

Lecture1 Introduction

2102874 Speech Processing

Dr Nisachon Tangsangiumvisai

Course Information

Lecturer Dr Nisachon Tangsangiumvisai

ดร. นิสาชณ ตั้งเสงี่ยมวิสัย

Contact

office 02-2186909

E-mail : Nisachon.T@Chula.ac.th

webpage : www.ee.eng.chula.ac.th/~ntang

Course Information (II)

Time Wednesday 9.00 – 12.00 pm

Schedule as given in the course outline

Evaluation Assignments 30 %

Mid-term Examination 40 %

Project 30 %



Books

■ **Speech and Audio Signal Processing :
Processing and perception of speech and music**
B. Gold and N. Morgan
Wiley, 2000.

■ **Digital Speech : Coding for low bit rate
communication systems**
A. M. Kondoz
Wiley, 2000.

Course Contents

Part I : Speech Processing

- Fundamentals of the course
- Speech Modeling
- Speech Analysis and Synthesis
- Speech Coding
- Human Recognition
- Speaker Verification
- Text-to-Speech Synthesis
- Speech Prosody

Course Contents (II)

Part II : Speech Enhancement

- Noise Reduction Techniques

Overview

Speech Communication

■ What is speech communication ?

The transfer of information from one person to another via *speech*, which consists of variations in pressure coming from the mouth of a speaker.

■ What is sound ?

Sound is an acoustic *wave* that results when a vibrating source (e.g vocal cords) disturbs an elastic medium (e.g air).

Speech Communication (II)

- When a sound wave reaches a listener's ear drum, the vibrations are transmitted to the inner ear (or cochlea), where mechanical displacements are converted to neural pulses that are sent to the brain and result in the sensation of sound.
- **Speech chain** consists of speech mechanism in the speaker, transmission through medium, and a speech perception process in the ear and brain of the listener.
- In many applications of speech processing, part of the speech chain is implemented by a simulation device.

Speech Processing

- **Speech Analysis**
 - Speech Recognition
 - Speaker Verification
 - Speaker Identification
- **Speech Synthesis**
 - Text-to-speech (TTS)
- **Speech Codec**

Factors

- **Analogue** form of Speech for Telecommunication
 - transmission power
 - spectral utilization
- **Digital** Transmission of Speech
 - low cost
 - consistent quality
 - security
 - spectral efficiency

Digitization of Speech Signals

- **Sampling** : Nyquist Criterion
- **Quantization** : Number of quantizer levels is proportional to reconstruction at the receiver
- **Speech Transmission Systems**
 - In mobile communication systems
 - ⇒ for reduction in PCM bit rate
 - (64 kbps; restricts the desired spectral efficiency)

Digitization of Speech Signals (II)

■ Digital Coding of Speech Signals :

The parameters for speech production (sensitive to corruption) are encoded and transmitted through degraded channels.

→ problem in maintaining speech quality

■ Bit Rate Reduction :

- difficult to maintain speech quality as the bit rate falls
- the algorithm becomes complex; not practical for real-time implementation
- results in excessive delay → echo control problem

Digitization of Speech Signals (III)

■ With the use of DSP chips (faster & more reliable)

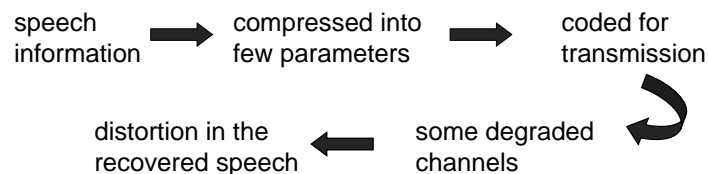
→ possible for real-time implementation of highly complex algorithms

■ Digitally encoded Speech :

- condense down to a binary sequence
 - advantages of digital systems can be exploited (flexibility, security, integration into Integrated Service Digital Networks (ISDN))
- BUT extra BW for transmission → signal compression

Speech Coding

■ Low bit rate (LBR) speech coders



- To identify the sensitivity of each speech parameters for error control at the decoder (For power efficiency : only the most sensitive bits are **error protected**)

Audio Coding for Wireless Communication

Criteria to be considered :

- High quality digital voice
- Bit rate as low as 4 kbps (due to low-rate applications e.g. cellular phones)
- Low delay speech coding at 16 kbps (standard CCITT G.728)
- Implementational Complexity

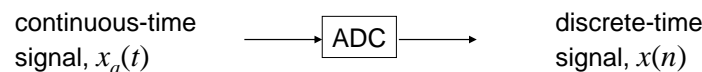
Let's Begin!



Fundamentals of the course

- Digital Signals
- Spectrum Analysis of Digital Signals
 - *Fourier Transform*
 - *Z-transform*
- Digital Systems
- Digital Filters

Digital Signals



Analog-to-Digital Converter (ADC) performs :

- *sampling* of the amplitudes of the analog signal $x_a(t)$ on an equidistant grid along the horizontal time axis

- *quantization* of the amplitudes to fixed samples represented by numbers $x(n)$ along the vertical amplitude axis

Sampling

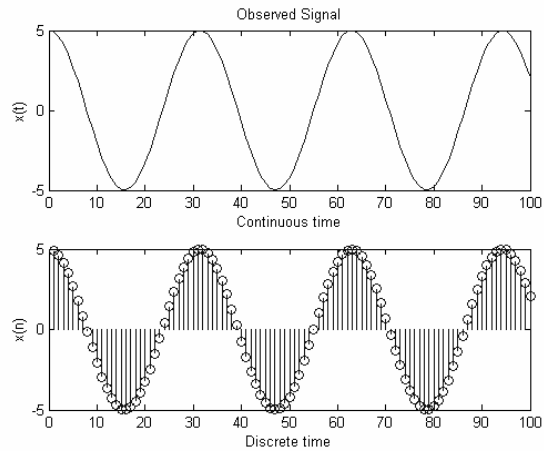
$$x(n) = x_a(nT), \quad -\infty < n < \infty \quad \dots(1)$$

↑ sampling time
↓ analogue signal
Integer (discrete time index)

To avoid aliasing

- input has to be band-limited
- sampling frequency must satisfy the nyquist criterion

Example of signals



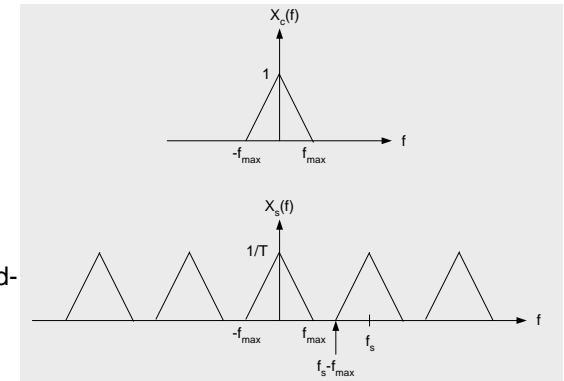
Sampling (II)

$$f_s \geq 2 f_{\max} \quad \dots(2)$$

$$(T \leq \frac{1}{2} T_{\max}) \quad \dots(3)$$

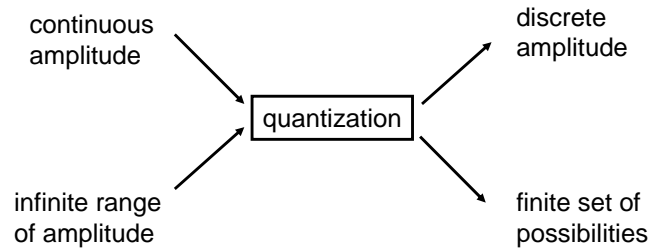
In telecommunication networks, analogue speech signals are band-limited to 300-3400 Hz,

$$f_s = 8 \text{ kHz}$$



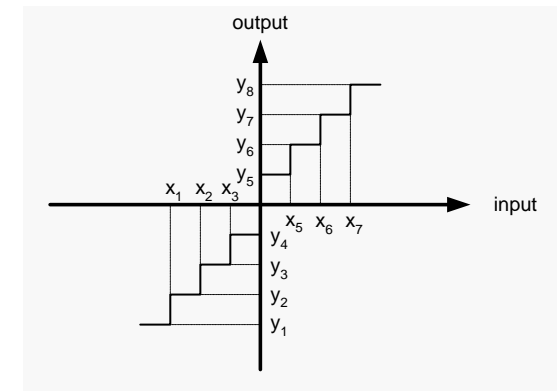
Effect of sampling

Quantization



Scalar Quantization

Each of a set of discrete values is quantized separately.



Characteristics of uniform quantization

Scalar Quantization (II)

- Quantisation step-size (Δ)
: distance between the finite sets of amplitude levels
- Each Δ_n is represented by a code-word for transmission.
- At digital receiver \rightarrow De-quantizer
: indicate which discrete amplitude to be used.
- Channel transmission bit rate :

$$T_c = B f_s \quad (\text{bit/sec}) \quad \dots(4)$$

B : number of bits representing the discrete amplitudes. (length of $c(n)$)

Scalar Quantization (III)

- From Eq.(4), if f_s is fixed,
 T_c can be reduced by reducing B
- $\downarrow B \rightarrow \uparrow \Delta \rightarrow \downarrow$ quality of reconstructed signal

- In the i^{th} interval,

$$x_i = \frac{\Delta_i}{2} \leq x_a(nT) < x_i + \frac{\Delta_i}{2} \quad \dots(5)$$

$x_a(nT)$ is represented by the quantized amplitude.

Scalar Quantization (IV)

- The quantization of the amplitudes to fixed numbers in the range between $-32768 \dots 32767$ is based on a 16-bit representation of the sample amplitudes.
- It allows 2^{16} quantized values in the range of -2^{15} to $2^{15}-1$.
- Thus, for B -bit representation, the number range is $-2^{(B-1)}$ to $2^{(B-1)}-1$.

Scalar Quantization (V)

- Instantaneous squared error is

$$(x_a(nT) - x_i)^2 \quad \dots(6)$$

- Mean Squared Error of the signal

$$E_i^2 = \int_{x_i - \frac{\Delta_i}{2}}^{x_i + \frac{\Delta_i}{2}} (\boxed{x} - x_i)^2 \underbrace{p(x) dx}_{\text{p.d.f. of } x} \quad \dots(7)$$

$x = x_a(nT)$

- Assume fine resolution (small Δ_i), $p(x)$ is flat within $[x_i - \frac{\Delta_i}{2}, x_i + \frac{\Delta_i}{2}] \rightarrow$ use its center value $p(x_i)$

Scalar Quantization (VI)

- Hence, MSE becomes

$$E_i^2 = \frac{\Delta_i^3}{12} p(x_i) \quad \dots(8)$$

- Prob. of signal falling in this interval is

$$\Gamma_i = \int_{x_i - \frac{\Delta_i}{2}}^{x_i + \frac{\Delta_i}{2}} p(x) dx = p(x_i) \Delta_i \quad \dots(9)$$

- Substituting for $p(x_i)$ in Eq.(8), $E_i^2 = \frac{\Delta_i^2}{12} \Gamma_i \quad \dots(10)$

- Hence, total MSE is $E^2 = \frac{1}{12} \sum_{i=1}^N \Gamma_i \Delta_i^2 \quad \dots(11)$
(N = no. of levels in the quantizer)

Scalar Quantization (VII)

- Assume uniform quantisation step-size

$$E^2 = \frac{\Delta^2}{12} \quad \dots(12)$$

- For B-bit binary code words

$$N = 2^B \quad \dots(13)$$

- Assuming $|x| \leq X_{\max}$ and $p(x)$ is symmetrical

$$2X_{\max} = \Delta 2^B \quad \dots(14)$$

Scalar Quantization (VIII)

- Thus, the step-size can be found from

$$\frac{2X_{\max}}{2^B} = \Delta \quad \dots(15)$$

- Quantization error $e_q(n)$ is bounded by

$$-\frac{\Delta}{2} \leq e_q(n) \leq \frac{\Delta}{2} \quad \dots(16)$$

- Uniform quantizer assumes uniform pdf. of constant height $\frac{1}{2X_{\max}}$.

Scalar Quantization (IX)

- Hence, input power is given by

$$P_x = \int_{-X_{\max}}^{X_{\max}} \frac{x^2}{2X_{\max}} p(x) dx \quad \dots(17)$$

- SNR

$$\frac{P_x}{P_n} = \frac{X_{\max}^2 / 3}{\Delta^2 / 12} = 2^{2B} \rightarrow 6.02B \text{ (dB)} \quad \dots(18)$$

Vector Quantization

To characterize the spectrum of a speech signal (**spectral analysis**), consider the following example :

Uncompressed Signal

- 10 kHz sampled speech with 16-bit speech amplitudes
- Information rate is 160,000 bps (required for storage of speech analysis)

Compressed Signal

- spectral vector \mathbf{v}_l , $l=1,2,\dots,L$
- vectors of dimension $p=10$ using 100 spectral vectors/s.
- by representing each spectral component to 16-bit precision, the required storage is 100x10x16 bps

➡ **10 times reduction**

Vector Quantization (II)

It is needed a **single** spectral representation for each basic speech unit.

- ➡ A finite number of **unique** spectral vectors. (each vector corresponds to one of the basic speech units.)
- ➡ impossible, due to time-varying properties of the spectra of the signal.
- ➡ build a codebook of **distinct** analysis vectors. (common techniques for vector quantization (VQ) methods)

Vector Quantization (III)

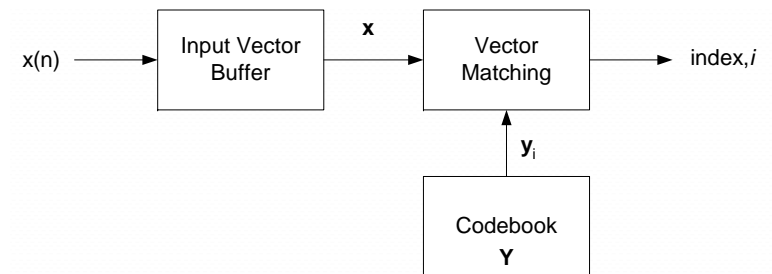
◆ So, if we require to a codebook with 1024 unique spectral vectors, then we need a 10-bit number.

◆ Assume a rate of 100 spectral vectors/s, a total bit rate of 1 kbps is required to represent the spectral vectors of a speech signal. (this is 1/16 time the rate required by the continuous spectral vectors.)

◆ So, the VQ method is an **efficient** representation of the spectral information in the speech signal.

Vector Quantization (IV)

- ◆ Block Quantization / Pattern Matching Quantization
- ◆ Signals are quantized as a single vector.



Block diagram showing a vector quantization

Vector Quantization (V)

Advantages

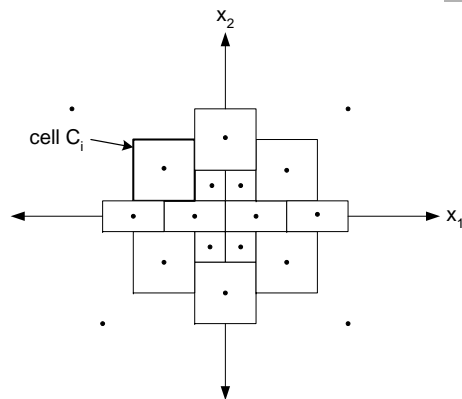
- ♦ reduced storage
- ♦ reduced computation for determining similarity of spectral analysis vectors (use of a table lookup)
- ♦ discrete representation of speech sounds (use of codebook)

Vector Quantization (VI)

Disadvantages

- ♦ since there is a *finite* number of codebook vectors, the “**best**” representation of a given spectral vector is obtained with a certain level (non-zero) of quantization error.
(As the size of the codebook increases, the size of the quantization error decreases.)
- ♦ To reduce the quantization error, larger codebook is needed ➡ more storage!!

Vector Quantization (VII)



Partitioning a 2-dimensional space into 18 cells

Comparison

■ Scalar Quantization

- different cells have same shape

■ Vector Quantization

- different cells have different shapes
- at very low bit rate, performs better than the scalar method, but at computational and storage costs.

Comparison (II)

Vector quantization systems tend to be less robust to random channel errors than scalar method.

E.g. 10-bit vector quantizer
10 1-bit scalar quantizer

Assume channel error rate that causes 1-bit error, thus
scalar quantizer : 1-bit in error \rightarrow affect 1 value in 1 dimension
vector quantizer : 1-bit in error \rightarrow affect 10-bit vector

Signal Reconstruction

discrete-time signal, $x(n)$ \rightarrow **DAC** \rightarrow continuous-time signal, $x_a(t)$

Digital-to-Analog Converter is used for reconstruction of analog signal.

Spectrum Analysis of Digital Signals

■ Discrete-time Fourier Transform (DTFT)

$$X(e^{j\omega}) = F[x(n)] = \sum_{n=-\infty}^{\infty} x(n)e^{-j\omega n} \quad \dots (19)$$

a function of continuous ω

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega}) e^{jn\omega} d\omega \quad \dots (20)$$

Spectrum Analysis (II)

■ Discrete Fourier Transform (DFT)

$$X(k) = \sum_{n=0}^{N-1} x(n)e^{-j2\pi kn/N}, \quad k = 0, 1, \dots, N-1 \quad \dots (21)$$

for digital computation

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k)e^{j2\pi kn/N}, \quad n = 0, 1, \dots, N-1 \quad \dots (22)$$

Spectrum Analysis (III)

- The fast version of the DFT is called the Fast Fourier Transform (FFT).

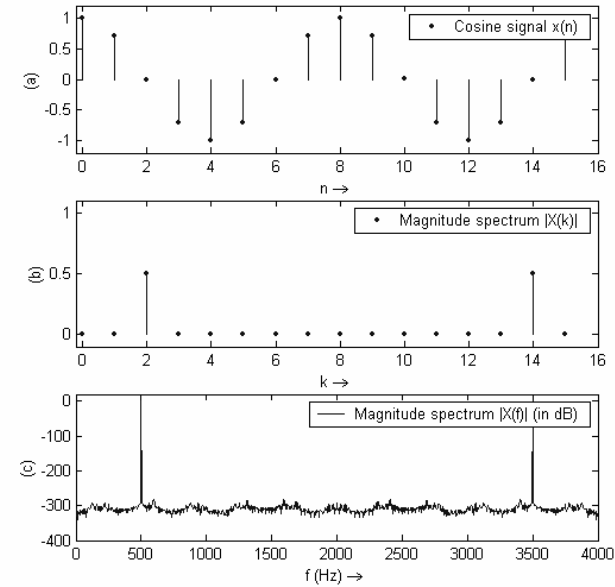
>> Xmag = abs(fft(x,N)); Xphase = angle(fft(x,N));

- From Eq.(21),

$$X(k) = X_R(k) + jX_I(k), \quad k = 0, 1, \dots, N-1$$

Magnitude spectrum : $|X(k)| = \sqrt{X_R^2(k) + X_I^2(k)}$... (23)

Phase : $\varphi(k) = \arctan\left(\frac{X_I(k)}{X_R(k)}\right)$... (24)

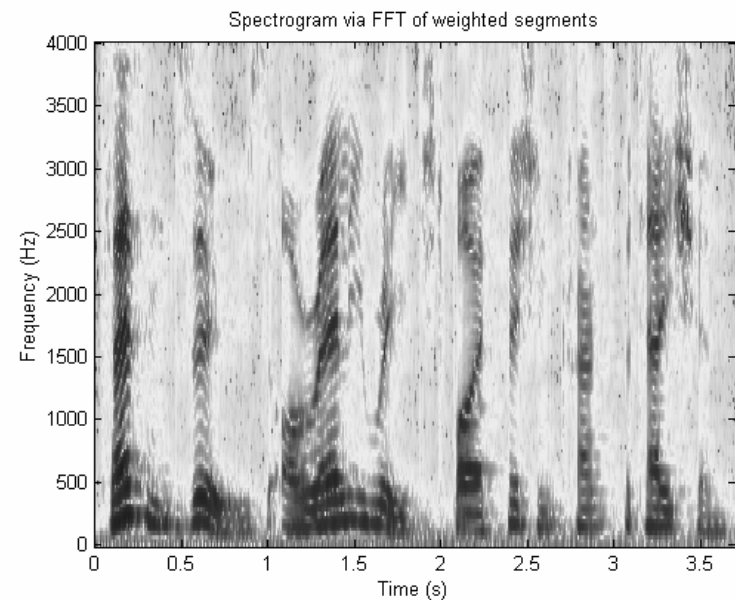


Spectrum Analysis (V)

Time-frequency representation → Spectrogram

- It is an estimate of the short-time, time-localized frequency content of the signal.
- The signal is split into segments of length N
- and are multiplied by a window
- and an FFT is performed
- An overlap of the weighted segment is used to increase the time-localization of the short-time spectra.

>> [B,F,T] = spectrogram(x,NFFT,Fs,WINDOW,NOVERLAP);



Spectrum Analysis (VII)

■ z-transform

$$X(z) = Z[x(n)] \quad \dots (25)$$

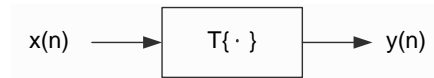
$$= \sum_{n=-\infty}^{\infty} x(n) z^{-n} \quad \dots (26)$$

therefore

$$x(n) = \frac{1}{2\pi j} \oint_C X(z) z^{n-1} dz \quad \dots (27)$$

$$X(z) \Big|_{z=e^{j\omega}} = X(e^{j\omega}) = \sum_{n=-\infty}^{\infty} x(n) e^{-j\omega n} \quad \dots (28)$$

Digital Systems



$$y(n) = T\{x(n)\} \quad \dots(29)$$

Properties of discrete-time systems:

- linearity
- shift-invariance
- causality
- invertibility

Digital Systems (II)

■ Linearity

$$T\{ax_1(n) + bx_2(n)\} = aT\{x_1(n)\} + bT\{x_2(n)\} \quad \dots(30)$$

a, b are constants.

■ Shift-invariance

$$\text{If } x(n) \rightarrow y(n) \quad \dots(31)$$

$$\text{then } x(n - n_0) \rightarrow y(n - n_0) \quad \dots(32)$$

Digital Systems (III)

■ Causality

For any n_0 , the response of the system at time n_0 depends only upon the values of the input for $n < n_0$.

Examples :

$$y(n) = x(n) + x(n-1) \quad \dots(33)$$

$$y(n) = x(n) + x(n+1) \quad \dots(34)$$

■ Invertibility

The input to the system maybe uniquely determined by observing the output.

Digital Systems (IV)

■ Invertibility (contd.)

$$x_1(n) \leftrightarrow y_1(n)$$

$$x_2(n) \leftrightarrow y_2(n)$$

Examples :

when $x_1(n) \neq x_2(n)$ and $y_1(n) \neq y_2(n)$

If $y(n) = x(n)h(n)$... (35)

→ $x(n) = y(n)/h(n)$... (36)

Digital Systems (V)

■ Linear-time-invariant (LTI) systems

■ Discrete Convolution

$x(n)$: input signal

$h(n)$: impulse response of a digital system

$y(n)$: output signal

$$y(n) = x(n) * h(n) = \sum_{k=-\infty}^{\infty} x(k)h(n-k) \quad \dots(37)$$

>> $y = \text{conv}(x,h);$

Digital Systems (VI)

■ FIR system : a system with a finite impulse response

$$H(z) = \frac{Y(z)}{X(z)} = 1 + b_1z^{-1} + b_2z^{-2} \quad \dots(38)$$

■ IIR system : a system with an infinite impulse response

$$H(z) = \frac{Y(z)}{X(z)} = \frac{1}{1 + a_1z^{-1} + a_2z^{-2}} \quad \dots(39)$$

>> $y = \text{filter}(B,A,x);$

>> $[H,W] = \text{freqz}(B,A,N);$

Digital Filters

■ Lowpass filter (LPF)

- select low frequencies up to the cutoff frequency f_c
- Attenuate frequencies higher than f_c

■ Highpass filter (HPF)

- select frequencies higher than f_c
- attenuate frequencies below f_c

■ Bandpass filter (BPF)

- select frequencies between a lower f_{c1} and a higher f_{c2}
- attenuate frequencies outside this range

Digital Filters (II)

- **Bandreject filter**
 - attenuate frequencies between a lower f_{c1} and a higher f_{c2}
 - frequencies outside this range are passed
- **Notch filter**
 - attenuate frequencies in a narrow bandwidth around the cutoff frequency f_c
- **Allpass filter**
 - pass all frequencies but modify the phase of the input signal

Summary

- Introduction to the course
- Fundamentals of the course

- Next Lecture
Lecture 2 : Speech Modeling

