

2.5. Acoustic features

The data for this project consists of audio cough recordings. These audio recordings are converted into spectrogram images so that image processing techniques can be applied. A spectrogram is an image that captures the frequencies and amplitudes of audio signals using a ~~Discrete Fourier Transform (DFT)~~. Figure 2.15 shows how an audio signal is converted into a spectrogram.

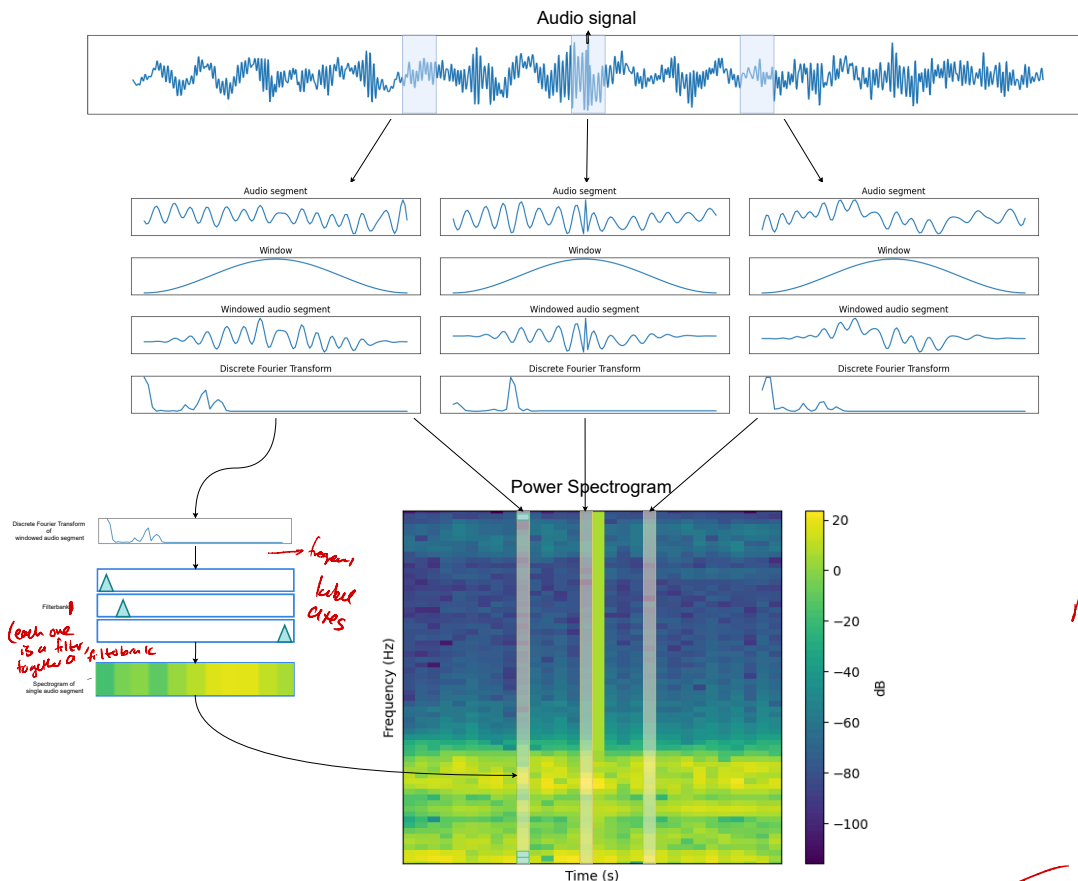


Figure 2.15: The process to convert a raw audio signal to a spectrogram. The audio signal is split into segments. Each segment is multiplied with a window function. A Discrete Fourier Transform (DFT) is applied to each of the windowed segments. The DFT is multiplied with filterbanks and the result is stacked to form a spectrogram image.

The audio signal is divided into ~~audio~~ segments. A window is created using a windowing function such as a Hamming window! The audio segment is multiplied by the window to minimise the frequency leakage between overlapping signals. A Fast Fourier Transform is then applied to the windowed audio signal to indicate the probabilities that a certain frequency is present in that audio signal. Each FFT of the audio segments are then multiplied with filterbanks to summarise the information. Finally, the resulting values are then stacked together to form a contour plot with time on the x axis, frequency on the y axis and the magnitude (in decibels) on the z axis. The contour plot is shown as

Nice figure!

is then applied to each frame in order to reduce the spectral leakage caused by the truncation at the frame boundaries.

often referred to as frames.

Short, consecutive and sometimes overlapping

compute

the squared magnitude taken

probabilities that

Each FFT of the audio segments are then multiplied with filterbanks to summarise the information. Finally, the resulting values are then stacked together to form a contour plot with time on the x axis, frequency on the y axis and the magnitude (in decibels) on the z axis. The contour plot is shown as

The resulting magnitude spectrum is

reduce the frequency resolution.

NOTE: a spectrogram generally does not include the filterbank step, but your feature extraction does.

NOTE
The result of applying the linear filterbank to the squared magnitude FFT is a vector of spectral energies
each one is a filter together a filterbank
label axes
Frequency (Hz)
Time (s)
dB
20
0
-20
-40
-60
-80
-100

a 2D image with the colours indicating the magnitude of the frequencies. *for each frame*

Note you can also just use a mel scale on the frequency axis (ie map each lin FFT freq to a mel freq).

A mel spectrogram is a spectrogram where the linear filterbanks are replaced by mel filterbanks using the mel-scale as shown in Equation 2.20. The mel scale is used to mimic how humans perceive sounds. It is based on a logarithmic scale and is calculated such that each unit sounds equal in pitch difference.

how does something 'sound equal in pitch difference'?

$$m = 1127 \times \log\left(1 + \frac{f}{700}\right) \quad (2.20)$$

Mel-frequency cepstral coefficients (MFCCs) are *calculated* by applying a Discrete Cosine Transform (DCT) to the mel *scale spectral energies* spectrograms. MFCCs capture uncorrelated information from audio signals that are useful to machine learning models.

2.6. Conclusion

This chapter discussed various theoretical concepts needed to understand contrastive learning and the terminology used throughout this project. The next chapter is a literature review of the existing work done in the field of TB classification and contrastive learning.

NOTE MFCCs are useful for speech because they can be set up to be fairly independent of the speaker's pitch and also fairly independent of the channel (room acoustics, microphone characteristics etc)

How useful they are for cough remains to be clearly established.

(I think the pitch independence is not important, but possibly the channel independence is)

However it is not correct to say that they are 'useful to me' in general. Rather, they extract useful info from speech for me to use.