Automatic Speech Recognition handout (1) Jan - Mar 2009

Revision: 1.1

Speech Signal Processing and Feature Extraction

Hiroshi Shimodaira (h.shimodaira@ed.ac.uk)

Speech Production



Vocal Organs & Vocal Tract

 $s(t) = h(t) * v(t) = \int_0^\infty h(\tau) x(t-\tau) d\tau$ time domain: *Fourier transform* frequency domain: $S(\Omega) = H(\Omega)V(\Omega)$

ASR (H. Shimodaira)



Speech Communication



Automatic Speech Recognition



ASR (H. Shimodaira)

Signal Analysis for ASR



Convert acoustic signal into a sequence of feature vectors



Feature parameters for ASR

Features should

- contain sufficient information to distinguish phonemes / phones
 - good time-resolutions [e.g. 10ms]
 - good frequency-resolutions [e.g. 20 channels/Bark-scale]
- **not contain (or be separated from)** F_0 and its harmonics
- **be robust against speaker variation**
- **be robust against noise / channel distortions**
- have good characteristics in terms of pattern recognition
 - The number of features is as few as possible
 - Features are independent of each other

 \Rightarrow A large number of features have been proposed

ASR (H. Shimodaira)

Converting analogue signals to machine readable form

- **Discretisation (digitising)** $x_c(t) \rightarrow x[n]$
 - continuous time \Rightarrow discrete time
 - continuous amplitude ⇒discrete amplitude

1:4

Sampling of continuous-time signals



- **Continuous-time signal:** $x_c(t)$
- Modulated signal by a periodic impulse train:

 $x_s(t) = x_c(t) \sum_{n=-\infty}^{\infty} \delta(t - nT_s) = \sum_{n=-\infty}^{\infty} x_c(nT_s)\delta(t - nT_s)$

Sampled signal: $x[n] = x_s(nT_s) \quad \cdots$ discrete-time signal

 T_s : Sampling interval

1:8

ASR (H. Shimodaira)

Sampling of continuous-time signals(cont. 2)

Q: Is the C/D conversion invertible ?

$$x_c(t) \xrightarrow{C/D} x[n] \xrightarrow{D/C} x_c(t)$$
?

Sampling of continuous-time signals(cont. 3)

Q: Is the C/D conversion invertible ?

$$x_c(t) \xrightarrow{C/D} x[n] \xrightarrow{D/C} x_c(t)$$
?

A: "No" in general, but "Yes" under a special condition: "Nyquist sampling theorem"

If $x_c(t)$ is band-limited (i.e. no frequency components > $F_s/2$), then $x_c(t)$ can be fully reconstructed by x[n].

$$\begin{aligned} x_{c}(t) &= h_{T_{s}}(t) * \sum_{k=-\infty}^{\infty} x[k]\delta(t-kT_{s}) = \sum_{k=-\infty}^{\infty} x[k]h_{T_{s}}(t-kT_{s}) \\ h_{T_{s}}(t) &= \operatorname{sinc}(t/T_{s}) = \frac{\sin(\pi t/T_{s})}{\pi t/T_{s}} \end{aligned}$$

 $F_s/2$: Nyquist Frequency, $F_s = 1/T_s$: Sampling Frequency

ASR (H. Shimodaira)

I:10

Sampling of continuous-time signals(cont. 4)

Interpretation in frequency domain:



ASR (H. Shimodaira)



Questions

- **1.** What sampling frequencies (F_s) are used for ASR ?
 - **microphone voice:** $12kHz \sim 20kHz$
 - **telephone voice:** $\sim 8kHz$
- 2. What are the advantages / disadvantages of using higher F_s ?
- 3. Why is pre-emphasis (+6dB/oct.) employed?

$$x[n] = x_0[n] - ax_0[n-1], \quad a = 0.95 \sim 0.97$$

ASR (H. Shimodaira)

```
I:12
```

An interpretation of FT

Inner product between two vectors (Linear Algebra)

2-dimensional case

$$\mathbf{a} = (a_1, a_2)^t$$

$$\mathbf{b} = (b_1, b_2)^t$$

$$\mathbf{a} \cdot \mathbf{b} = \mathbf{a}^t \mathbf{b} = a_1 b_1 + a_2 b_2$$

$$= \parallel \mathbf{a} \parallel \parallel \mathbf{b} \parallel \cos \theta$$



$$\begin{aligned} \mathbf{x} &\triangleq \{x[n]\}_{-\infty}^{\infty} \\ \mathbf{e}_{\omega} &\triangleq \left\{\mathbf{e}^{j\omega n}\right\}_{-\infty}^{\infty} = \left\{\cos(\omega n) + j\sin(\omega n)\right\}_{-\infty}^{\infty} \\ &\triangleq \cos_{\omega} + j\sin_{\omega} \\ X(e^{j\omega}) &= \sum_{n=-\infty}^{\infty} x[n]e^{-j\omega n} = \mathbf{x} \cdot \mathbf{e}^{j\omega n} = \mathbf{x} \cdot \cos_{\omega} + j\mathbf{x} \cdot \sin_{\omega} \end{aligned}$$

 $m{x}\cdot\cos_{\omega}$: proportion of how much \cos_{ω} component is contained in $m{x}$

ASR (H. Shimodaira)

I:14

Spectral analysis: Fourier Transform

FT for continuous-time signals (& continuous-frequency)

$$\begin{array}{lll} X_c(\Omega) &=& \int_{-\infty}^{\infty} x_c(t) e^{-j\Omega t} dt & (\text{time domain} \to \text{freq. domain}) \\ x_c(t) &=& \frac{1}{2\pi} \int_{-\infty}^{\infty} X_c(\Omega) e^{j\Omega t} d\Omega & (\text{freq. domain} \to \text{time domain}) \end{array}$$

FT for discrete-time signals (& continuous-frequency)

$$\begin{array}{lll} X(e^{j\omega}) &=& \sum\limits_{n=-\infty}^{\infty} x[n]e^{-j\omega n} \\ x[n] &=& \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega})e^{j\omega n}d\omega \end{array}$$

 $|X(e^{j\omega})|^2 \cdots$ Power spectrum $\log |X(e^{j\omega})|^2 \cdots$ Log power spectrum

 $\begin{array}{ll} \mbox{where} & \omega = 2\pi f, \; f = 1/T, \; \; \omega = T_s \Omega, \\ & e^{-j\omega n} = \cos(\omega n) + j \sin(\omega n), \quad j: \; \mbox{the imaginary unit} \end{array}$

Short-time Spectrum Analysis

Problem with FT

- Assuming signals are stationary: signal properties do not change over time
- If signals are non-stationary ⇒ loses information on time varying features
- ⇒ Short-time Fourier transform (STFT) (Time-dependent Fourier transform)

11

Divide signals into short-time segments (frames) and apply FT to each frame.

Short-time Spectrum Analysis(cont. 2)



ASR (H. Shimodaira)

I:16

Short-time Spectrum Analysis(cont. 3)

Trade-off problem of short time spectrum analysis



 \Rightarrow a compromise:

window width (frame width): $20 \sim 30$ ms window shift (frame shift): $5 \sim 15$ ms

The Effect of Windowing in STFT

Time domain:

 $y_k[n] = w_k[n]x[n], \quad w_k[n]$: time-window for k-th frame Simply cutting out a short segment (frame) from x[n] implies applying a rectangular window on to x[n].

 \Rightarrow causes discontinuities at the edges of the segment. Instead, a tapered window is usually used.. e.g. Hamming ($\alpha = 0.46164$) or Hanning ($\alpha = 0.5$) window)



The Effect of Windowing in STFT(cont. 2)

Frequency domain:

$$Y_k(e^{j\omega}) = \frac{1}{2\pi} \int_{-\pi}^{\pi} W_k(e^{j\theta}) X(e^{j(\omega-\theta)}) d\theta \quad \cdots \ \text{Periodic convolution}$$

- Power spectrum of the frame is given as a periodic convolution between the power spectra of x[n] and wk[n].
- If we want Y_k(e^{jω}) = X(e^{jω}), the necessary and sufficient condition for this is W_k(e^{jω}) = δ(ω),

i.e. $w_k[n] = \mathcal{F}^{-1}\delta(\omega) = 1$, which means the length of $w_k[n]$ is infinite.

 \Rightarrow there is no window function of finite length that causes no distortion.



Problems with STFT

- **The estimated power spectrum contains harmonics of** F_0 , which makes it difficult to estimate the envelope of the spectrum.
- Frequency bins of STFT are highly correlated each other, i.e. power spectrum representation is highly redundant.



Cepstrum Analysis(cont. 2) x[n] = h[n] * v[n] $\downarrow \mathcal{F}$ (Fourier transform)

Log spectrum

 $\log |X(e^{j\omega})| = \log |H(e^{j\omega})| +$ $\log |V(e^{j\omega})|$ (spectral envelope) (spectral fine structure)

 $X(e^{j\omega}) = H(e^{j\omega})V(e^{j\omega})$

↓ log

Cepstrum

$$c(\tau) = \mathcal{F}^{-1} \left\{ \log |X(e^{j\omega})| \right\}$$
$$= \mathcal{F}^{-1} \left\{ \log |H(e^{j\omega})| \right\} + \mathcal{F}^{-1} \left\{ \log |V(e^{j\omega})| \right\}$$

 \mathcal{F}^{-1}

ASR (H. Shimodaira)

h[n]: vocal tract

v[n]: glottal sounds

LPC Analysis

Linear Predictive Coding (LPC): a model-based / parametric spectrum estimation Assume a "linear system" for human speech production sound source $x[n] \Rightarrow |vocal tract| \Rightarrow speech y[n]$

$$x[n] \longrightarrow h[n] \longrightarrow y[n] \qquad h[n]: \text{ impulse response}$$
$$y[n] = h[n] * x[n] = \sum_{k=0}^{\infty} h[k] x[n-k]$$

Using a model enables us to

- estimate a spectrum of vocal tract from small amount of observations
- represent the spectrum with a small number of parameters
- synthesise speech with the parameters

ASR (H. Shimodaira)

LPC analysis in detail

Predict y[n] from $y[n-1], y[n-2], \cdots$ $\hat{y}[n] = \sum_{k=1}^{N} a_k y[n-k]$ $e[n] = y[n] - \hat{y}[n] = y[n] - \sum_{k=1}^{N} a_k y[n-k] \cdots$ prediction error

Optimisation problem —

Find $\{a_k\}$ that minimises the mean square (MS) error:

$$P_{e} = E\left\{e^{2}[n]\right\} = E\left\{\left(y[n] - \sum_{k=1}^{N} a_{k}y[n-k]\right)^{2}\right\}$$

I:25

1:24

ASR (H. Shimodaira)

Spectrum estimated by FT & LPC



ASR (H. Shimodaira)



LPC summary

- **Spectrum can be modelled/coded with around** 14LPCs.
- LPC family
 - PARCOR (Partial Auto-Correlation Coefficient)
 - LSP (Line Spectral Pairs) / LSF (Line Spectrum Frequencies)
 - CSM (Composite Sinusoidal Model)
- LPC can be used to predict log-area ratio coefficients lossless tube model
- LPC-(Mel)Cepstrum: LPC based cepstrum.
- Drawback:
 - LPC assumes AR model which does not suit to model nasal sounds that have zeros in spectrum.
 - **Difficult to determine the prediction order** *N***.**

l:27

 $^{\{}a_k\}$: LPC coefficients

Taking into Perceptual Attributes

Physical quality	Perceptual quality
Intensity	Loudness
Fundamental frequency	Pitch
Spectral shape	Timbre
Onset/offset time	Timing
Phase difference in binaural hearing	Location

Technical terms

- equal-loudness contours
- masking

ASR (H. Shimodaira)

- **auditory** filters (critical-band filters)
- critical bandwidth

Taking into Perceptual Attributes(cont. 3)

Non-linear frequency scale

Bark scale

 $b(f) = 13 \arctan(0.00076f) + 3.5 \arctan((f/7500)^2)$ [Bark]

Mel scale



```
ASR (H. Shimodaira)
```

Taking into Perceptual Attributes(cont. 2)



Filter Bank Analysis



$$x_i[n] = h_i[n] * x[n] = \sum_{k=0}^{M_i-1} h_i[k]x[n-k]$$

 $h_i[n]$: Impulse response of Bandpass filter i

ASR (H. Shimodaira)

1:28

Filter Bank Analysis(cont. 2)



Trade-off problem

ASR (H. Shimodaira)

1:32

MFCC

MFCC: Mel-frequency Cepstrum Coefficients c[n]

$$x[n] \xrightarrow{\mathbf{DFT}} X[k] \to |X[k]|^2 \xrightarrow{\mathbf{Mel-frequency} \atop \mathbf{fi \, lterbank}} \log |S[m]| \xrightarrow{\mathbf{DCT}} c[n]$$

DCT: $c[n] = \sqrt{\frac{2}{N}} \sum_{i=i}^{N} s[i] \cos \left(\frac{\pi n(i-0.5)}{N}\right)$, where $s[i] = \log |S[i]|$

- MFCCs are widely used in HMM-based ASR systems.
- **The first 12 MFCCs** ($c[1] \sim c[12]$) are generally used.

ASR (H. Shimodaira)

Filter Bank Analysis(cont. 3)

Another implementation: apply a mel-scale filter bank to STFT power spectrum to obtain mel-scale power spectrum

MFCC(cont. 2)

- MFCCs are less correlated each other than DCT/Filter-bank based spectrum.
- **Good compression rate.**

Feature	dimensionality / frame
Speech wave	400
DCT Sepctrum	$64 \sim 256$
Filter-bank	$10 \sim 20$
MFCC	12

where $F_s = 16 k H z$, frame-width = 25 m s, frame-shift = 10 m s are assumed.

MFCCs show better ASR performance than filter-bank features, but MFCCs are not robust against noises.

Perceptually-based Linear Prediction (PLP)

Using temporal features: dynamic features

In SP lab-sessions on speech recognition using HTK,

- MFCCs, and energy
- Δ MFCCs, Δ energy
- Δ^2 MFCCs, Δ^2 energy
- $\Rightarrow \Delta *, \Delta^2 *:$ delta features

(dynamic features / time derivatives) [Furui, 1986]


```
ASR (H. Shimodaira)
```

I:38

Other features with low dimensionality

Formants (F_1, F_2, F_3, \cdots)

They are not used in modern ASR systems, but why?

Using temporal features: dynamic features(cont. 2)

Using temporal features: dynamic features(cont. 3)

- An acoustic feature vector, eg MFCCs, representing part of a speech signal is highly correlated with its neighbours.
- HMM based acoustic models assume there is no dependency between the observations.
- Those correlations can be captured to some extent by augmenting the original set of static acoustic features, eg. MFCCs, with dynamic features.

ASR	(H.	Shimodaira)
-----	-----	-------------

I:40

SUMMARY(cont. 2)

- Front-end analysis has a great influence on ASR performance.
- For robust ASR in real environments, various techniques for front-end processing have been proposed. e.g. spectral subtraction (SS), cepstral mean normalisation (CMN)
- **Do not believe what you've got in spectral analysis.**

You are not seeing the true one.

You are looking at speech signals / features through a pin hole.

- sampled
- windowed

ASR (H. Shimodaira)

I:42

SUMMARY

- Nyquist Sampling theory
- Short-time Spectrum Analysis
 - Non-parametric method
 - Short-time Fourier Transform
 - Cepstrum, MFCC
 - Filter bank
 - Parametric methods
 - LPC, PLP
 - Windowing effect: trade-off between time and frequency resolutions
- **Dynamic features (delta features)**
- There is no best feature that can be used for any purposes, but MFCC is widely used for ASR and TTS.

References

- John N. Holmes, Wendy J. Holmes, "Speech Synthesis and Recognition", Taylor and Francis (2001), 2nd edition (chapter 2, 4, 10)
- http://mi.eng.cam.ac.uk/comp.speech/
- http://mi.eng.cam.ac.uk/~ajr/SpeechAnalysis/
- http://cslu.cse.ogi.edu/HLTsurvey/
- B. Gold, N. Morgan, "Speech and Audio Signal Processing: Processing and Perception of Speech and Music", John Wiley and Sons (1999).
- "Spoken language processing: a guide to theory, algorithm, and system development", Xuedong Huang, Alex Acero and Hsiao-Wuen Hon, Prentice Hall (2001). isbn: 0130226165

References(cont. 2)

- "Robusness in Automatic Speech Recognition", J-C Junqua and J-P Hanton, , Kluwer Academic Publications (1996). isbn: 0-7923-9646-4
- "A Comparative Study of Traditional and Newly Proposed Features for Recognition of Speech Under Stress", Sahar Bou-Ghazale and John H.L. Hansen, IEEE Trans SAP, vol. 8, no. 4, pp.429–442, July 2000.

I:44