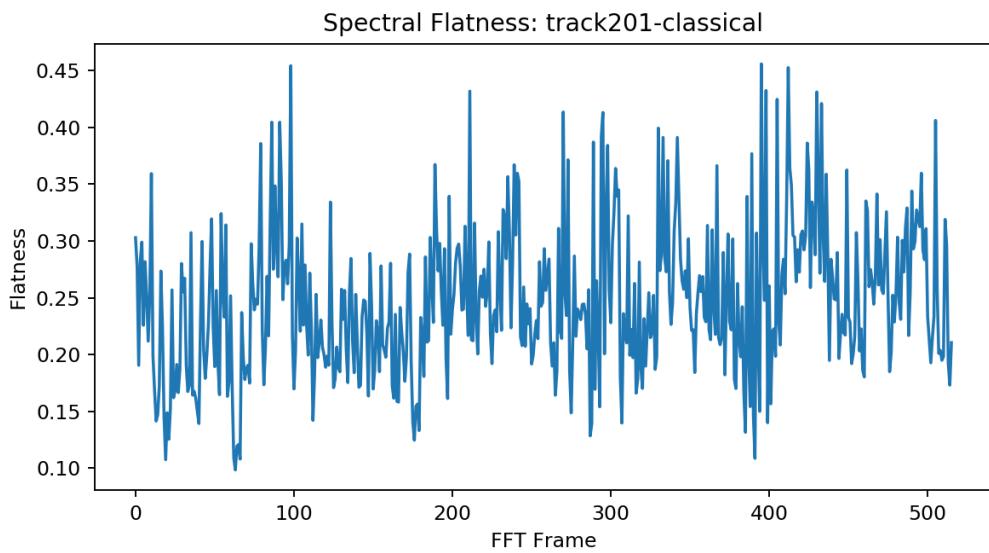
2.3 Spectral Flatness

Spectral flatness is the ratio between the geometric and arithmetic means of the magnitude of the Fourier transform,

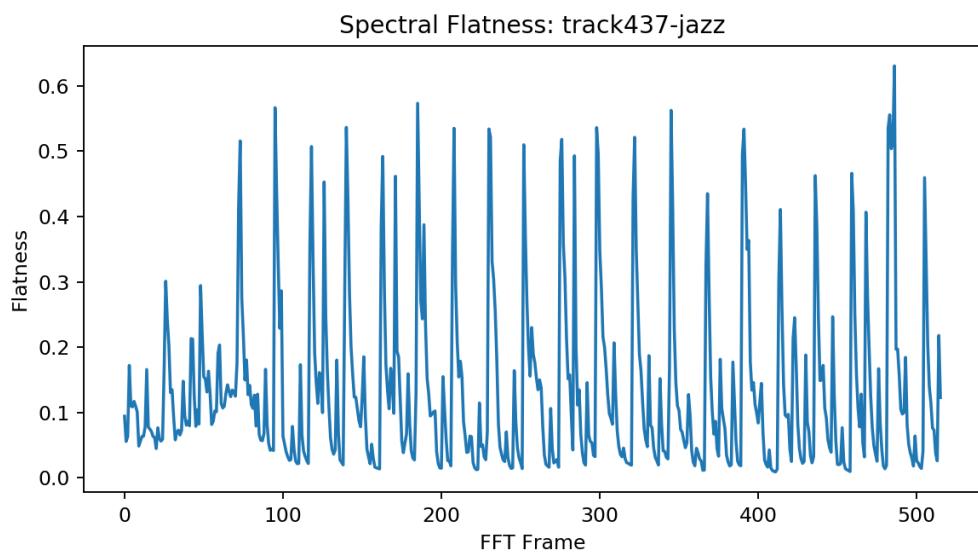
$$SF(n) = \frac{\left(\prod_{k=0}^K |X_n(k)|\right)^{1/K}}{\frac{1}{K} \sum_{k=0}^K |X_n(k)|} \quad (13)$$

A very small flatness corresponds to the presence of tonal components while a flatness equal to one corresponds to a very noisy signal. Thus, flatness is a measure of the noisiness of the spectrum.

Method This implementation was done using some low-level `scipy` and `numpy` methods, specifically `scipy.stats.mstats.gmean` and `numpy.mean` to calculate the geometric and arithmetic means, respectively. Again, this analysis was performed on classical, jazz, and metal samples. For the full Python implementation, see section 5.4.3.

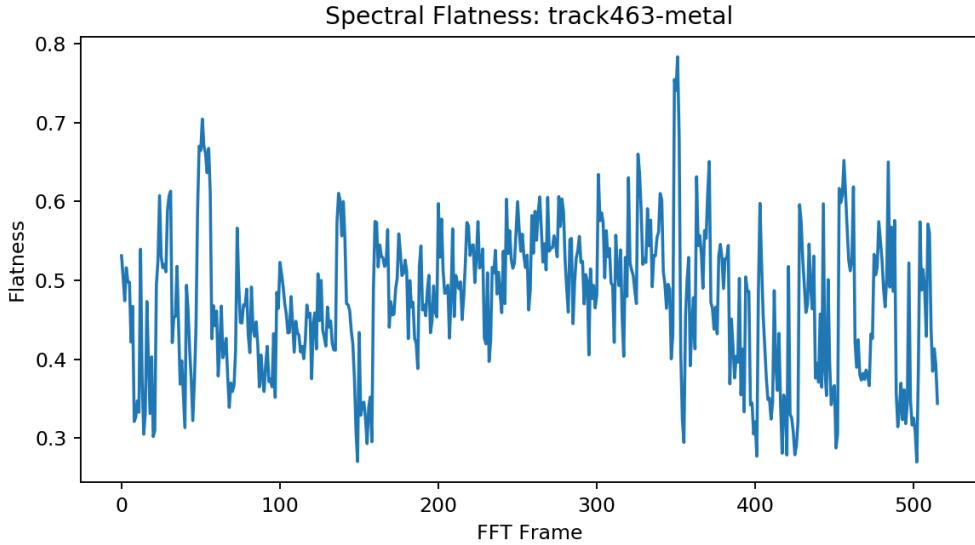


(a) Spectral flatness of 'track201-classical' as a function of frame number.



(b) Spectral flatness of 'track437-jazz' as a function of frame number.

Figure 9: Flatness



(c) Spectral flatness of 'track463-metal' as a function of frame number.

Figure 9: Flatness

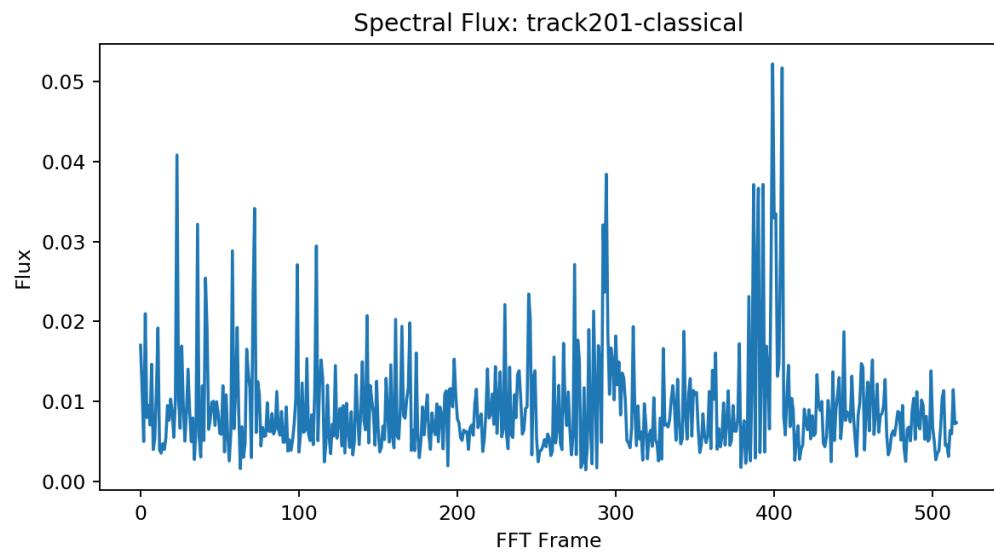
Comments Notice the magnitude of flatness across the entire metal sample (figure 9c) is almost double that of classical (figure 9a). This is a direct indication that the metal track is much noisier. Then comparing the flatness curve of this jazz sample, we see the same alternating magnitudes as in the centroid (figure 8c) and spread (figure 8d) curves, again indicating the noise inherent in the high-pitch components of this sample.

2.4 Spectral Flux

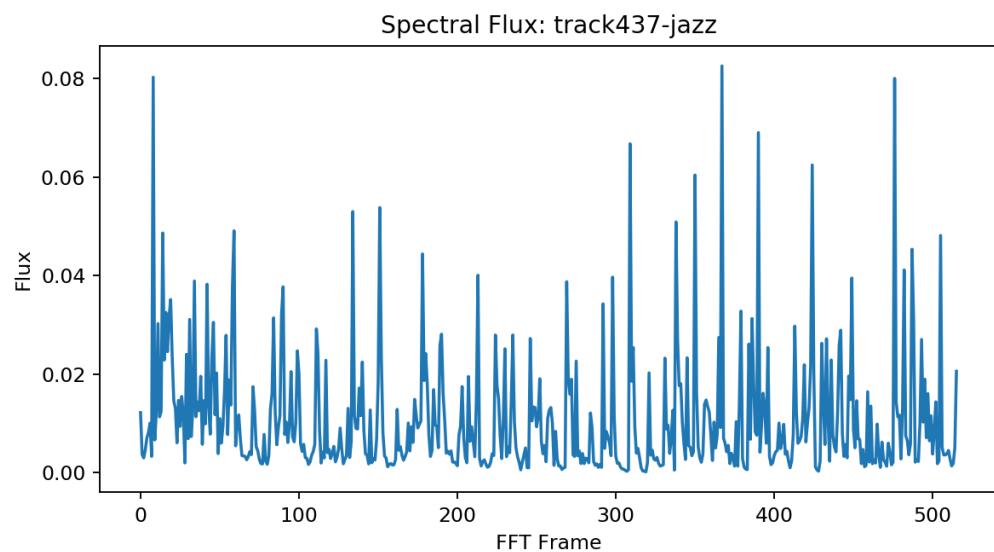
The spectral flux is a global measure of the spectral changes between two adjacent frames, $n - 1$ and n ,

$$F_n = \sum_{k=0}^K (P_n(k) - P_{n-1}(k))^2 \quad (14)$$

Where $P_n(k)$ is the normalized frequency distribution for frame n , given by 10. Applying 14 to the same three samples, the plots of figure 10 were produced. See section 5.4.4 for Python implementation.

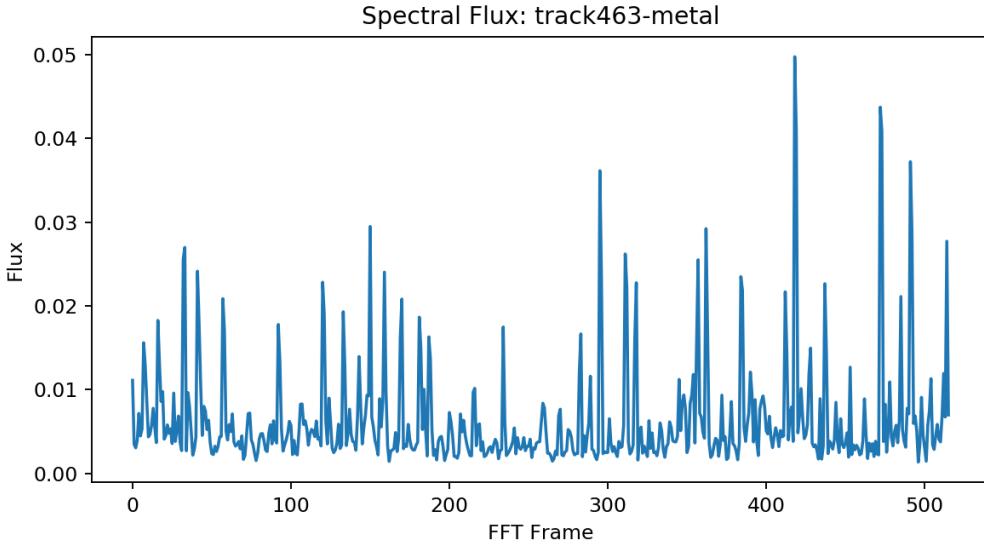


(a) Spectral flux of 'track201-classical' as a function of frame number.



(b) Spectral flux of 'track437-jazz' as a function of frame number.

Figure 10: Flux



(c) Spectral flux of 'track463-metal' as a function of frame number.

Figure 10: Flux

Comments Considering the similarity of figures 10b and 10c, spectral flux is not a good indication of genre.

3 Psychoacoustics

Background Psychoacoustics involves the study of the human auditory system and our perception of audio. In this lab, we will focus on some spectral content that can be defined mathematically. The human auditory system is sensitive to a range of frequencies with a logarithmic relationship. One very common model of human sound perception is called the mel/Bark scale.

3.1 The mel/Bark Scale

The Bark scale is defined as

$$z = 7 \operatorname{arcsinh}(f/650) = 7 \log \left(f/650 + \sqrt{1 + (f/650)^2} \right),$$

where f is measured in Hz. However, in this lab, we will use a modified definition,

$$m = 1127.01048 * \log(1 + f/700). \quad (15)$$

3.2 The Cochlear Filterbank

Additionally, the human auditory system behaves like a sequence of filters (filterbank) with overlapping frequency responses. Perception of pitch can be quantified using the total energy at the output of each filter, summing the spectral energy that falls into one critical band (the frequency band within which a second tone will interfere with the first). Our model is simple, with $N_B = 40$ logarithmically spaced triangle filters centered at frequencies Ω_p , which are implicitly defined as

$$\text{mel}_p = 1127.01048 * \log(1 + \Omega_p/700). \quad (16)$$

The set of frequencies chosen to correspond to each Ω_p is chosen to be equally spaced on the mel scale. Letting the indexing of the center frequencies of the filters start with $p = 1$, we define

$$\text{mel}_p = p \frac{\text{mel}_{\max} - \text{mel}_{\min}}{N_B + 1} + \text{mel}_{\min} \quad (17)$$

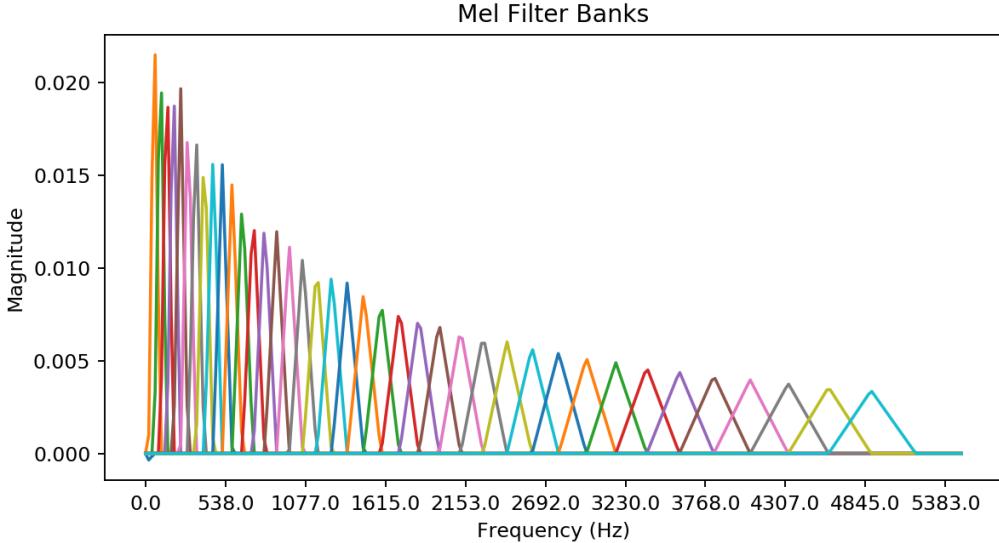


Figure 11: A collection of $N_B = 40$ filterbanks to be used in the calculation of mfcc coefficients when convolved with an audio sample.

where N_B is the number of filters in the filterbank and

$$mel_{min} = 1127.01048 * \log(1 + 20/700),$$

$$mel_{max} = 1127.01048 * \log(1 + 0.5 * f_s/700),$$

taking f_s to be the sampling frequency of the audio. This definition gives mel_{max} to be defined as the highest frequency preserved in the audio file and mel_{min} to be 20Hz (a good baseline for the limit of human hearing). We then define H_p , the hat filter corresponding to each filter to be centered around Ω_p :

$$H_p(f) = \begin{cases} \frac{2}{\Omega_{p+1} - \Omega_{p-1}} \frac{f - \Omega_{p-1}}{\Omega_p - \Omega_{p-1}} & \text{if } f \in [\Omega_{p-1}, \Omega_p), \\ \frac{2}{\Omega_{p+1} - \Omega_{p-1}} \frac{\Omega_{p+1} - f}{\Omega_{p+1} - \Omega_p} & \text{if } f \in [\Omega_p, \Omega_{p+1}). \end{cases} \quad (18)$$

Each triangular filter in 18 is normalized such that the integral of each filter is 1 and the filters overlap so that the frequency where H_p is maximum is at the starting frequency of the next filter, H_{p+1} and ending frequency of H_{p-1} .

Finally, the implementation of the above derivation produces the plot shown in figure 11.

Comments The triangle filters in figure 11 appear to decrease and increase on the left side of the x-axis, however this is an anomaly due to the quantization of values used to generate each filter. This can be seen in the offset peaks in a few of the filters. In reality, each filter maintains a unity integral over the range N_B . Python implementation for this can be found in section 5.5.

4 Conclusion

Spectral analysis of audio is definitively more useful than time-domain analysis, however this lab showed just how difficult it still can be to differentiate genres analytically. I think this was a good introduction to some of the useful statistical analysis techniques that can be used to programmatically sort signals based on characteristic features. Additionally, the script to run this analysis was relatively slow. With the digital age we live in, delay from input to output is very important, so some optimization would most likely help to make this sort of analysis even more useful. Next, I would like to apply this analysis to some of my favorite songs and see what I can gain from it. Eventually, I envision using this type of functionality to create a highly correlated audio visualizer or an advanced tool to equalize concert halls.

5 Appendix

5.1 Appendix A: Main

```
1 """
2 Created on Wed Jan 16 15:24:09 2019
3
4 @author: Andrew Teta
5 """
6
7 import scipy.io.wavfile
8 import scipy.signal
9 import scipy.stats
10 import numpy as np
11 import math
12 import matplotlib.pyplot as plt
13 from tkinter import filedialog
14
15 # UI dialog to select files -> selection of multiple files will run all functions for
16 # each file
16 files = filedialog.askopenfilenames()
17 # Loop over files selected
18 for f in files:
19     filename = str.split(f, '/')[-1]
20     filename = str.split(filename, '.') [0]
21     filepath = f
22
23 # Extract 24s sample from center of 60s clip
24 fs,songSample = extractSample(filepath,24,60)
25
26 # Calculate loudness (energy function returned just in case)
27 E,sigma,Nloudness = loudness(songSample,512)
28
29 # Save figures
30 plt.figure(figsize=(8,4),dpi=170)
31 plt.plot(range(Nloudness),sigma)
32 plt.xlabel('FFT Frame')
33 plt.ylabel('Loudness')
34 plt.title('Loudness: '+filename)
35 plt.savefig('figs/loudness_'+filename)
36 plt.close()
37
38 # Calculate zero-crossing-rate
39 Z,Nzcr = zcr(songSample,512)
40
41 plt.figure(figsize=(8,4),dpi=170)
42 plt.plot(range(Nzcr),Z)
43 plt.xlabel('FFT Frame')
44 plt.ylabel('ZCR')
45 plt.title('ZCR: '+filename)
46 plt.savefig('figs/zcr_'+filename)
47 plt.close()
48
49 # Calculate spectrograms for Hann and Blackman windows
50 spectroHann,spectroBlack,freqAxis,timeAxis = spectrogram(songSample,fs,512)
51 spectroHann = np.transpose(spectroHann)
52 spectroBlack = np.transpose(spectroBlack)
53
54 mel = mfcc(spectroBlack,freqAxis,1,fs,40)
55
56 # Generate frequencies axis labels
57 tick = np.arange(0,len(freqAxis),25)
58 label = np.around(tick*fs/512)
59
60 # save figures
61 plt.figure(figsize=(8,4),dpi=170)
```

```

62 plt.imshow(spectroHann, cmap='inferno')
63 plt.gca().invert_yaxis()
64 plt.colorbar(shrink=0.5).set_label('Magnitude (dB)')
65 plt.title('Hann Window Spectrogram: '+filename)
66 plt.xlabel('FFT Frame')
67 plt.ylabel('Frequency (Hz)')
68 plt.yticks(tick,label)
69 plt.savefig('figs/powerHann_'+filename)
70 plt.close()
71
72 plt.figure(figsize=(8,4),dpi=170)
73 plt.imshow(spectroBlack, cmap='inferno')
74 plt.gca().invert_yaxis()
75 plt.colorbar(shrink=0.5).set_label('Magnitude (dB)')
76 plt.title('Blackman Window Spectrogram: '+filename)
77 plt.xlabel('FFT Frame')
78 plt.ylabel('Frequency (Hz)')
79 plt.yticks(tick,label)
80 plt.savefig('figs/powerBlack_'+filename)
81 plt.close()
82
83 # calculate statistical vectors, centroid and spread
84 centroid,spread,P,nFrames = spectralCentroid(songSample,fs,512)
85
86 # save figures
87 plt.figure(figsize=(8,4),dpi=170)
88 plt.plot(centroid)
89 plt.title('Spectral Centroid: '+filename)
90 plt.ylabel('Frequency (Hz)')
91 plt.xlabel('FFT Frame')
92 plt.savefig('figs/centroid_'+filename)
93 plt.close()
94
95 plt.figure(figsize=(8,4),dpi=170)
96 plt.plot(spread)
97 plt.title('Spectral Spread: '+filename)
98 plt.ylabel('Spread')
99 plt.xlabel('FFT Frame')
100 plt.savefig('figs/spread_'+filename)
101 plt.close()
102
103 # calculate spectral flatness
104 flat,nFrames = flatness(songSample,fs,512)
105
106 # save figure
107 plt.figure(figsize=(8,4),dpi=170)
108 plt.plot(flat)
109 plt.ylabel('Flatness')
110 plt.xlabel('FFT Frame')
111 plt.title('Spectral Flatness: '+filename)
112 plt.savefig('figs/flatness_'+filename)
113 plt.close()
114
115 sflux = flux(P)
116
117 plt.figure(figsize=(8,4),dpi=170)
118 plt.plot(sflux)
119 plt.ylabel('Flux')
120 plt.xlabel('FFT Frame')
121 plt.title('Spectral Flux: '+filename)
122 plt.savefig('figs/flux_'+filename)
123 plt.close()
124
125 print('done')
126
```

5.2 Appendix B: Sample Extraction

```
1 # =====
2 # Input:
3 #   filepath = string containing absolute path to audio file
4 #   duration = length of sample to extract in seconds
5 #   location = start index of sample given in seconds
6 # Output:
7 #   arg1 = sampling frequency (Hz)
8 #   arg2 = sample of specified duration in seconds
9 # =====
10 def extractSample(filepath, duration, location):
11     fs, data = scipy.io.wavfile.read(filepath)
12     print('Read file: ' + filepath + '\n')
13     Tsample = 1/fs
14     startIndex = int(location/Tsample)
15     endIndex = startIndex + int(duration/Tsample)
16     return fs, data[startIndex:endIndex]
17
```

5.3 Appendix C: Time-domain Analysis

5.3.1 Loudness

```
1 # =====
2 # Input:
3 #   sample = audio clip
4 #   N = frame size
5 # Output:
6 #   arg1 = energy vector
7 #   arg2 = standard deviation vector
8 #   arg3 = num frames (for accurate plotting)
9 # =====
10 def loudness(sample, N):
11     nframes = int(len(sample)/N)
12     E = np.zeros(nframes)
13     sigma = np.zeros(nframes)
14     for n in range(nframes):
15         E[n] = (1/N)*sum(sample[n*N:n*N+(N-1)])
16     for n in range(nframes):
17         sigma[n] = np.sqrt((1/N)*sum((sample[n*N:n*N+(N-1)] - E[n])**2))
18     return E, sigma, nframes
19
```

5.3.2 ZCR

```
1 # =====
2 # Input:
3 #   sample = audio clip
4 #   N = frame size
5 # Output:
6 #   arg1 = zero-crossing-rate
7 #   arg2 = num frames (for accurate plotting)
8 # =====
9 def zcr(sample, N):
10     nframes = int(len(sample)/N)
11     Z = np.zeros(nframes)
12     for n in range(nframes):
13         Z[n] = 1/(N-1) * sum(0.5*abs(np.sign(sample[n*N+1:n*N+(N-1)]) - np.sign(sample[n*N:n*N+(N-2)])))
14     return Z, nframes
15
```

5.4 Appendix D: Frequency-domain Analysis

5.4.1 Spectrogram

```
1 # _____
2 # Input:
3 #   sample = audio clip
4 #   fs = sampling frequency
5 #   N = frame size
6 # Output:
7 #   arg1 = spectrogram using Hann window
8 #   arg2 = spectrogram using Blackman window
9 #   arg3 = vector of axis points in frequency (Hz)
10 #  arg4 = vector of axis points in time (s)
11 #
12 def spectrogram(sample, fs, N):
13     # generate a Hann window
14     wHann = scipy.signal.get_window('hann', N-1, False)
15     # Generate Blackman window
16     wBlack = scipy.signal.get_window('blackman', N-1, False)
17     # Save figure
18     plt.figure(dpi=170)
19     plt.plot(wHann)
20     plt.title('Hann Window')
21     plt.xlabel('Samples')
22     plt.ylabel('Magnitude')
23     plt.savefig('figs/hann')
24     plt.close()
25     # save figure
26     plt.figure(dpi=170)
27     plt.plot(wBlack)
28     plt.title('Blackman Window')
29     plt.xlabel('Samples')
30     plt.ylabel('Magnitude')
31     plt.savefig('figs/blackman')
32     plt.close()
33     # Number of FFT frames
34     nframes = int(len(sample)/N)
35     SHann = np.zeros([nframes, int(N/2)])
36     SBlack = np.zeros([nframes, int(N/2)])
37     time = np.zeros(nframes)
38     # Loop over FFT frames
39     for n in range(nframes):
40         # Calculate a frequency power spectrum for Hann window
41         f, t, sH = np.transpose(scipy.signal.spectrogram(sample[n*N:n*N+(N-1)], fs, wHann,
42 mode='magnitude'))
43         # throw away values smaller than 10^-3 and square values
44         sH = sH**2
45         sH = np.where(sH < 10**-3, 10**-3, sH)
46         # Repeat for Blackman window
47         f, t, sB = np.transpose(scipy.signal.spectrogram(sample[n*N:n*N+(N-1)], fs, wBlack,
48 mode='magnitude'))
49         sB = sB**2
50         sB = np.where(sB < 10**-3, 10**-3, sB)
51         # Reshape a bit
52         SHann[n, :] = np.transpose(sH*fs/N)
53         SBlack[n, :] = np.transpose(sB*fs/N)
54         time[n] = t
55     # convert linear values to dB scale
56     SHann = 20*np.log10(SHann/np.amax(SHann))
57     SBlack = 20*np.log10(SBlack/np.amax(SBlack))
58     return SHann, SBlack, f, time
```

5.4.2 Spectral Centroid and Spread

```
1 # _____
2 # Input:
3 #   sample = audio clip
4 #   fs = sampling frequency
5 #   N = frame size
6 # Output:
7 #   arg1 = vector of size nframes of spectral centroids
8 #   arg2 = vector of size nframes of spectral spreads
9 #   arg3 = length of output vectors = nframes
10 #
11 def spectralCentroid(sample, fs, N):
12     nframes = int(len(sample)/N)
13     SP = np.zeros([nframes, int(N/2)])
14     # Generate Blackman window
15     wBlack = scipy.signal.get_window('blackman', N-1, False)
16     # Generate a spectrogram matrix
17     for n in range(nframes):
18         f, t, sp = np.transpose(scipy.signal.spectrogram(sample[n*N:n*N+(N-1)], fs, wBlack,
19         mode='magnitude'))
20         SP[n, :] = np.transpose(sp)
21         frange = np.shape(sp)[0]
22         P = np.zeros([nframes, frange])
23         mu = np.zeros(nframes)
24         spread = np.zeros(nframes)
25         # Loop over FFT frames
26         for n in range(nframes):
27             # calculate sum of all frequency components in this frame
28             sfreq = sum(SP[n, 0:frange])
29             for k in range(frange):
30                 # calculate the relative 'probability' of each frequency bin
31                 P[n, k] = SP[n, k] / sfreq
32                 # calculate spectral centroid (center of mass) of each frame
33                 mu[n] = mu[n] + (k * fs / N) * P[n, k]
34                 # calculate spectral spread of each frame
35                 spread[n] = np.sqrt(spread[n] + P[n, k] * (k - mu[n]) ** 2)
36     return mu, spread, P, nframes
```

5.4.3 Spectral Flatness

```
1 # _____
2 # Input:
3 #   sample = audio clip
4 #   fs = sampling frequency
5 #   N = frame size
6 # Output:
7 #   arg1 = vector of length nframes of spectral flatness
8 #   arg2 = length of output vectors = nframes
9 #
10 def flatness(sample, fs, N):
11     nframes = int(len(sample)/N)
12     K = int(N/2)
13     SP = np.zeros([nframes, int(N/2)])
14     wBlack = scipy.signal.get_window('blackman', N-1, False)
15     geo = np.zeros(nframes)
16     arith = np.zeros(nframes)
17     SF = np.zeros(nframes)
18     for n in range(nframes):
19         f, t, sp = np.transpose(scipy.signal.spectrogram(sample[n*N:n*N+(N-1)], fs, wBlack,
20 mode='magnitude'))
21         SP[n, :] = np.transpose(sp)
22         geo[n] = scipy.stats.mstats.gmean(SP[n, :])
23         arith[n] = np.mean(SP[n, :])
24     SF = geo / arith
25     return SF, nframes
```

5.4.4 Spectral Flux

```
1 # _____
2 # Input:
3 #   P = probability matrix of dimension [nframes, N/2]
4 # Output:
5 #   arg1 = vector of length nframes of spectral flux
6 #
7 def flux(P):
8     nframes = np.shape(P)[0]
9     frange = np.shape(P)[1]
10    F = np.zeros(nframes)
11    for n in range(nframes):
12        F[n] = sum((P[n, :] - P[n-1, :]) ** 2)
13    return F
14
```

5.5 Psychoacoustics

5.5.1 Mel Filter Banks and Coefficients

```
1 # _____
2 # Input:
3 # SP = audio sample spectrogram
4 # f = list of frequency values with size = range of frequency bins in SP
5 # n = frame index of SP
6 # fs = sampling frequency of audio sample
7 # nBanks = number of mel filter banks to generate
8 # Output:
9 # arg1 = N.B x K (K = len(f)) matrix of filter values for H-p
10 # arg2 = matrix of filter banks, size = nBanks x len(f)
11 #
12 def mfcc(SP,f,n,fs,nBanks):
13     mel_min = 1127.01048*math.log(1 + 20/700)
14     omega_min = (math.e**((mel_min/1127.01048) - 1) * 700
15     mel_max = 1127.01048 * math.log(1 + 0.5*fs/700)
16     omega_max = (math.e**((mel_max/1127.01048) - 1) * 700
17     # Generate frequency array in Hz
18     linfreq = np.arange(np.round(omega_min,0),np.round(omega_max,0)+1,1)
19     # Map frequency array to mel array
20     melfreq = 1127.01048*np.log(1+linfreq/700)
21     # Generate nbanks+2 mel frequencies uniformly spaced between mel_min and mel_max,
22     # inclusive
23     mel = np.linspace(mel_min, mel_max, nBanks+2)
24     # Array of center frequencies of each filter
25     omega = np.zeros(nBanks+2,np.float64)
26     # Populate array
27     for i in range(nBanks+2):
28         idx = np.argmin(np.abs(melfreq-mel[i]))
29         omega[i] = linfreq[idx]
30     # Generate filters
31     nframes = np.shape(SP)[1]
32     frange = np.shape(SP)[0]
33     mfcc = np.zeros([nBanks,frange])
34     h = np.zeros([nBanks,frange])
35     plt.figure(figsize=(8,4),dpi=170)
36     plt.ylabel('Magnitude')
37     plt.xlabel('Frequency (Hz)')
38     plt.title('Mel Filter Banks')
39     tick = np.arange(0,len(f),25)
40     label = np.around(tick*fs/512)
41     plt.xticks(tick,label)
42     for p in range(nBanks):
43         omegaC = omega[p]
44         omegaL = omega[p-1]
45         omegaR = omega[p+1]
46         for k in range(frange):
47             freq = f[k]
48             if (freq >= omegaL and freq < omegaC):
49                 h[p,k] = (2/(omegaR-omegaL)) * ((freq-omegaL)/(omegaC-omegaL))
50             elif (freq >= omegaC and freq < omegaR):
51                 h[p,k] = (2/(omegaR-omegaL)) * ((omegaR-freq)/(omegaR-omegaC))
52             else:
53                 pass
54             plt.plot(range(len(f)),h[p])
55             plt.savefig('figs/melBanks')
56             #mfcc[p-1,:] = np.sum((np.abs(h*SP[:,n]))**2)
57             return mfcc,h
```

References

- [1] “Signals & Systems.” Signals & Systems, by Alan V. Oppenheim et al., The McGraw-Hill, 1997, p. 78.