# ECE462 – Lecture 23

### Speech Compression

Use a speech model of speech synthesis



At transmitter: Analyze speech signal to determine



At receiver: Synthesize signal based on speech model

- Channel recorders/Linear predictive coders
- Code excited linear prediction (CELP)

## Linear Predictive Coding

A model for speech synthesis





G: gain,{b<sub>i</sub>} parameters of vocal filter

### LPC-10 Standard: 2.4 Kbs

>Input signal is sampled at 8,000 samples/sec

Signal is broken into segments of 180 samples/segment

(22.5 msec of speech / segment)

Voiced vs Unvoiced decision

>Voiced Speech contains more energy (higher signal values)

>Unvoiced speech contains higher frequency values

(more zero crossings compared to voiced speech: note that both types of signal are approximately zero mean)

In LPC-10 signal is filtered (LP) with a bandwidth of 1KHZ and then the energy and zero crossings are calculated to decide between voiced - unvoiced

#### Pitch period estimation

Computationally intensive process

≻Human pitch period is usually between 2.5 to 19.5 msecs

Autocorrelation Approach: Given a voiced speech segment
 X(n), n= n<sub>0</sub>,...,n<sub>0</sub>+N

Autocorrelation function:  $R_{xx}(k) = \sum_{n_0 \neq k+k}^{n_0 + N} x(n) x(n-k)$ 

 $k \ge 0$ 

 $R_{xx}(k)$  will have a maximum when k is equal to the pitch period

#### Pitch period estimation

2. Average magnitude difference Approach (AMDF)

AMDF(P) = 
$$\frac{1}{N} \sum_{i=n_0+1}^{n_0+N} |x_i - x(i-p)|$$



#### • <u>Calculating the parameters of the Vocal Filter</u>

According to the model  $y_n = \sum_{i=1}^m b_i y_{n-\mathbb{I}^x} + G\varepsilon_n$ 

This implies that  $y_n$  is estimated (or predicted) by m past values. In practice of course does not hold exactly, so to estimate the parameters  $\{b_i\}$ , i=1,...,m and the gain G we minimize the average error

$$e_n^2 = (y_n - \sum_{i=1}^m b_i y_{n-2i} - G\varepsilon_n)^2$$

That is we minimize  $E\{e_n^2\}$  w.r.t.  $b_i$ 

Taking the derivative of  $E\{e_n^2\}$  w.r.t.  $b_i$ 

$$\frac{\partial E\{en^2\}}{\partial b_j} = \frac{\partial}{\partial b_j} E\{(y_n - \sum_{i=1}^m b_i y_{n-\frac{4}{2}i} - G\varepsilon_n)\} = 0$$

$$= -2E[y_n - \sum_{i=1}^m b_i y_{n-\frac{4}{2}i} - G\varepsilon_n / y_{n-j}] = 0$$

=> 
$$\sum_{i=1}^{m} b_i E\{y_{n-i} y_{n-j}\} = E\{y_n y_{n-j}\}$$
 <sup>(2)</sup>  
j=1,2,....,m  
because E{ $\epsilon_n y_{n-j}$ } = 0 j≠ 0

Note that assuming the speech signal is "stationary" that is its statistics do not change with time, then

$$\begin{split} & \mathsf{E}\{\mathsf{y}_{\mathsf{n}-\mathsf{i}} \; \mathsf{y}_{\mathsf{n}-\mathsf{j}}\} \; = \; \mathsf{R}_{\mathsf{y}\mathsf{y}}(|\mathsf{i}-\mathsf{j}|) \\ & \mathsf{and} \qquad \mathsf{E}\{\mathsf{y}_{\mathsf{n}} \; \mathsf{y}_{\mathsf{n}-\mathsf{j}}\} \; = \; \mathsf{R}_{\mathsf{y}\mathsf{y}}(|\mathsf{j}|) \end{split}$$

where once again the autocorrelation function is estimated as

$$\widehat{\mathsf{R}}_{yy}(\mathsf{k}) = \sum_{n_0+l+k}^{n_0+N} \mathsf{y}_n \, \mathsf{y}_{n-k}$$

Then from (2) above the following linear system of equations is obtained



Once  $\{b_i\}_{i=1,\dots,m}$  have been obtained from 1 with j=0 we obtain

$$\Rightarrow E\{ y_n^2 \} - \sum_{i=1}^M b_i E\{y_{n-i}, y_n\} = GE\{\epsilon_n y_n\} = G$$

$$\int G = R(0) - \sum_{i=1}^M b_i R(i) \int b_i definition$$

- Note: The matrix R has a nice structure. It is Toeptitz: (All diagonal elements and those parallel to the diagonal are equal -> Efficient solution of system exists (Levinson – Durbin Algorithm)
- Levinson Durbin Algorithm

1. Initialize 
$$E_0 = \widehat{R_{yy}}(0), i = 0$$

- 2. Increment i by 1
- 3. Calculate  $k_i = \sum_{j=1}^{i-1} b_j^{(i-1)} R_{yy}(i-j+1) R_{yy}(i)$ 4. Set  $b_i^{(i)} = k_i$ ,  $(i-1) \in i_{j-1}$ 5. For  $j = 1, 2, \dots, i-1$ Calculate  $b_j^{(i)} = b_j^{(i-1)} + k_i b_{i-j}$ 6.  $E_i = (1 - k_i^2) E_{i-1}$ 7. If icm so to step 2. Note: the k are known 1 as reflective

7. If i<m go to step 2. Note: the k<sub>i</sub> are known 1 as reflection coefficients

- If  $|k_i| < 1 \rightarrow$  filter stable
- For effective reconstruction of voiced segment m > 10
- Given  $k_i, i = 1, \cdots, m \leftrightarrow \{b_i\}, i = 1, \cdots, m$

Since stability is better controlled by  $k_i$ , quantize and transmit  $\{k_i\}$ .

#### Parameters to be transmitted

- 1. v/uv decision (1 bit)
- 2. Pitch period (quantized to 1 of 60 different values) 6 bits
- 3. Vocal track filter parameters

(m=10 for voice)

(m=5 for unvoiced)

Plus G (5 bit log compounded quant) For  $k_1, k_2: g_i = \frac{1+k_i}{1-k_i} \rightarrow 5$  bit uniform quantized  $k_3, k_4: 5$  bit uniform quantized  $k_5, k_8: 4$  bit uniform quantized  $k_9: 2$  bits,  $k_{10}: 3$  bits + 1 bit for synchronization -> 54 bits/frame

- So for LPC-10 coder
- 54 bits/frame of 180 samples 8,000 samples/sec
- So bit rate:(8000/180)x54 = 2400 bps.
- This encoding approach generates low bit rate speech, with "artificial voice quality."