

LPC para voz humana

$$x(n) = \sum_{k=1}^p a_k x(n-k) + G u(n)$$

AR - Model
(IIR Filter)

$$x(z) = \sum_{k=1}^p a_k x(z) z^{-k} + G u(z)$$

$$x(z) \left(1 - \sum_{k=1}^p a_k z^{-k} \right) = G u(z)$$

$$H(z) = \frac{x(z)}{u(z)} = \frac{G}{1 - \sum_{k=1}^p a_k z^{-k}}$$

All poles must
for voice
tract.

Predicción LINEAL:

$$\hat{x}(n) = \sum_{k=1}^p a_k x(n-k)$$

predicir muestra actual
de datos anteriores

$$\begin{aligned} e(n) &= x(n) - \hat{x}(n) \\ &= x(n) - \sum_{k=1}^p a_k x(n-k) = G u(n) \end{aligned}$$

obj el error de predicción es la entrada
al sistema, i.e., corresponde al sonido
que se produce en las cuerdas vocales

Una estimación óptima de $\hat{x}(n)$ es aquella que hace
que $E[e^2(n)]$ sea mínimo (i.e., MMSE)

$$E[e^2(n)] = E[(x(n) - \hat{x}(n))^2] \\ = E[x^2(n) - 2x(n)\hat{x}(n) + \hat{x}^2(n)]$$

es conveniente expresar

$$\hat{x}(n) = \sum_{k=1}^p a_k x(n-k) = \underline{a}^T \underline{x}_n \quad \text{donde}$$

$$\underline{a} = [a_1, a_2, \dots, a_p]^T$$

$$\underline{x}_n = [x(n-1), x(n-2), \dots, x(n-p)]^T$$

OBS) $\hat{x}(n) = \underline{a}^T \underline{x}_n$ es 1×1
 $= \underline{x}_n^T \underline{a}$ es equivalente

$$\Rightarrow E[e^2(n)] = E[x^2(n) - 2x(n)\underline{a}^T \underline{x}_n + \underline{a}^T \underline{x}_n \underline{x}_n^T \underline{a}] \\ = E[x^2(n)] - 2\underline{a}^T \underbrace{E[x(n)\underline{x}_n]}_{\substack{\underline{r}_x \\ (p \times 1)}} + \underline{a}^T \underbrace{E[\underline{x}_n \underline{x}_n^T]}_{\substack{\underline{R}_x \\ (p \times p)}} \underline{a}$$

$$\underline{r}_x = \begin{bmatrix} E[x(n) \cdot x(n-1)] \\ E[x(n) \cdot x(n-2)] \\ \vdots \\ E[x(n) \cdot x(n-p)] \end{bmatrix} = \begin{bmatrix} r_{xx}(1) \\ r_{xx}(2) \\ \vdots \\ r_{xx}(p) \end{bmatrix}$$

$$\underline{R_x} = \begin{bmatrix} E[x^{(n-1)} x^{(n-1)}] & E[x^{(n-1)} x^{(n-2)}] & \dots & E[x^{(n-1)} x^{(n-p)}] \\ E[x^{(n-2)} x^{(n-1)}] & E[x^{(n-2)} x^{(n-2)}] & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \\ E[x^{(n-p)} x^{(n-1)}] & \dots & \dots & E[x^{(n-p)} x^{(n-p)}] \end{bmatrix}$$

$$= \begin{bmatrix} r_{xx}(0) & r_{xx}(1) & \dots & r_{xx}(p-1) \\ r_{xx}(1) & r_{xx}(0) & & \\ \vdots & & \ddots & \\ r_{xx}(p-1) & & & r_{xx}(0) \end{bmatrix} \quad \begin{array}{l} \text{Toeplitz y} \\ \text{simétrica} \end{array}$$

OBS) $E[C \times c_n] \times [C \times m] = r_{xx}(n, m)$ auto correlación
 si x es WSS $\Rightarrow r_{xx}(n, m) = r_{xx}(|n-m|)$

$$\Rightarrow E\{e^2(n)\} = r_{xx}(0) - 2\underline{a}^T \underline{r_x} + \underline{a}^T \underline{R_x} \underline{a}$$

$$\nabla E\{e^2(n)\} = 0 \Rightarrow -2\underline{r_x} + 2\underline{R_x} \underline{a} = 0$$

(derivar con respecto a "a")

$$\Rightarrow \boxed{\underline{R_x} \underline{a} = \underline{r_x}}$$

Esto se resuelve eficientemente con Levinson - Durbin

OBS) • los coeficientes "a" sólo representan la predicción lineal, pero no el filtro AR. Este filtro se compone por:

$$\hat{a} = [1, -a_1, -a_2, \dots, -a_p]$$

• la estimación de r_{xx} no es trivial. Una forma es usar

$$\hat{r}_{xx}(m) = \frac{1}{N} \sum_{n=0}^{N-m-1} x(n) \cdot x(n+m) \quad 0 \leq m \leq p$$

lo que se obtiene con "xcorr" en matlab.

• Una estimación adecuada de r_{xx} requiere de al menos 6 ciclos de la señal $x(n)$

¿Para qué se usa LPC en voz humana?

- Telefonía celular y VoIP : CELP y variaciones
- Reconocimiento del habla : MFCC y MF-LPC
- Reconocimiento de locutor : MFCC y e(n)
- Aplicaciones climáticas : Filtrado inverso : e(n)
- Aplicaciones fonéticas : Frecuencia de LP polos
- Síntesis fisiológica de voz : mejora sobre síntesis con autoexcitación
- Reducción de ruido : Kalman LPC

OTROS CAMPOS :

- Señales de tiempo y señales : modulaciones que afectan a señales que pueden ser estimadas y renovadas con LPC

Dificultades típicas:

- Glotis varía su condición de borde temporalmente (no es un tubo cerrado)
- Resonancias nasales \rightarrow ceros en Hz)
- Acoplamiento resonancias subglóticas \rightarrow ceros en Hz) y polos \hookrightarrow perturbaciones
- Acoplamiento no lineal entre Fuente y Filtro
- Problemas de estimación de LPC en condiciones de ruido de fondo alto.
- El residual o Fuente en las cuerdas no siempre está en las cuerdas ni tiene mucho valor fisiológicamente o físicamente relevante.
- No siempre LPC está bien definido para consonantes e imputencias de la voz
- La Fuente puede ceros y polos en su estructura y es difícil diferenciar éstos de los efectos del tracto vocal.

TRANSFORMADA DISCRETA DE COSENOS

$$\begin{aligned}
 X_{\text{DCT}}^{(N)}(k) &= 2 \sum_{n=0}^{N-1} x(n) \cos\left(\frac{\pi k}{2N} (2n+1)\right) \quad k=0, 1, \dots, N-1 \\
 (\text{DCT-II}) & \\
 &= 2 \sum_{n=0}^{N-1} x(n) \cos\left(\frac{2\pi k}{2N} \left(n + \frac{1}{2}\right)\right)
 \end{aligned}$$

OBS)

$$\text{re} \left\{ X^{(2N)}(k) \right\} = \text{re} \left\{ \sum_{n=0}^{N-1} x(n) e^{-j \frac{2\pi k}{N} n} \right\}$$

PORTE REAL DE UNA DFT

$$= \sum_{n=0}^{N-1} x(n) \cos\left(\frac{2\pi k}{N} n\right)$$

RELACION ENTRE DCT Y DFT

$$\Rightarrow X_{\text{DCT}}^{(N)}(k) = e^{j \frac{2\pi k}{2N}} X^{(2N)}(k) + e^{-j \frac{2\pi k}{2N}} X^{(2N)}(k)$$

- OBS)
- DCT es siempre real
 - Tiene menores discontinuidades (asume simetría por) lo que hace que la energía se distribuya en pocos valores de k (menos que en la DFT)
 - Se utiliza mucho en compresión de señales por este motivo.

¿ De qué sirve cepstrum ?

Suponga de una señal esta filtrada de modo que

$$x(n) = u(n) * h(n)$$

$$X(\omega) = U(\omega) \cdot H(\omega)$$

$$\begin{aligned} \log(x(\omega)) &= \log(u(\omega) \cdot H(\omega)) \\ &= \log(u(\omega)) + \log(H(\omega)) \end{aligned}$$

oBSJ • Bajo ciertas condiciones las características de $u(\omega)$ y $H(\omega)$ son suficientemente distintas para ser separadas directamente por el logaritmo (ie, no tienen componentes cruzadas en sus rangos de frecuencia)

- si $h(n)$ es de bajo orden, su energía se concentra en los primeros coef. cepstrum
- si $u(n)$ oscila a mayor frecuencia, la energía de $u(n)$ se concentra en componentes superiores del cepstrum
- Esto sucede en voz y análisis de señales sísmicas.

$$\begin{aligned} \Rightarrow \hat{h}(n) &= \text{cepstrum real } (1 : L) \\ \hat{x}(n) &= \text{cepstrum real } (L+1 : \frac{N}{2}) \end{aligned}$$

oBSJ • Cepstrum es simétrico con respecto a su punto medio $(\frac{N}{2})$.

- L es el largo del "lifter", N es el largo de la señal.