

Pattern Recognition Assignment 2: Speech Recognition

Hailey Yu Haihong
Levi Leyh

May 10, 2023

Contents

1	Plot of female and male speech and the music signal	1
2	Plot of the spectrograms	3
3	Comparison of the spectrogram and the (normalized) cepstrogram	4
4	4. Correlation of the spectral and cepstral coefficient series	6
5	Questions in the text	7
6	Further thoughts	7

1 Plot of female and male speech and the music signal

This is the used code:

Listing 1: Plot of female and male speech and the music signal

```
1  """
2  Authors: Hailey Yu Haihong, Levi Leyh
3
4  This code loads and visualizes three audio signals:
5  female speech, male speech, and music.
6  It then zooms in on a range of 20 ms in different regions of the signals
7  and plots the zoomed-in signals in separate subplots.
8  """
9  import matplotlib.pyplot as plt
10 import numpy as np
11 from scipy.io import wavfile
12
13 # Load the female and male speech signal and the music signal
14 # Get the sampling frequency
15 female_fs, female_speech = wavfile.read('Sounds/female.wav')
16 male_fs, male_speech = wavfile.read('Sounds/male.wav')
17 music_fs, music = wavfile.read('Sounds/music.wav')
18
19 # Plot the female speech signal
20 t_female_speech = np.arange(len(female_speech))/female_fs # Time axis in seconds
21 plt.figure()
22 plt.plot(t_female_speech, female_speech)
23 plt.xlabel('Time (s)')
24 plt.ylabel('Amplitude')
25 plt.title('Female Speech Signal')
26
27 # Plot the male speech signal
28 t_male_speech = np.arange(len(male_speech))/male_fs # Time axis in seconds
29 plt.figure()
30 plt.plot(t_male_speech, male_speech)
31 plt.xlabel('Time (s)')
32 plt.ylabel('Amplitude')
33 plt.title('Male Speech Signal')
34
35 # Plot the music signal
36 t_music = np.arange(len(music))/music_fs # Time axis in seconds
37 plt.figure()
38 plt.plot(t_music, music)
39 plt.xlabel('Time (s)')
40 plt.ylabel('Amplitude')
41 plt.title('Music Signal')
42
43 # Zoom in on a range of 20 ms in different regions of the signals
44 start_time = [0.5, 1, 1.5] # Start time of the zoomed-in range
45 duration = 0.02 # Duration of the zoomed-in range in seconds
46
47
48 # Create a figure for female speech signal
49 fig_f, axes_f = plt.subplots(nrows=len(start_time), figsize=(8, 6))
50 # Adjust spacing between subplots
51 plt.subplots_adjust(hspace=0.5)
52
53 # Zoom in on the female speech signal
54 for i, start in enumerate(start_time):
55     # Calculate indices and extract zoomed signal
56     idx_start = int(start * female_fs)
57     idx_end = idx_start + int(duration * female_fs)
58     t_zoom = t_female_speech[idx_start:idx_end]
59     female_speech_zoom = female_speech[idx_start:idx_end]
60
61     # Plot zoomed signal in corresponding subplot
62     ax = axes_f[i]
63     ax.plot(t_zoom, female_speech_zoom)
64     ax.set_xlabel('Time (s)')
65     ax.set_ylabel('Amplitude')
66     ax.set_title(f'Female Speech Signal (Zoomed in at {start:.2f} s)')
67
68
69 # Create figure for male speech signal
70 fig_m, axes_m = plt.subplots(nrows=len(start_time), figsize=(8, 6))
71 # Adjust spacing between subplots
72 plt.subplots_adjust(hspace=0.5)
73
74 # Zoom in on the male speech signal
75 for i, start in enumerate(start_time):
76     # Calculate indices and extract zoomed signal
77     idx_start = int(start * male_fs)
78     idx_end = idx_start + int(duration * male_fs)
79     t_zoom = t_male_speech[idx_start:idx_end]
80     male_speech_zoom = male_speech[idx_start:idx_end]
81
82     # Plot zoomed signal in corresponding subplot
83     ax = axes_m[i]
84     ax.plot(t_zoom, male_speech_zoom)
85     ax.set_xlabel('Time (s)')
86     ax.set_ylabel('Amplitude')
87     ax.set_title(f'Male Speech Signal (Zoomed in at {start:.2f} s)')
```

Listing 2: Plot of female and male speech and the music signal

```

88 # Create a figure with multiple subplots for the music signal
89 fig_mu, axes_mu = plt.subplots(nrows=len(start_time), figsize=(8, 6))
90 # Adjust spacing between subplots
91 plt.subplots_adjust(hspace=0.5)
92
93 # Zoom in on the music signal
94 for i, start in enumerate(start_time):
95     # Calculate indices and extract zoomed signal
96     idx_start = int(start * music_fs)
97     idx_end = idx_start + int(duration * music_fs)
98     t_zoom = t_music[idx_start:idx_end]
99     music_zoom = music[idx_start:idx_end]
100
101     # Plot zoomed signal in corresponding subplot
102     ax = axes_mu[i]
103     ax.plot(t_zoom, music_zoom)
104     ax.set_xlabel('Time (s)')
105     ax.set_ylabel('Amplitude')
106     ax.set_title(f'Music Signal (Zoomed in at {start:.2f} s)')
107
108 plt.show()
```

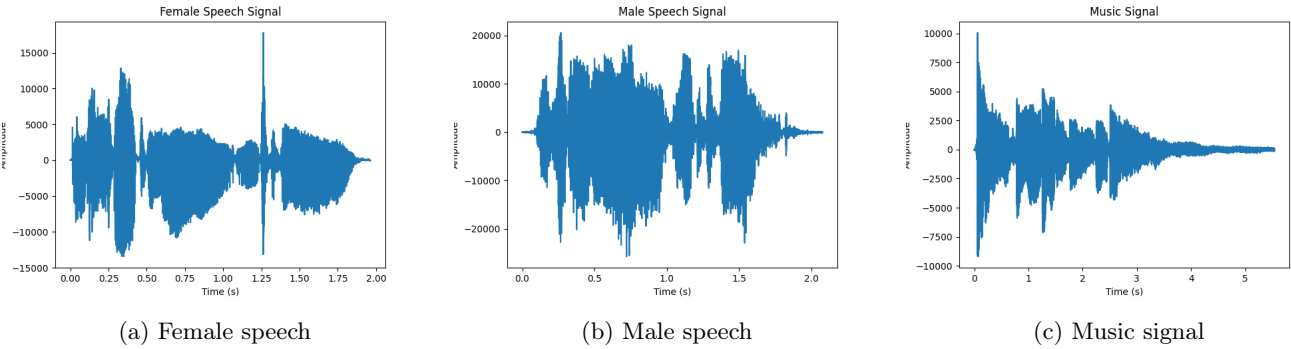


Figure 1: Plot of given audiofiles

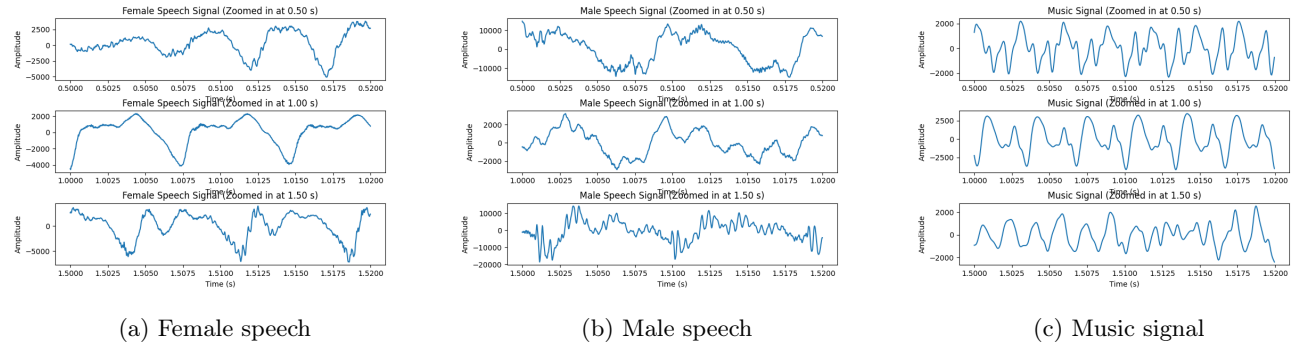


Figure 2: Zoom on representative signal patches

2 Plot of the spectrograms

This is the used code:

Listing 3: Plot of spectrograms

```
1  """
2  Authors: Hailey Yu Haihong, Levi Leyh
3
4  This code analyzes speech and music signals by computing and plotting spectrograms and Mel-
5  frequency cepstral coefficients (MFCCs).
6  It also normalizes the cepstrum and generates a cepstrogram displaying the normalized MFCC
7  coefficients for each signal.
8  Correlation matrices for both the spectral and cepstral coefficient series are computed and
9  plotted,
10 allowing for direct comparison between the two.
11 """
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import scipy.io.wavfile as wavfile
15 from scipy.signal import hann, spectrogram
16 from scipy.fftpack import dct
17 from python_speech_features import mfcc
18 from python_speech_features import logfbank
19
20 # Load the female and male speech signal and the music signal
21 female_fs, female_speech = wavfile.read('Sounds/female.wav')
22 male_fs, male_speech = wavfile.read('Sounds/male.wav')
23 music_fs, music = wavfile.read('Sounds/music.wav')
24
25 # Set spectrogram parameters
26 windowSize = 0.03
27 overlap = 0.015
28
29 # Compute spectrogram
30 female_windowSize = round(windowSize * female_fs )
31 female_overlap = round(female_windowSize * 0.5)
32 female_F, female_T, female_S = spectrogram(female_speech, fs=female_fs, window=hann(
33     female_windowSize),
34     nperseg=female_windowSize, noverlap=female_overlap
35 )
36
37 male_windowSize = round(windowSize * male_fs )
38 male_overlap = round(male_windowSize * 0.5)
39 male_F, male_T, male_S = spectrogram(male_speech, fs=male_fs, window=hann(male_windowSize),
40     nperseg=male_windowSize, noverlap=male_overlap
41 )
42
43 music_windowSize = round(windowSize * music_fs )
44 music_overlap = round(music_windowSize * 0.5)
45 music_F, music_T, music_S = spectrogram(music, fs=music_fs, window=hann(music_windowSize),
46     nperseg=music_windowSize, noverlap=music_overlap
47 )
48
49 # Plot spectrograms
50 plt.subplot(3,3,1)
51 plt.imshow(20*np.log10(abs(female_S)), cmap='jet', aspect='auto', origin='lower',
52     extent=[female_T.min(), female_T.max(), female_F.min(), female_F.max()])
53 plt.colorbar()
54 plt.xlabel('Time (s)')
55 plt.ylabel('Frequency (Hz)')
56 plt.title('spectrogram female speech')
57
58 plt.subplot(3,3,2)
59 plt.imshow(20*np.log10(abs(male_S)), cmap='jet', aspect='auto', origin='lower',
60     extent=[male_T.min(), male_T.max(), male_F.min(), male_F.max()])
61 plt.colorbar()
62 plt.xlabel('Time (s)')
63 plt.ylabel('Frequency (Hz)')
64 plt.title('spectrogram male speech')
65
66 plt.subplot(3,3,3)
67 plt.imshow(20*np.log10(abs(music_S)), cmap='jet', aspect='auto', origin='lower',
68     extent=[music_T.min(), music_T.max(), music_F.min(), music_F.max()])
69 plt.colorbar()
70 plt.xlabel('Time (s)')
71 plt.ylabel('Frequency (Hz)')
72 plt.title('spectrogram music')
```

3 Comparison of the spectrogram and the (normalized) cepstrogram

This is the used code:

Listing 4: Plot of female and male speech and the music signal (continued)

```
68
69
70 # Compute MFCC coefficients
71 female_coeffs = mfcc(female_speech, female_fs, winlen=windowSize, winstep=overlap, numcep=13)
72 male_coeffs = mfcc(male_speech, male_fs, winlen=windowSize, winstep=overlap, numcep=13)
73 music_coeffs = mfcc(music, music_fs, winlen=windowSize, winstep=overlap, numcep=13)
74
75 print(len(female_coeffs), len(female_coeffs[0]))
76 print(len(female_S), len(female_S[0]))
77
78
79 # Normalize cepstrum
80 def normalize(mfcc):
81     mean_cep = np.mean(mfcc, axis=0)
82     std_cep = np.std(mfcc, axis=0)
83     cepstrogram_centered = mfcc - mean_cep
84     return cepstrogram_centered / std_cep
85
86 female_coeffs_norm = normalize(female_coeffs)
87 male_coeffs_norm = normalize(male_coeffs)
88 music_coeffs_norm = normalize(music_coeffs)
89
90 print(np.min(female_coeffs))
91
92
93 # Plot the MFCC coefficients and the normalized cepstrogram
94 plt.subplot(3, 3, 4)
95 plt.imshow(female_coeffs.T, aspect='auto', cmap='jet', origin='lower',
96            extent=[female_T.min(), female_T.max(), np.min(female_coeffs), np.max(female_coeffs)])
97 plt.xlabel('Time (s)')
98 plt.ylabel('MFCC coefficient')
99 plt.title('cepstrogram female speech')
100 plt.colorbar()
101
102 plt.subplot(3, 3, 5)
103 plt.imshow(male_coeffs.T, aspect='auto', cmap='jet', origin='lower',
104            extent=[male_T.min(), male_T.max(), np.min(male_coeffs), np.max(male_coeffs)])
105 plt.xlabel('Time (s)')
106 plt.ylabel('MFCC coefficient')
107 plt.title('cepstrogram male speech')
108 plt.colorbar()
109
110 plt.subplot(3, 3, 6)
111 plt.imshow(music_coeffs.T, aspect='auto', cmap='jet', origin='lower',
112            extent=[music_T.min(), music_T.max(), np.min(music_coeffs), np.max(music_coeffs)])
113 plt.xlabel('Time (s)')
114 plt.ylabel('MFCC coefficient')
115 plt.title('cepstrogram music')
116 plt.colorbar()
117
118 # Plot the normalized mfcc
119 plt.subplot(3, 3, 7)
120 plt.imshow(female_coeffs_norm.T, aspect='auto', cmap='jet', origin='lower',
121            extent=[female_T.min(), female_T.max(), np.min(female_coeffs_norm), np.max(
122                female_coeffs_norm)])
123 plt.xlabel('Time (s)')
124 plt.ylabel('MFCC coefficient normalized')
125 plt.title('normalized cepstrogram female speech')
126 plt.colorbar()
127
128 plt.subplot(3, 3, 8)
129 plt.imshow(male_coeffs_norm.T, aspect='auto', cmap='jet', origin='lower',
130            extent=[male_T.min(), male_T.max(), np.min(male_coeffs_norm), np.max(male_coeffs_norm)
131                ])
132 plt.xlabel('Time (s)')
133 plt.ylabel('MFCC coefficient normalized')
134 plt.title('normalized cepstrogram male speech')
135 plt.colorbar()
136
137 plt.subplot(3, 3, 9)
138 plt.imshow(music_coeffs_norm.T, aspect='auto', cmap='jet', origin='lower',
139            extent=[music_T.min(), music_T.max(), np.min(music_coeffs_norm), np.max(
140                music_coeffs_norm)])
141 plt.xlabel('Time (s)')
142 plt.ylabel('MFCC coefficient normalized')
143 plt.title('cepstrogram music')
144 plt.colorbar()
145
146 plt.show()
```

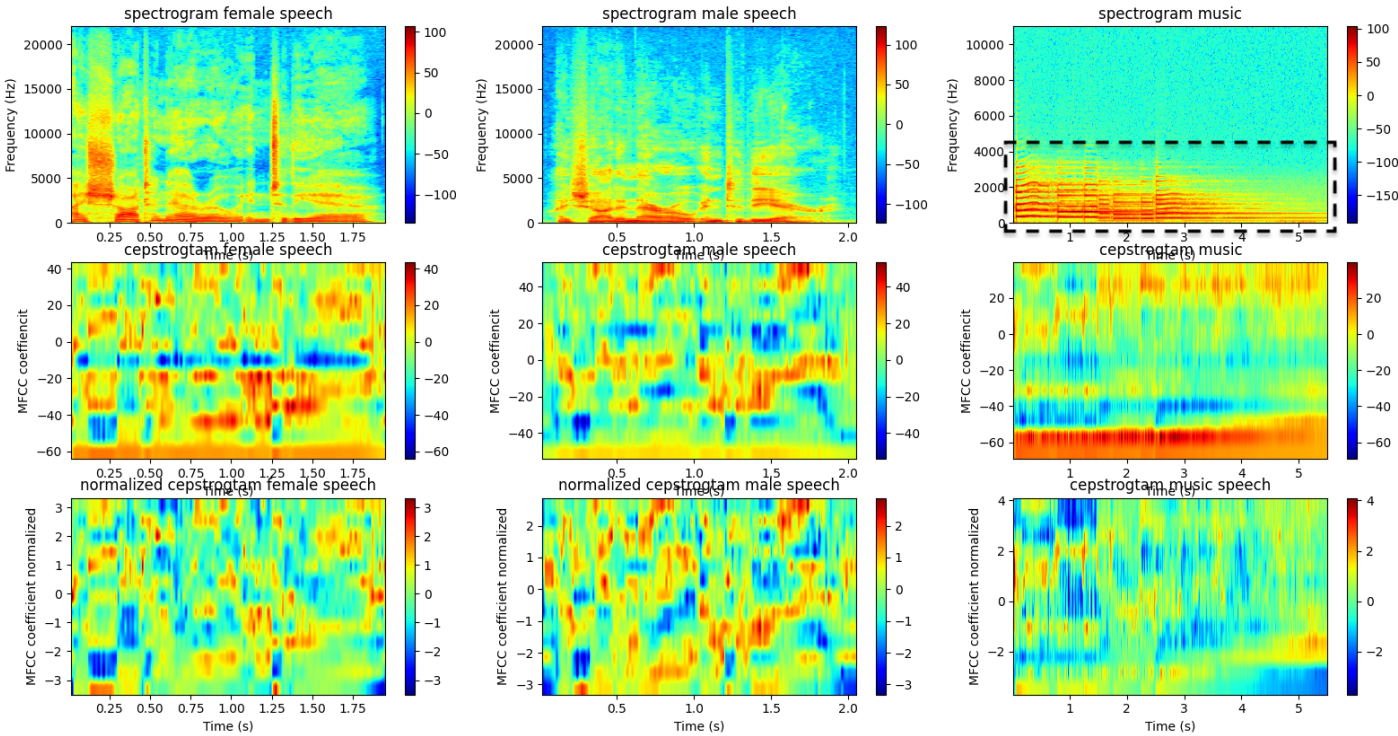


Figure 3: Plotted spectrograms and cepstrograms

In the music sample, harmonics are identified and marked by a dotted square. It’s easy to see there are many horizontal lines in that area, which corresponds to harmonic components.

4 4. Correlation of the spectral and cepstral coefficient series

This is the used code:

Listing 5: Correlation of spectrograms and cepstrogram (continued)

```
144
145
146 # Compute the correlations
147 female_corr_spec = np.corrcoef(20*np.log10(abs(female_S)), rowvar=True)
148 male_corr_spec   = np.corrcoef(20*np.log10(abs(male_S)),   rowvar=True)
149 music_corr_spec  = np.corrcoef(20*np.log10(abs(music_S)),  rowvar=True)
150
151 female_corr_ceps = np.corrcoef(female_coeffs.T, rowvar=True)
152 male_corr_ceps   = np.corrcoef(male_coeffs.T,   rowvar=True)
153 music_corr_ceps  = np.corrcoef( music_coeffs.T, rowvar=True)
154
155 female_corr_ceps_norm = np.corrcoef(female_coeffs_norm.T, rowvar=True)
156 male_corr_ceps_norm   = np.corrcoef(male_coeffs_norm.T,   rowvar=True)
157 music_corr_ceps_norm  = np.corrcoef( music_coeffs_norm.T, rowvar=True)
158
159 # Plot the correlations
160 plt.subplot(3,3,1)
161 plt.imshow(np.abs(female_corr_spec), aspect='auto', origin='lower', cmap='gray')
162 plt.title('spectrum correltion female')
163 plt.colorbar()
164
165 plt.subplot(3,3,2)
166 plt.imshow(np.abs(male_corr_spec), aspect='auto', origin='lower', cmap='gray')
167 plt.title('spectrum correltion male')
168 plt.colorbar()
169
170 plt.subplot(3,3,3)
171 plt.imshow(np.abs(music_corr_spec), aspect='auto', origin='lower', cmap='gray')
172 plt.title('spectrum correltion music')
173 plt.colorbar()
174
175 plt.subplot(3,3,4)
176 plt.imshow(np.abs(female_corr_ceps), aspect='auto', origin='lower', cmap='gray')
177 plt.title('cepstrum correltion female')
178 plt.colorbar()
179
180 plt.subplot(3,3,5)
181 plt.imshow(np.abs(male_corr_ceps), aspect='auto', origin='lower', cmap='gray')
182 plt.title('cepstrum correltion male')
183 plt.colorbar()
184
185 plt.subplot(3,3,6)
186 plt.imshow(np.abs(music_corr_ceps), aspect='auto', origin='lower', cmap='gray')
187 plt.title('cepstrum correltion music')
188 plt.colorbar()
189
190 plt.subplot(3,3,7)
191 plt.imshow(np.abs(female_corr_ceps_norm), aspect='auto', origin='lower', cmap='gray')
192 plt.title('normalized cepstrum correltion female')
193 plt.colorbar()
194
195 plt.subplot(3,3,8)
196 plt.imshow(np.abs(male_corr_ceps_norm), aspect='auto', origin='lower', cmap='gray')
197 plt.title('normalized cepstrum correltion male')
198 plt.colorbar()
199
200 plt.subplot(3,3,9)
201 plt.imshow(np.abs(music_corr_ceps_norm), aspect='auto', origin='lower', cmap='gray')
202 plt.title('norlamized cepstrum correltion music')
203 plt.colorbar()
204
205 plt.show()
```

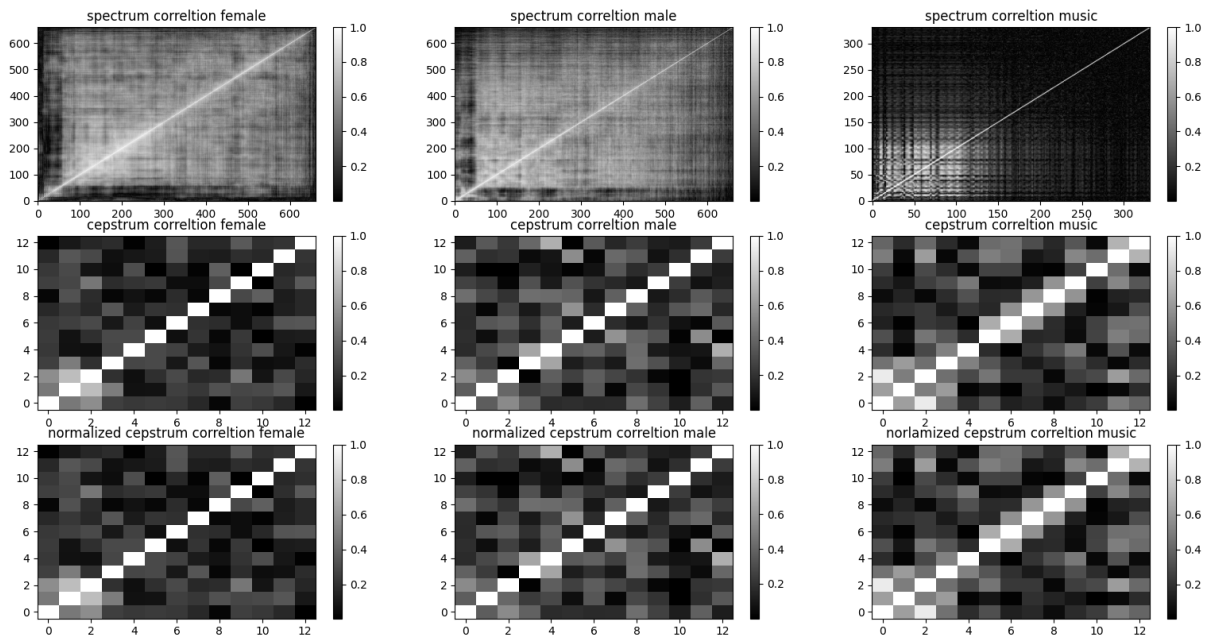


Figure 4: Plotted correlation matrices

5 Questions in the text

Which representation do you think is the easiest for you, as a human, to interpret, and why?

The easiest representation for a human to interpret is the spectrogram because it provides a visual representation of the frequency content of a signal over time. This allows for easy identification of patterns and variations within the signal. For example, it's very easy to find the harmonics in spectrogram by searching for horizontal lines. On the other hand, the cepstrogram shows the spectral envelope of a signal, which is a less intuitive concept for most people.

Can you see that the male and female spectrograms represent the same phrase? Could a computer discover this? Why/why not?

Yes, you can see that the male and female spectrograms represent the same phrase from the figures and can be easily distinguished from the music signal. A computer can discover this as well depending on the quality of the samples and the algorithms used for analysis. Humans are generally better at recognizing patterns and making judgments based on context, while computers are better at analyzing large amounts of data quickly and accurately. When comparing spectrograms and checking for similar phases in two audio files, humans can interpret the content of the audio and identify patterns in the spectrograms based on their knowledge of language and speech and their perception of sound, while computers are better suited for tasks like identifying specific patterns or comparing large datasets but may struggle to interpret the context or meaning behind the audio and may not recognize subtle variations or nuances.

Can you see that the male and female cepstrograms represent the same phrase now? What about a computer?

The cepstrogram is more difficult for human to interpret. Even looking at the normalized cepstrogram, it's hard for human-eye to capture similarities in this complex pattern. However, it would be easier for computers, which can convert visual data into numbers and calculations. We can compute the distance between two cepstral vectors using techniques such as Euclidean distance or cosine similarity. If the distance between two cepstrograms is small, it indicates that they are similar. Additionally, we can use machine learning algorithms such as k-NN, SVM, or neural networks to train a model to classify cepstral vectors as similar or dissimilar.

Which matrix, spectral or cepstral, looks the most diagonal to you?

For speech signals, the cepstral matrix looks more diagonal. There's a distinct contrast between the diagonal line (bright) and the rest (dark), which means that the correlation between the coefficients are very small. In contrast, the spectral matrix is quite bright not only in the diagonal line, implying that its coefficients are strongly correlated. As for the music signal, there's not much difference between spectral and cepstral. Maybe this is because MFCC is more for speech signals rather than music.

6 Further thoughts

Some thoughts on the possibility of confusing the MFCC representation in a speech recognizer. Can you think of a case where two utterances have noticeable differences to a human listener, and may come with different interpretations or connotations, but still have very similar MFCCs? (Hint: Think about what information the MFCCs remove.) What about the opposite situation—are there two signals that sound very similar to humans, but have substantially different MFCCs?

Yes, two utterances can have noticeable differences to humans but still have very similar MFCCs. This is because MFCCs remove information that is important for human perception but not necessarily for speech recognition. Conversely, two signals that sound similar to humans may have substantially different MFCCs, due to differences in acoustic characteristics that are not perceptually salient to humans but are captured by the MFCCs.