



**ΕΛΛΗΝΙΚΗ ΔΗΜΟΚΡΑΤΙΑ
ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ**

Ψηφιακή Επεξεργασία Φωνής

Ενότητα 6η: Κωδικοποίηση Φωνής

Στυλιανού Ιωάννης

Τμήμα Επιστήμης Υπολογιστών

CS578- SPEECH SIGNAL PROCESSING

LECTURE 7: SPEECH CODING

Yannis Stylianou



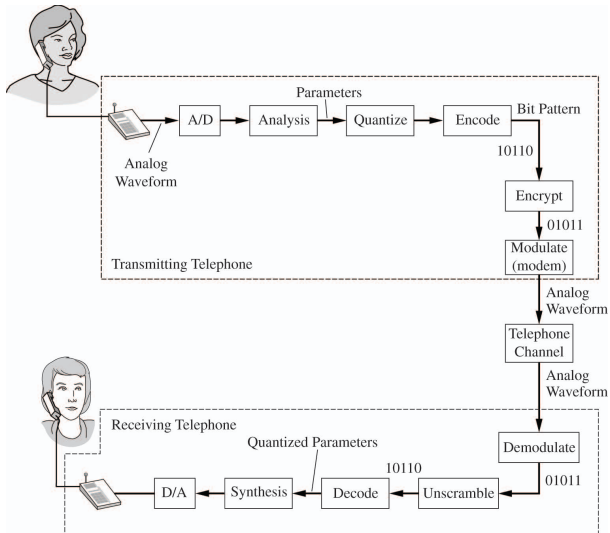
University of Crete, Computer Science Dept., Multimedia Informatics Lab
yannis@csd.uoc.gr

Univ. of Crete

OUTLINE

- 1 INTRODUCTION
- 2 STATISTICAL MODELS
- 3 SCALAR QUANTIZATION
 - Max Quantizer
 - Companding
 - Adaptive quantization
 - Differential and Residual quantization
- 4 VECTOR QUANTIZATION
 - The k-means algorithm
 - The LBG algorithm
- 5 MODEL-BASED CODING
 - Basic Linear Prediction, LPC
 - Mixed Excitation LPC (MELP)
- 6 ACKNOWLEDGMENTS

DIGITAL TELEPHONE COMMUNICATION SYSTEM



CATEGORIES OF SPEECH CODERS

- Waveform coders (16-64 kbps, $f_s = 8000\text{Hz}$)
- Hybrid coders (2.4-16 kbps, $f_s = 8000\text{Hz}$)
- Vocoders (1.2-4.8 kbps, $f_s = 8000\text{Hz}$)

CATEGORIES OF SPEECH CODERS

- Waveform coders (16-64 kbps, $f_s = 8000\text{Hz}$)
- Hybrid coders (2.4-16 kbps, $f_s = 8000\text{Hz}$)
- Vocoder (1.2-4.8 kbps, $f_s = 8000\text{Hz}$)

CATEGORIES OF SPEECH CODERS

- Waveform coders (16-64 kbps, $f_s = 8000\text{Hz}$)
- Hybrid coders (2.4-16 kbps, $f_s = 8000\text{Hz}$)
- Vocoders (1.2-4.8 kbps, $f_s = 8000\text{Hz}$)

- Closeness of the processed speech waveform to the original speech waveform
- Naturalness
- Background artifacts
- Intelligibility
- Speaker identifiability

- Closeness of the processed speech waveform to the original speech waveform
- Naturalness
 - Background artifacts
 - Intelligibility
 - Speaker identifiability

- Closeness of the processed speech waveform to the original speech waveform
- Naturalness
- Background artifacts
- Intelligibility
- Speaker identifiability

- Closeness of the processed speech waveform to the original speech waveform
- Naturalness
- Background artifacts
- Intelligibility
- Speaker identifiability

- Closeness of the processed speech waveform to the original speech waveform
- Naturalness
- Background artifacts
- Intelligibility
- Speaker identifiability

MEASURING SPEECH QUALITY

▷ Subjective tests:

- Diagnostic Rhyme Test (DRT)
- Diagnostic Acceptability Measure (DAM)
- Mean Opinion Score (MOS)

▷ Objective tests:

- Segmental Signal-to-Noise Ratio (SNR)
- Articulation Index

MEASURING SPEECH QUALITY

▷ Subjective tests:

- Diagnostic Rhyme Test (DRT)
- Diagnostic Acceptability Measure (DAM)
- Mean Opinion Score (MOS)

▷ Objective tests:

- Segmental Signal-to-Noise Ratio (SNR)
- Articulation Index

MEASURING SPEECH QUALITY

▷ Subjective tests:

- Diagnostic Rhyme Test (DRT)
- Diagnostic Acceptability Measure (DAM)
- Mean Opinion Score (MOS)

▷ Objective tests:

- Segmental Signal-to-Noise Ratio (SNR)
- Articulation Index

MEASURING SPEECH QUALITY

▷ Subjective tests:

- Diagnostic Rhyme Test (DRT)
- Diagnostic Acceptability Measure (DAM)
- Mean Opinion Score (MOS)

▷ Objective tests:

- Segmental Signal-to-Noise Ratio (SNR)
- Articulation Index

MEASURING SPEECH QUALITY

▷ Subjective tests:

- Diagnostic Rhyme Test (DRT)
- Diagnostic Acceptability Measure (DAM)
- Mean Opinion Score (MOS)

▷ Objective tests:

- Segmental Signal-to-Noise Ratio (SNR)
- Articulation Index

- Statistical models of speech (preliminary)
- Scalar quantization (i.e., waveform coding)
- Vector quantization (i.e., subband and sinusoidal coding)

- Statistical models of speech (preliminary)
- Scalar quantization (i.e., waveform coding)
- Vector quantization (i.e., subband and sinusoidal coding)

- Statistical models of speech (preliminary)
- Scalar quantization (i.e., waveform coding)
- Vector quantization (i.e., subband and sinusoidal coding)

LINEAR PREDICTION CODING, LPC

- Classic LPC
- Mixed Excitation Linear Prediction, MELP
- Multipulse LPC
- Code Excited Linear Prediction (CELP)

LINEAR PREDICTION CODING, LPC

- Classic LPC
- Mixed Excitation Linear Prediction, MELP
- Multipulse LPC
- Code Excited Linear Prediction (CELP)

LINEAR PREDICTION CODING, LPC

- Classic LPC
- Mixed Excitation Linear Prediction, MELP
- Multipulse LPC
- Code Excited Linear Prediction (CELP)

LINEAR PREDICTION CODING, LPC

- Classic LPC
- Mixed Excitation Linear Prediction, MELP
- Multipulse LPC
- Code Excited Linear Prediction (CELP)

OUTLINE

- 1 INTRODUCTION
- 2 STATISTICAL MODELS
- 3 SCALAR QUANTIZATION
 - Max Quantizer
 - Companding
 - Adaptive quantization
 - Differential and Residual quantization
- 4 VECTOR QUANTIZATION
 - The k-means algorithm
 - The LBG algorithm
- 5 MODEL-BASED CODING
 - Basic Linear Prediction, LPC
 - Mixed Excitation LPC (MELP)
- 6 ACKNOWLEDGMENTS

PROBABILITY DENSITY OF SPEECH

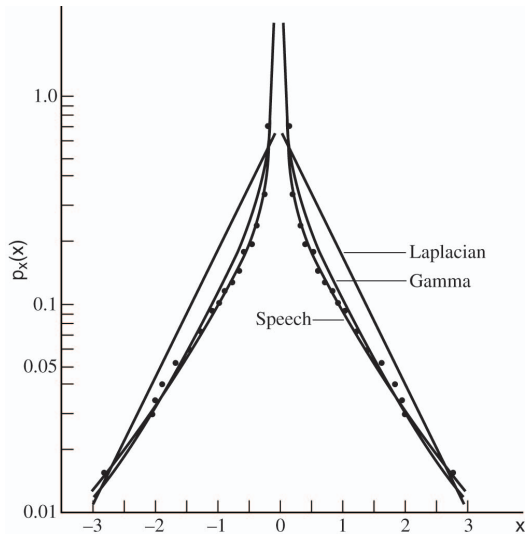
By setting $x[n] \rightarrow x$, the histogram of speech samples can be approximated by a *gamma density*:

$$p_X(x) = \left(\frac{\sqrt{3}}{8\pi\sigma_x|x|} \right)^{1/2} e^{-\frac{\sqrt{3}|x|}{2\sigma_x}}$$

or by a simpler *Laplacian density*:

$$p_X(x) = \frac{1}{\sqrt{2}\sigma_x} e^{-\frac{\sqrt{3}|x|}{\sigma_x}}$$

DENSITIES COMPARISON



OUTLINE

1 INTRODUCTION

2 STATISTICAL MODELS

3 SCALAR QUANTIZATION

- Max Quantizer
- Companding
- Adaptive quantization
- Differential and Residual quantization

4 VECTOR QUANTIZATION

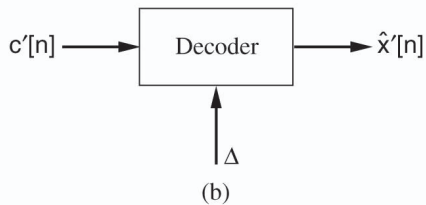
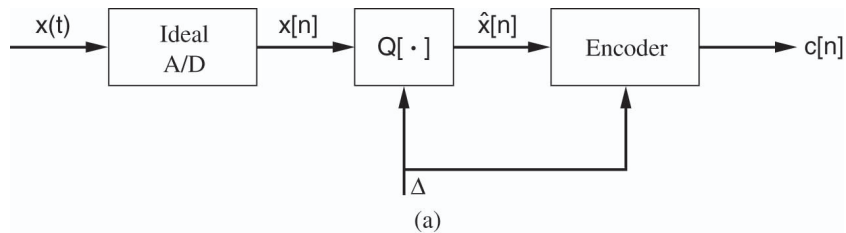
- The k-means algorithm
- The LBG algorithm

5 MODEL-BASED CODING

- Basic Linear Prediction, LPC
- Mixed Excitation LPC (MELP)

6 ACKNOWLEDGMENTS

CODING AND DECODING



FUNDAMENTALS OF SCALAR CODING

- Let's quantize a single sample speech value, $x[n]$ into M *reconstruction* or *decision* levels:

$$\hat{x}[n] = \hat{x}_i = Q(x[n]), \quad x_{i-1} < x[n] \leq x_i$$

with $1 \leq i \leq M$ and x_k denotes the M decision levels with $0 \leq k \leq M$.

- Assign a *codeword* in each reconstruction level. Collection of codewords makes a *codebook*.
- Using B -bit binary codebook we can represent each 2^B different *quantization* (reconstruction) levels.
- *Bit rate*, I , is defined as: $I = Bf_s$

FUNDAMENTALS OF SCALAR CODING

- Let's quantize a single sample speech value, $x[n]$ into M *reconstruction* or *decision* levels:

$$\hat{x}[n] = \hat{x}_i = Q(x[n]), \quad x_{i-1} < x[n] \leq x_i$$

with $1 \leq i \leq M$ and x_k denotes the M decision levels with $0 \leq k \leq M$.

- Assign a *codeword* in each reconstruction level. Collection of codewords makes a *codebook*.
- Using B -bit binary codebook we can represent each 2^B different *quantization* (reconstruction) levels.
- *Bit rate*, I , is defined as: $I = Bf_s$

FUNDAMENTALS OF SCALAR CODING

- Let's quantize a single sample speech value, $x[n]$ into M *reconstruction* or *decision* levels:

$$\hat{x}[n] = \hat{x}_i = Q(x[n]), \quad x_{i-1} < x[n] \leq x_i$$

with $1 \leq i \leq M$ and x_k denotes the M decision levels with $0 \leq k \leq M$.

- Assign a *codeword* in each reconstruction level. Collection of codewords makes a *codebook*.
- Using B -bit binary codebook we can represent each 2^B different *quantization* (reconstruction) levels.
- *Bit rate*, I , is defined as: $I = Bf_s$

FUNDAMENTALS OF SCALAR CODING

- Let's quantize a single sample speech value, $x[n]$ into M *reconstruction* or *decision* levels:

$$\hat{x}[n] = \hat{x}_i = Q(x[n]), \quad x_{i-1} < x[n] \leq x_i$$

with $1 \leq i \leq M$ and x_k denotes the M decision levels with $0 \leq k \leq M$.

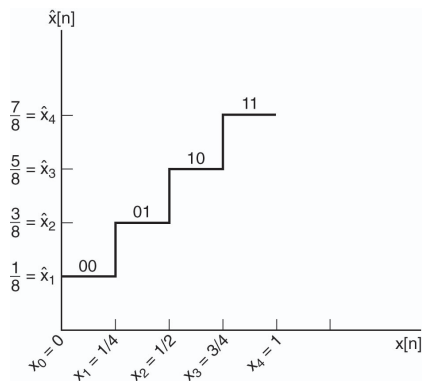
- Assign a *codeword* in each reconstruction level. Collection of codewords makes a *codebook*.
- Using B -bit binary codebook we can represent each 2^B different *quantization* (reconstruction) levels.
- *Bit rate*, I , is defined as: $I = Bf_s$

UNIFORM QUANTIZATION

$$x_i - x_{i-1} = \Delta, \quad 1 \leq i \leq M$$
$$\hat{x}_i = \frac{x_i + x_{i-1}}{2}, \quad 1 \leq i \leq M$$

Δ is referred to as uniform quantization step size.

▷ Example of a 2-bit *uniform quantization*:



UNIFORM QUANTIZATION: DESIGNING DECISION REGIONS

- Signal range: $-4\sigma_x \leq x[n] \leq 4\sigma_x$
- Assuming B-bit binary codebook, we get 2^B quantization (reconstruction) levels
- Quantization step size, Δ :

$$\Delta = \frac{2x_{max}}{2^B}$$

- Δ and *quantization noise*.

CLASSES OF QUANTIZATION NOISE

There are two classes of quantization noise:

- Granular Distortion:

$$\hat{x}[n] = x[n] + e[n]$$

where $e[n]$ is the quantization noise, with:

$$-\frac{\Delta}{2} \leq e[n] \leq \frac{\Delta}{2}$$

- Overload Distortion: *clipped samples*

ASSUMPTIONS

- Quantization noise is an ergodic white-noise random process:

$$\begin{aligned}r_e[m] &= E(e[n]e[n+m]) \\ &= \sigma_e^2, \quad m = 0 \\ &= 0, \quad m \neq 0\end{aligned}$$

- Quantization noise and input signal are uncorrelated:

$$E(x[n]e[n+m]) = 0 \quad \forall m$$

- Quantization noise is uniform over the quantization interval

$$\begin{aligned}p_e(e) &= \frac{1}{\Delta}, \quad -\frac{\Delta}{2} \leq e \leq \frac{\Delta}{2} \\ &= 0, \quad \text{otherwise}\end{aligned}$$

ASSUMPTIONS

- Quantization noise is an ergodic white-noise random process:

$$\begin{aligned}r_e[m] &= E(e[n]e[n+m]) \\ &= \sigma_e^2, \quad m = 0 \\ &= 0, \quad m \neq 0\end{aligned}$$

- Quantization noise and input signal are uncorrelated:

$$E(x[n]e[n+m]) = 0 \quad \forall m$$

- Quantization noise is uniform over the quantization interval

$$\begin{aligned}p_e(e) &= \frac{1}{\Delta}, \quad -\frac{\Delta}{2} \leq e \leq \frac{\Delta}{2} \\ &= 0, \quad \text{otherwise}\end{aligned}$$

ASSUMPTIONS

- Quantization noise is an ergodic white-noise random process:

$$\begin{aligned} r_e[m] &= E(e[n]e[n+m]) \\ &= \sigma_e^2, \quad m = 0 \\ &= 0, \quad m \neq 0 \end{aligned}$$

- Quantization noise and input signal are uncorrelated:

$$E(x[n]e[n+m]) = 0 \quad \forall m$$

- Quantization noise is uniform over the quantization interval

$$\begin{aligned} p_e(e) &= \frac{1}{\Delta}, \quad -\frac{\Delta}{2} \leq e \leq \frac{\Delta}{2} \\ &= 0, \quad \text{otherwise} \end{aligned}$$

DEFINITION (DITHERING)

We can force $e[n]$ to be white and uncorrelated with $x[n]$ by adding noise to $x[n]$ before quantization!

SIGNAL-TO-NOISE RATIO

- To quantify the severity of the quantization noise, we define the *Signal-to-Noise Ratio* (SNR) as:

$$\begin{aligned} SNR &= \frac{\sigma_x^2}{\sigma_e^2} \\ &= \frac{E(x^2[n])}{E(e^2[n])} \\ &\approx \frac{\frac{1}{N} \sum_{n=0}^{N-1} x^2[n]}{\frac{1}{N} \sum_{n=0}^{N-1} e^2[n]} \end{aligned}$$

- For uniform pdf and quantizer range $2x_{max}$:

$$\begin{aligned} \sigma_e^2 &= \frac{\Delta^2}{12} \\ &= \frac{x_{max}^2}{3 \cdot 2^{2B}} \end{aligned}$$

- Or

$$SNR = \frac{3 \cdot 2^{2B}}{\left(\frac{x_{max}}{\sigma_x}\right)^2}$$

- and in dB:

$$SNR(dB) \approx 6B + 4.77 - 20 \log_{10} \left(\frac{x_{max}}{\sigma_x} \right)$$

- and since $x_{max} = 4\sigma_x$:

$$SNR(dB) \approx 6B - 7.2$$

SIGNAL-TO-NOISE RATIO

- To quantify the severity of the quantization noise, we define the *Signal-to-Noise Ratio* (SNR) as:

$$\begin{aligned} SNR &= \frac{\sigma_x^2}{\sigma_e^2} \\ &= \frac{E(x^2[n])}{E(e^2[n])} \\ &\approx \frac{\frac{1}{N} \sum_{n=0}^{N-1} x^2[n]}{\frac{1}{N} \sum_{n=0}^{N-1} e^2[n]} \end{aligned}$$

- For uniform pdf and quantizer range $2x_{max}$:

$$\begin{aligned} \sigma_e^2 &= \frac{\Delta^2}{12} \\ &= \frac{x_{max}^2}{3 \cdot 2^{2B}} \end{aligned}$$

- Or

$$SNR = \frac{3 \cdot 2^{2B}}{\left(\frac{x_{max}}{\sigma_x}\right)^2}$$

- and in dB:

$$SNR(dB) \approx 6B + 4.77 - 20 \log_{10} \left(\frac{x_{max}}{\sigma_x} \right)$$

- and since $x_{max} = 4\sigma_x$:

$$SNR(dB) \approx 6B - 7.2$$

SIGNAL-TO-NOISE RATIO

- To quantify the severity of the quantization noise, we define the *Signal-to-Noise Ratio* (SNR) as:

$$\begin{aligned} SNR &= \frac{\sigma_x^2}{\sigma_e^2} \\ &= \frac{E(x^2[n])}{E(e^2[n])} \\ &\approx \frac{\frac{1}{N} \sum_{n=0}^{N-1} x^2[n]}{\frac{1}{N} \sum_{n=0}^{N-1} e^2[n]} \end{aligned}$$

- For uniform pdf and quantizer range $2x_{max}$:

$$\begin{aligned} \sigma_e^2 &= \frac{\Delta^2}{12} \\ &= \frac{x_{max}^2}{3 \cdot 2^{2B}} \end{aligned}$$

- Or

$$SNR = \frac{3 \cdot 2^{2B}}{\left(\frac{x_{max}}{\sigma_x}\right)^2}$$

- and in dB:

$$SNR(dB) \approx 6B + 4.77 - 20 \log_{10} \left(\frac{x_{max}}{\sigma_x} \right)$$

- and since $x_{max} = 4\sigma_x$:

$$SNR(dB) \approx 6B - 7.2$$

SIGNAL-TO-NOISE RATIO

- To quantify the severity of the quantization noise, we define the *Signal-to-Noise Ratio* (SNR) as:

$$\begin{aligned} SNR &= \frac{\sigma_x^2}{\sigma_e^2} \\ &= \frac{E(x^2[n])}{E(e^2[n])} \\ &\approx \frac{\frac{1}{N} \sum_{n=0}^{N-1} x^2[n]}{\frac{1}{N} \sum_{n=0}^{N-1} e^2[n]} \end{aligned}$$

- For uniform pdf and quantizer range $2x_{max}$:

$$\begin{aligned} \sigma_e^2 &= \frac{\Delta^2}{12} \\ &= \frac{x_{max}^2}{3 \cdot 2^{2B}} \end{aligned}$$

- Or

$$SNR = \frac{3 \cdot 2^{2B}}{\left(\frac{x_{max}}{\sigma_x}\right)^2}$$

- and in dB:

$$SNR(dB) \approx 6B + 4.77 - 20 \log_{10} \left(\frac{x_{max}}{\sigma_x} \right)$$

- and since $x_{max} = 4\sigma_x$:

$$SNR(dB) \approx 6B - 7.2$$

SIGNAL-TO-NOISE RATIO

- To quantify the severity of the quantization noise, we define the *Signal-to-Noise Ratio* (SNR) as:

$$\begin{aligned} SNR &= \frac{\sigma_x^2}{\sigma_e^2} \\ &= \frac{E(x^2[n])}{E(e^2[n])} \\ &\approx \frac{\frac{1}{N} \sum_{n=0}^{N-1} x^2[n]}{\frac{1}{N} \sum_{n=0}^{N-1} e^2[n]} \end{aligned}$$

- For uniform pdf and quantizer range $2x_{max}$:

$$\begin{aligned} \sigma_e^2 &= \frac{\Delta^2}{12} \\ &= \frac{x_{max}^2}{3 \cdot 2^{2B}} \end{aligned}$$

- Or

$$SNR = \frac{3 \cdot 2^{2B}}{\left(\frac{x_{max}}{\sigma_x}\right)^2}$$

- and in dB:

$$SNR(dB) \approx 6B + 4.77 - 20 \log_{10} \left(\frac{x_{max}}{\sigma_x} \right)$$

- and since $x_{max} = 4\sigma_x$:

$$SNR(dB) \approx 6B - 7.2$$

PULSE CODE MODULATION, PCM

- B bits of information per sample are transmitted as a codeword
- instantaneous coding
- not signal-specific
- 11 bits are required for “toll quality”
- what is the rate for $f_s = 10\text{kHz}$?
- For CD, with $f_s = 44100$ and $B = 16$ (16-bit PCM), what is the SNR?

PULSE CODE MODULATION, PCM

- B bits of information per sample are transmitted as a codeword
- instantaneous coding
 - not signal-specific
 - 11 bits are required for “toll quality”
 - what is the rate for $f_s = 10\text{kHz}$?
 - For CD, with $f_s = 44100$ and $B = 16$ (16-bit PCM), what is the SNR?

PULSE CODE MODULATION, PCM

- B bits of information per sample are transmitted as a codeword
- instantaneous coding
- not signal-specific
- 11 bits are required for “toll quality”
- what is the rate for $f_s = 10\text{kHz}$?
- For CD, with $f_s = 44100$ and $B = 16$ (16-bit PCM), what is the SNR?

PULSE CODE MODULATION, PCM

- B bits of information per sample are transmitted as a codeword
- instantaneous coding
- not signal-specific
- 11 bits are required for “toll quality”
- what is the rate for $f_s = 10\text{kHz}$?
- For CD, with $f_s = 44100$ and $B = 16$ (16-bit PCM), what is the SNR?

PULSE CODE MODULATION, PCM

- B bits of information per sample are transmitted as a codeword
- instantaneous coding
- not signal-specific
- 11 bits are required for “toll quality”
- what is the rate for $f_s = 10\text{kHz}$?
- For CD, with $f_s = 44100$ and $B = 16$ (16-bit PCM), what is the SNR?

PULSE CODE MODULATION, PCM

- B bits of information per sample are transmitted as a codeword
- instantaneous coding
- not signal-specific
- 11 bits are required for “toll quality”
- what is the rate for $f_s = 10\text{kHz}$?
- For CD, with $f_s = 44100$ and $B = 16$ (16-bit PCM), what is the SNR?

OPTIMAL DECISION AND RECONSTRUCTION LEVEL

if $x[n] \mapsto p_x(x)$ we determine the optimal decision level, x_i and the reconstruction level, \hat{x} , by minimizing:

$$\begin{aligned} D &= E[(\hat{x} - x)^2] \\ &= \int_{-\infty}^{\infty} p_x(x)(\hat{x} - x)^2 dx \end{aligned}$$

and assuming M reconstruction levels $\hat{x} = Q[x]$:

$$D = \sum_{i=1}^M \int_{x_{i-1}}^{x_i} p_x(x)(\hat{x}_i - x)^2 dx$$

So:

$$\begin{aligned} \frac{\partial D}{\partial \hat{x}_k} &= 0, \quad 1 \leq k \leq M \\ \frac{\partial D}{\partial x_k} &= 0, \quad 1 \leq k \leq M - 1 \end{aligned}$$

OPTIMAL DECISION AND RECONSTRUCTION LEVEL

if $x[n] \mapsto p_x(x)$ we determine the optimal decision level, x_i and the reconstruction level, \hat{x} , by minimizing:

$$\begin{aligned} D &= E[(\hat{x} - x)^2] \\ &= \int_{-\infty}^{\infty} p_x(x)(\hat{x} - x)^2 dx \end{aligned}$$

and assuming M reconstruction levels $\hat{x} = Q[x]$:

$$D = \sum_{i=1}^M \int_{x_{i-1}}^{x_i} p_x(x)(\hat{x}_i - x)^2 dx$$

So:

$$\begin{aligned} \frac{\partial D}{\partial \hat{x}_k} &= 0, \quad 1 \leq k \leq M \\ \frac{\partial D}{\partial x_k} &= 0, \quad 1 \leq k \leq M - 1 \end{aligned}$$

OPTIMAL DECISION AND RECONSTRUCTION LEVEL

if $x[n] \mapsto p_x(x)$ we determine the optimal decision level, x_i and the reconstruction level, \hat{x} , by minimizing:

$$\begin{aligned} D &= E[(\hat{x} - x)^2] \\ &= \int_{-\infty}^{\infty} p_x(x)(\hat{x} - x)^2 dx \end{aligned}$$

and assuming M reconstruction levels $\hat{x} = Q[x]$:

$$D = \sum_{i=1}^M \int_{x_{i-1}}^{x_i} p_x(x)(\hat{x}_i - x)^2 dx$$

So:

$$\begin{aligned} \frac{\partial D}{\partial \hat{x}_k} &= 0, \quad 1 \leq k \leq M \\ \frac{\partial D}{\partial x_k} &= 0, \quad 1 \leq k \leq M - 1 \end{aligned}$$

OPTIMAL DECISION AND RECONSTRUCTION LEVEL, *cont.*

- The minimization of D over decision level, x_k , gives:

$$x_k = \frac{\hat{x}_{k+1} + \hat{x}_k}{2}, \quad 1 \leq k \leq M - 1$$

- The minimization of D over reconstruction level, \hat{x}_k , gives:

$$\begin{aligned}\hat{x}_k &= \int_{x_{k-1}}^{x_k} \left[\frac{p_x(x)}{\int_{x_{k-1}}^{x_k} p_x(s) ds} \right] x dx \\ &= \int_{x_{k-1}}^{x_k} \tilde{p}_x(x) x dx\end{aligned}$$

OPTIMAL DECISION AND RECONSTRUCTION LEVEL, *cont.*

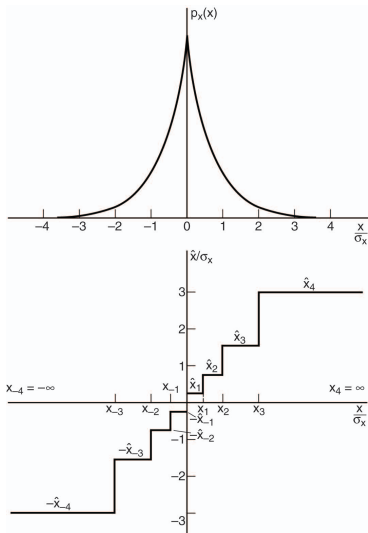
- The minimization of D over decision level, x_k , gives:

$$x_k = \frac{\hat{x}_{k+1} + \hat{x}_k}{2}, \quad 1 \leq k \leq M - 1$$

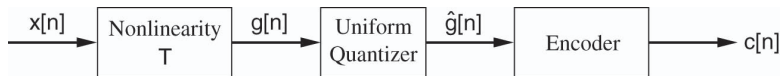
- The minimization of D over reconstruction level, \hat{x}_k , gives:

$$\begin{aligned}\hat{x}_k &= \int_{x_{k-1}}^{x_k} \left[\frac{p_x(x)}{\int_{x_{k-1}}^{x_k} p_x(s) ds} \right] x dx \\ &= \int_{x_{k-1}}^{x_k} \tilde{p}_x(x) x dx\end{aligned}$$

EXAMPLE WITH LAPLACIAN PDF



PRINCIPLE OF COMPANDING



(a)



(b)

COMPANDING EXAMPLES

Companding examples:

- Transformation to a uniform density:

$$g[n] = T(x[n]) = \int_{-\infty}^{x[n]} p_x(s) ds - \frac{1}{2}, \quad \frac{-1}{2} \leq g[n] \leq \frac{1}{2}$$
$$= 0 \quad \text{elsewhere}$$

- μ -law:

$$T(x[n]) = x_{max} \frac{\log(1 + \mu \frac{|x[n]|}{x_{max}})}{\log(1 + \mu)} \text{sign}(x[n])$$

COMPANDING EXAMPLES

Companding examples:

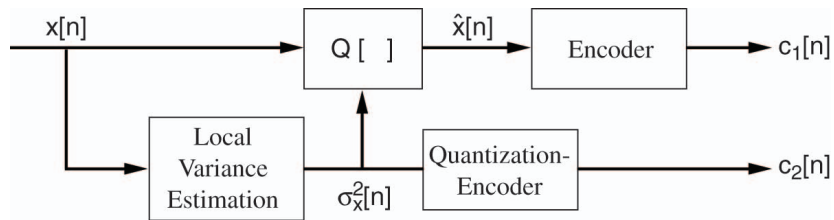
- Transformation to a uniform density:

$$g[n] = T(x[n]) = \int_{-\infty}^{x[n]} p_x(s) ds - \frac{1}{2}, \quad \frac{-1}{2} \leq g[n] \leq \frac{1}{2}$$
$$= 0 \quad \text{elsewhere}$$

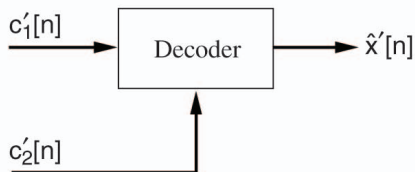
- μ -law:

$$T(x[n]) = x_{max} \frac{\log(1 + \mu \frac{|x[n]|}{x_{max}})}{\log(1 + \mu)} \text{sign}(x[n])$$

ADAPTIVE QUANTIZATION

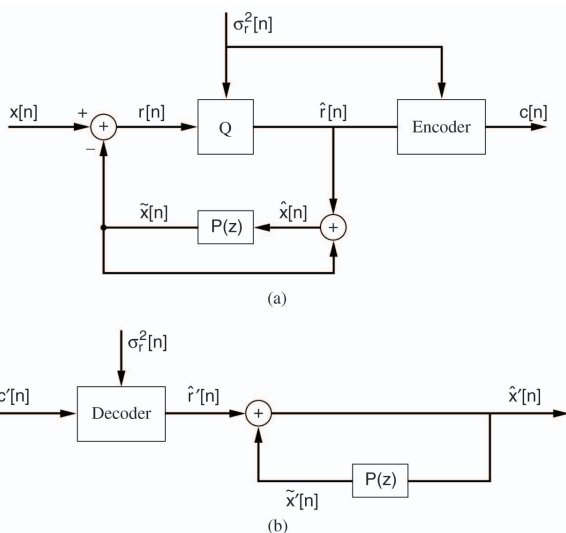


(a)



(b)

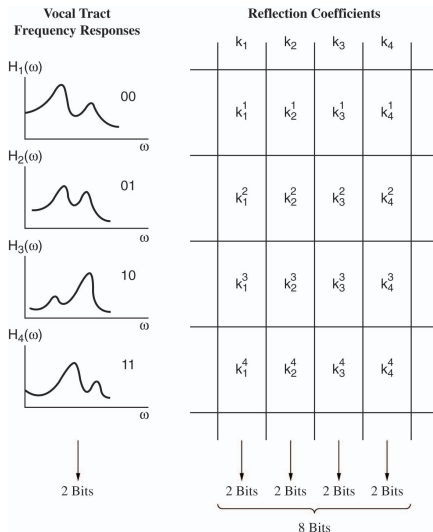
DIFFERENTIAL AND RESIDUAL QUANTIZATION



OUTLINE

- 1 INTRODUCTION
- 2 STATISTICAL MODELS
- 3 SCALAR QUANTIZATION
 - Max Quantizer
 - Companding
 - Adaptive quantization
 - Differential and Residual quantization
- 4 VECTOR QUANTIZATION
 - The k-means algorithm
 - The LBG algorithm
- 5 MODEL-BASED CODING
 - Basic Linear Prediction, LPC
 - Mixed Excitation LPC (MELP)
- 6 ACKNOWLEDGMENTS

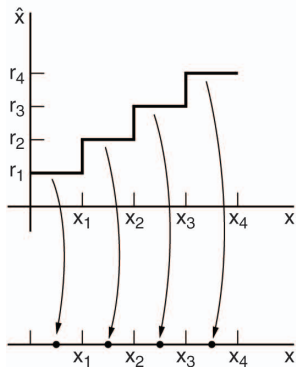
MOTIVATION FOR VQ



COMPARING SCALAR AND VECTOR QUANTIZATION

Max quantizer (1-D)

$$\hat{x} = Q[x]$$

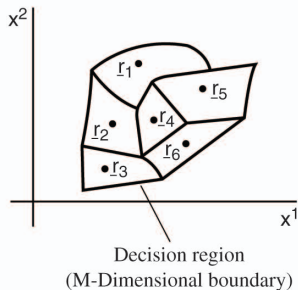


• = Centroid over the decision interval

$$D = E [(\hat{x} - x)^2]$$

Vector quantizer (2-D)

$$\hat{\underline{x}} = VQ[\underline{x}]$$



• = Centroid over the decision region

$$D = E [(\hat{\underline{x}} - \underline{x})^2(\hat{\underline{x}} - \underline{x})]$$

DISTORTION IN VQ

Here we have a multidimensional pdf $p_{\mathbf{x}}(\mathbf{x})$:

$$\begin{aligned} D &= E[(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x})] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \\ &= \sum_{i=1}^M \int \int_{\mathbf{x} \in \mathcal{C}_i} \cdots \int (\mathbf{r}_i - \mathbf{x})^T (\mathbf{r}_i - \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

Two constraints:

- A vector \mathbf{x} must be quantized to a reconstruction level \mathbf{r}_i that gives the smallest distortion:

$$\mathcal{C}_i = \{\mathbf{x} : \|\mathbf{x} - \mathbf{r}_i\|^2 \leq \|\mathbf{x} - \mathbf{r}_l\|^2, \forall l = 1, 2, \dots, M\}$$

- Each reconstruction level \mathbf{r}_i must be the centroid of the corresponding decision region, i.e., of the cell \mathcal{C}_i :

$$\mathbf{r}_i = \frac{\sum_{\mathbf{x}_m \in \mathcal{C}_i} \mathbf{x}_m}{\sum_{\mathbf{x}_m \in \mathcal{C}_i} 1} \quad i = 1, 2, \dots, M$$

DISTORTION IN VQ

Here we have a multidimensional pdf $p_{\mathbf{x}}(\mathbf{x})$:

$$\begin{aligned} D &= E[(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x})] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \\ &= \sum_{i=1}^M \int \int_{\mathbf{x} \in \mathcal{C}_i} \cdots \int (\mathbf{r}_i - \mathbf{x})^T (\mathbf{r}_i - \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

Two constraints:

- A vector \mathbf{x} must be quantized to a reconstruction level \mathbf{r}_i that gives the smallest distortion:

$$\mathcal{C}_i = \{\mathbf{x} : \|\mathbf{x} - \mathbf{r}_i\|^2 \leq \|\mathbf{x} - \mathbf{r}_l\|^2, \forall l = 1, 2, \dots, M\}$$

- Each reconstruction level \mathbf{r}_i must be the centroid of the corresponding decision region, i.e., of the cell \mathcal{C}_i :

$$\mathbf{r}_i = \frac{\sum_{\mathbf{x}_m \in \mathcal{C}_i} \mathbf{x}_m}{\sum_{\mathbf{x}_m \in \mathcal{C}_i} 1} \quad i = 1, 2, \dots, M$$

THE K-MEANS ALGORITHM

- S1:

$$D = \frac{1}{N} \sum_{k=0}^{N-1} (\hat{\mathbf{x}}_k - \mathbf{x}_k)^T (\hat{\mathbf{x}}_k - \mathbf{x}_k)$$

- S2: Pick an initial guess at the reconstruction levels $\{\mathbf{r}_i\}$
- S3: For each \mathbf{x}_k elect an \mathbf{r}_i closest to \mathbf{x}_k . Form clusters (*clustering step*)
- S4: Find the mean of \mathbf{x}_k in each cluster which gives a new \mathbf{r}_i . Compute D .
- S5: Stop when the change in D over two consecutive iterations is insignificant.

THE LBG ALGORITHM

- Set the *desired* number of cells: $M = 2^B$
- Set an initial codebook $\mathcal{C}^{(0)}$ with *one* codevector which is set as the average of the entire training sequence, \mathbf{x}_k , $k = 1, 2, \dots, N$.
- Split the codevector into two and get an *initial* new codebook $\mathcal{C}^{(1)}$.
- Perform a k-means algorithm to optimize the codebook and get the *final* $\mathcal{C}^{(1)}$
- Split the final codevectors into four and repeat the above process until the desired number of cells is reached.

THE LBG ALGORITHM

- Set the *desired* number of cells: $M = 2^B$
- Set an initial codebook $\mathcal{C}^{(0)}$ with *one* codevector which is set as the average of the entire training sequence, \mathbf{x}_k , $k = 1, 2, \dots, N$.
- Split the codevector into two and get an *initial* new codebook $\mathcal{C}^{(1)}$.
- Perform a k-means algorithm to optimize the codebook and get the *final* $\mathcal{C}^{(1)}$
- Split the final codevectors into four and repeat the above process until the desired number of cells is reached.

THE LBG ALGORITHM

- Set the *desired* number of cells: $M = 2^B$
- Set an initial codebook $\mathcal{C}^{(0)}$ with *one* codevector which is set as the average of the entire training sequence, \mathbf{x}_k , $k = 1, 2, \dots, N$.
- Split the codevector into two and get an *initial* new codebook $\mathcal{C}^{(1)}$.
- Perform a k-means algorithm to optimize the codebook and get the *final* $\mathcal{C}^{(1)}$
- Split the final codevectors into four and repeat the above process until the desired number of cells is reached.

THE LBG ALGORITHM

- Set the *desired* number of cells: $M = 2^B$
- Set an initial codebook $\mathcal{C}^{(0)}$ with *one* codevector which is set as the average of the entire training sequence, $\mathbf{x}_k, k = 1, 2, \dots, N$.
- Split the codevector into two and get an *initial* new codebook $\mathcal{C}^{(1)}$.
- Perform a k-means algorithm to optimize the codebook and get the *final* $\mathcal{C}^{(1)}$
- Split the final codevectors into four and repeat the above process until the desired number of cells is reached.

THE LBG ALGORITHM

- Set the *desired* number of cells: $M = 2^B$
- Set an initial codebook $\mathcal{C}^{(0)}$ with *one* codevector which is set as the average of the entire training sequence, \mathbf{x}_k , $k = 1, 2, \dots, N$.
- Split the codevector into two and get an *initial* new codebook $\mathcal{C}^{(1)}$.
- Perform a k-means algorithm to optimize the codebook and get the *final* $\mathcal{C}^{(1)}$
- Split the final codevectors into four and repeat the above process until the desired number of cells is reached.

OUTLINE

- 1 INTRODUCTION
- 2 STATISTICAL MODELS
- 3 SCALAR QUANTIZATION
 - Max Quantizer
 - Companding
 - Adaptive quantization
 - Differential and Residual quantization
- 4 VECTOR QUANTIZATION
 - The k-means algorithm
 - The LBG algorithm
- 5 MODEL-BASED CODING
 - Basic Linear Prediction, LPC
 - Mixed Excitation LPC (MELP)
- 6 ACKNOWLEDGMENTS

BASIC CODING SCHEME IN LPC

- Vocal tract system function:

$$H(z) = \frac{A}{1 - P(z)}$$

where

$$P(z) = \sum_{k=1}^p a_k z^{-1}$$

- Input is binary: impulse/noise excitation.
- If frame rate is 100 frames/s and we use 13 parameters ($p = 10$, 1 for Gain, 1 for pitch, 1 for voicing decision) we need 1300 parameters/s, instead of 8000 samples/s for $f_s = 8000\text{Hz}$.

BASIC CODING SCHEME IN LPC

- Vocal tract system function:

$$H(z) = \frac{A}{1 - P(z)}$$

where

$$P(z) = \sum_{k=1}^p a_k z^{-1}$$

- Input is binary: impulse/noise excitation.
- If frame rate is 100 frames/s and we use 13 parameters ($p = 10$, 1 for Gain, 1 for pitch, 1 for voicing decision) we need 1300 parameters/s, instead of 8000 samples/s for $f_s = 8000\text{Hz}$.

BASIC CODING SCHEME IN LPC

- Vocal tract system function:

$$H(z) = \frac{A}{1 - P(z)}$$

where

$$P(z) = \sum_{k=1}^p a_k z^{-1}$$

- Input is binary: impulse/noise excitation.
- If frame rate is 100 frames/s and we use 13 parameters ($p = 10$, 1 for Gain, 1 for pitch, 1 for voicing decision) we need 1300 parameters/s, instead of 8000 samples/s for $f_s = 8000\text{Hz}$.

SCALAR QUANTIZATION WITHIN LPC

For 7200 bps:

- Voiced/unvoiced decision: 1 bit
- Pitch (if voiced): 6 bits (uniform)
- Gain: 5 bits (nonuniform)
- Poles d_i : 10 bits (nonuniform) [5 bits for frequency and 5 bits for bandwidth] \times 6 poles = 60 bits

So: $(1 + 6 + 5 + 60) \times 100 \text{ frames/s} = 7200 \text{ bps}$

SCALAR QUANTIZATION WITHIN LPC

For 7200 bps:

- Voiced/unvoiced decision: 1 bit
- Pitch (if voiced): 6 bits (uniform)
- Gain: 5 bits (nonuniform)
- Poles d_i : 10 bits (nonuniform) [5 bits for frequency and 5 bits for bandwidth] \times 6 poles = 60 bits

So: $(1 + 6 + 5 + 60) \times 100 \text{ frames/s} = 7200 \text{ bps}$

REFINEMENTS TO THE BASIC LPC CODING SCHEME

- Companding in the form of a logarithmic operator on pitch and gain
- Instead of poles use the reflection (or the PARCOR) coefficients k_i , (nonuniform)
- Companding of k_i :

$$\begin{aligned}g_i &= T[k_i] \\ &= \log\left(\frac{1-k_i}{1+k_i}\right)\end{aligned}$$

- Coefficients g_i can be coded at 5-6 bits each! (which results in 4800 bps for an order 6 predictor, and 100 frames/s)
- Reduce the frame rate by a factor of two (50 frames/s) gives us a bit rate of 2400 bps

REFINEMENTS TO THE BASIC LPC CODING SCHEME

- Companding in the form of a logarithmic operator on pitch and gain
- Instead of poles use the reflection (or the PARCOR) coefficients k_i , (nonuniform)
- Companding of k_i :

$$\begin{aligned}g_i &= T[k_i] \\ &= \log\left(\frac{1-k_i}{1+k_i}\right)\end{aligned}$$

- Coefficients g_i can be coded at 5-6 bits each! (which results in 4800 bps for an order 6 predictor, and 100 frames/s)
- Reduce the frame rate by a factor of two (50 frames/s) gives us a bit rate of 2400 bps

REFINEMENTS TO THE BASIC LPC CODING SCHEME

- Companding in the form of a logarithmic operator on pitch and gain
- Instead of poles use the reflection (or the PARCOR) coefficients k_i , (nonuniform)
- Companding of k_i :

$$\begin{aligned}g_i &= T[k_i] \\ &= \log\left(\frac{1-k_i}{1+k_i}\right)\end{aligned}$$

- Coefficients g_i can be coded at 5-6 bits each! (which results in 4800 bps for an order 6 predictor, and 100 frames/s)
- Reduce the frame rate by a factor of two (50 frames/s) gives us a bit rate of 2400 bps

REFINEMENTS TO THE BASIC LPC CODING SCHEME

- Companding in the form of a logarithmic operator on pitch and gain
- Instead of poles use the reflection (or the PARCOR) coefficients k_i , (nonuniform)
- Companding of k_i :

$$\begin{aligned}g_i &= T[k_i] \\ &= \log\left(\frac{1-k_i}{1+k_i}\right)\end{aligned}$$

- Coefficients g_i can be coded at 5-6 bits each! (which results in 4800 bps for an order 6 predictor, and 100 frames/s)
- Reduce the frame rate by a factor of two (50 frames/s) gives us a bit rate of 2400 bps

REFINEMENTS TO THE BASIC LPC CODING SCHEME

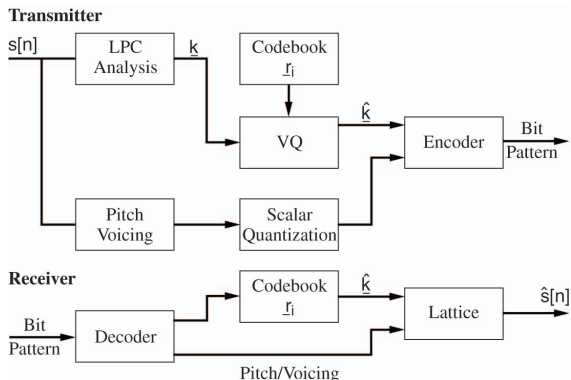
- Companding in the form of a logarithmic operator on pitch and gain
- Instead of poles use the reflection (or the PARCOR) coefficients k_i , (nonuniform)
- Companding of k_i :

$$\begin{aligned}g_i &= T[k_i] \\ &= \log\left(\frac{1-k_i}{1+k_i}\right)\end{aligned}$$

- Coefficients g_i can be coded at 5-6 bits each! (which results in 4800 bps for an order 6 predictor, and 100 frames/s)
- Reduce the frame rate by a factor of two (50 frames/s) gives us a bit rate of 2400 bps

VQ IN LPC CODING

- ▷ A 10-bit codebook (1024 codewords), 800 bps VQ provides a comparable quality to a 2400 bps scalar quantizer.
- ▷ A 22-bit codebook (4200000 codewords), 2400 bps VQ provides a higher output speech quality.



UNIQUE COMPONENTS OF MELP

- Mixed pulse and noise excitation
- Periodic or aperiodic pulses
- Adaptive spectral enhancements
- Pulse dispersion filter

LINE SPECTRAL FREQUENCIES (LSFs) IN MELP

LSFs for a p th order all-pole model are defined as follows:

- 1 Form two polynomials:

$$\begin{aligned}P(z) &= A(z) + z^{-(p+1)}A(z^{-1}) \\Q(z) &= A(z) - z^{-(p+1)}A(z^{-1})\end{aligned}$$

- 2 Find the roots of $P(z)$ and $Q(z)$, ω_i which are on the unit circle.
- 3 Exclude trivial roots at $\omega_i = 0$ and $\omega_i = \pi$.

For a 2400 bps:

- 34 bits allocated to scalar quantization of the LSFs
- 8 bits for gain
- 7 bits for pitch and overall voicing
- 5 bits for multi-band voicing
- 1 bit for the jittery state

which is 54 bits. With a frame rate of 22.5 ms, we get an 2400 bps coder.

OUTLINE

- 1 INTRODUCTION
- 2 STATISTICAL MODELS
- 3 SCALAR QUANTIZATION
 - Max Quantizer
 - Companding
 - Adaptive quantization
 - Differential and Residual quantization
- 4 VECTOR QUANTIZATION
 - The k-means algorithm
 - The LBG algorithm
- 5 MODEL-BASED CODING
 - Basic Linear Prediction, LPC
 - Mixed Excitation LPC (MELP)
- 6 ACKNOWLEDGMENTS

ACKNOWLEDGMENTS

Most, if not all, figures in this lecture are coming from the book:

T. F. Quatieri: Discrete-Time Speech Signal Processing,
principles and practice
2002, Prentice Hall

and have been used after permission from Prentice Hall

Τέλος Ενότητας



Ευρωπαϊκή Ένωση
Πρωτόκολλο Συνεταιρισμού



Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «**Ανοικτά Ακαδημαϊκά Μαθήματα στο Πανεπιστήμιο Κρήτης**» έχει χρηματοδοτήσει μόνο τη αναδιαμόρφωση του εκπαιδευτικού υλικού.
- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.



Σημειώματα

Σημείωμα αδειοδότησης

- Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά, Μη Εμπορική Χρήση, Όχι Παράγωγο Έργο 4.0 [1] ή μεταγενέστερη, Διεθνής Έκδοση. Εξαιρούνται τα αυτοτελή έργα τρίτων π.χ. φωτογραφίες, διαγράμματα κ.λ.π., τα οποία εμπεριέχονται σε αυτό και τα οποία αναφέρονται μαζί με τους όρους χρήσης τους στο «Σημείωμα Χρήσης Έργων Τρίτων».

[1] <http://creativecommons.org/licenses/by-nc-nd/4.0/>



- Ως **Μη Εμπορική** ορίζεται η χρήση:
 - που δεν περιλαμβάνει άμεσο ή έμμεσο οικονομικό όφελος από την χρήση του έργου, για το διανομέα του έργου και αδειοδόχο
 - που δεν περιλαμβάνει οικονομική συναλλαγή ως προϋπόθεση για τη χρήση ή πρόσβαση στο έργο
 - που δεν προσπορίζει στο διανομέα του έργου και αδειοδόχο έμμεσο οικονομικό όφελος (π.χ. διαφημίσεις) από την προβολή του έργου σε διαδικτυακό τόπο
- Ο δικαιούχος μπορεί να παρέχει στον αδειοδόχο ξεχωριστή άδεια να χρησιμοποιεί το έργο για εμπορική χρήση, εφόσον αυτό του ζητηθεί.

Σημείωμα Αναφοράς

Copyright Πανεπιστήμιο Κρήτης, Στυλιανού Ιωάννης. «Ψηφιακή Επεξεργασία Φωνής. Κωδικοποίηση Φωνής». Έκδοση: 1.0. Ηράκλειο/Ρέθυμνο 2015.
Διαθέσιμο από τη δικτυακή διεύθυνση: <http://www.csd.uoc.gr/~hy578>

Διατήρηση Σημειωμάτων

Οποιαδήποτε αναπαραγωγή ή διασκευή του υλικού θα πρέπει να συμπεριλαμβάνει:

- το Σημείωμα Αναφοράς
- το Σημείωμα Αδειοδότησης
- τη δήλωση Διατήρησης Σημειωμάτων
- το Σημείωμα Χρήσης Έργων Τρίτων (εφόσον υπάρχει)

μαζί με τους συνοδευόμενους υπερσυνδέσμους.

Σημείωμα Χρήσης Έργων Τρίτων

Το Έργο αυτό κάνει χρήση των ακόλουθων έργων:

Εικόνες/Σχήματα/Διαγράμματα/Φωτογραφίες

Εικόνες/σχήματα/διαγράμματα/φωτογραφίες που περιέχονται σε αυτό το αρχείο προέρχονται από το βιβλίο:

Τίτλος: *Discrete-time Speech Signal Processing: Principles and Practice*

Prentice-Hall signal processing series, ISSN 1050-2769

Συγγραφέας: Thomas F. Quatieri

Εκδότης: Prentice Hall PTR, 2002

ISBN: 013242942X, 9780132429429

Μέγεθος: 781 σελίδες

και αναπαράγονται μετά από άδεια του εκδότη.