Speech Signal Analysis 2

Hao Tang

Automatic Speech Recognition—ASR Lecture 3 22 January 2024

→ < ∃ →</p>

- Recap of spectrograms
- Auditory system
 - Masking
 - Mel filters
- Speech production model
 - Fundamental frequencies
 - Formants
- Mel Frequency Cepstral Coefficients



- dithering, removing DC offset, pre-emphasis
- windowing
- Discrete Fourier transform (DFT)
- Short-time Fourier transform (STFT)

Discrete Fourier Transform



Hao Tang Speech Signal Analysis 2

Discrete Fourier Transform



Hao Tang Speech Signal Analysis 2

$$X[k] = a + bi$$

• Real:
$$\Re{X[k]} = a$$

- Imaginary: $\mathfrak{Im}{X[k]} = b$
- Magnitude: $|X[k]| = \sqrt{a^2 + b^2}$

• Phase:
$$\angle X[k] = \arccos \frac{a}{\sqrt{a^2+b^2}}$$

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶

э

Spectrogram

Magnitude



Phase



- Spectrogram = Magnitude spectrogram = Power spectrogram
- Phase is not as important as magnitude for speech intelligibility.

Without log



æ

(4回) 4 回) 4 回)

Without log



With log



æ

-≣->

・日・ ・ ヨ・・



æ

æ

< Ξ



æ

æ

< Ξ



▲御▶ ▲ 臣▶ ▲ 臣▶



æ

-≣->

-

- One sound affects the presence of another sound.
- Both sounds are present, so masking is purely perceptual.
- Masking is a nonlinear effect.
- Many applications take advantage of masking (e.g., MP3).



æ

- 4 回 ト 4 三 ト 4 三 ト





▲御 ▶ ▲ 臣 ▶ ▲ 臣 ▶



▲御▶ ▲ 臣▶ ▲ 臣▶





- Triangle-shaped
- Asymmetric
- Sensitive to the amount of energy
- With larger bandwidth at higher frequency



- Triangle-shaped
- Asymmetric
- Sensitive to the amount of energy
- With larger bandwidth at higher frequency



- Triangle-shaped
- Asymmetric
- Sensitive to the amount of energy
- With larger bandwidth at higher frequency



- Triangle-shaped
- Asymmetric
- Sensitive to the amount of energy
- With larger bandwidth at higher frequency



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

300 Hz vs 310 Hz
2000 Hz vs 2010 Hz

▲ (日) ▶ ▲ (日)

æ

≣⇒

frequency (Mel)

< ロ > < 回 > < 回 > < 回 > < 回 >





< 同 ト < 三 ト < 三 ト



æ

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

$$Y[n] = \sum_{k=0}^{T-1} X[k] \cdot H_n[k]$$

- H_n is the *n*-th Mel filter.
- Mel filters are applied to the magnitude spectrum with dot product.
- The result is an *n*-dimensional vector for *n* Mel filters.

$$Y[n] = \sum_{k=0}^{T-1} X[k] \cdot H_n[k]$$

$$Y = \begin{bmatrix} H_1[0] & H_1[1] & \cdots \\ H_2[0] & H_2[1] & \cdots \\ \vdots & \vdots & \\ H_n[0] & H_n[1] & \cdots \end{bmatrix} \begin{bmatrix} X[0] \\ X[1] \\ \vdots \\ X[T-1] \end{bmatrix} = \begin{bmatrix} H_1 \\ H_2 \\ \vdots \\ H_n \end{bmatrix} X = HX$$

< ロ > < 回 > < 回 > < 回 > < 回 >

linear spectrogram



Mel spectrogram



linear spectrogram



Mel spectrogram







< ロ > < 部 > < き > < き >

Formants



< ロ > < 回 > < 回 > < 回 > < 回 >

Speech Production



э

breathing







・日・ ・ ヨ・・

æ



< ロ > < 回 > < 回 > < 回 > < 回 >



< ロ > < 回 > < 回 > < 回 > < 回 >



・ロト ・四ト ・ヨト ・ヨト

Resonance Frequency of A Tube



Frequency Response of A Tube



▲御▶ ▲ 臣▶ ▲ 臣▶



▲□ ▶ ▲ 臣 ▶ ▲ 臣 ▶ ○ 臣 ○ のへの



▲□ ▶ ▲ 臣 ▶ ▲ 臣 ▶ ● 臣 ● のへで

- Fundamental frequency
 - The first frequency component of the glottal pulse
 - Leading to pitch when perceived
- Harmonics
 - Subsequent frequency components of the glottal pulse
- Formants
 - Resonance frequencies of the vocal tract
 - Leading to the production and perception of certain phones, particularly vowels



æ



æ

-≣->

→ < Ξ → <</p>

Formants



< ロ > < 回 > < 回 > < 回 > < 回 >



포 🛌 포

Discrete Fourier Transform



Hao Tang Speech Signal Analysis 2

- Extract Mel spectrogram.
- Apply DFT to every spectrum.
- Truncate the high-frequency components.



Hao Tang Speech Signal Analysis 2

"All models are wrong, but some are useful."

-George Box, 1978

- dithering
- removing DC offset
- pre-emphasis
- windowing
- DFT
- Apply Mel filters
- DCT
- Truncate the high-frequency components

• Chapter 3–4, O'Shaughnessy, "Speech Communications: Human and Machine," 2000.