

Lección 3.1:

Uso del Cepstrum para el procesamiento de señales de VOZ

Dr. Alessandro Presacco
Dr. Jesús Savage

4 de mayo de 2021

Índice

- 1 Sistema homomórfico.
- 2 Cepstrum.
- 3 Uso del Cepstrum para el procesamiento de señales de voz.
- 4 Medidas de distancia usando Cepstrum.

Sistema homomórfico

Un sistema homomórfico es una clase de sistemas non-lineales que se puede definir con el principio generalizado de sobreposición

$$\begin{aligned}y[n] &= L\{x[n]\} = L\{x_1[n] + x_2[n]\} = \\ &= L\{x_1[n]\} + L\{x_2[n]\} = \\ &= y_1[n] + y_2[n]\end{aligned}$$

y

$$= L\{ax[n]\} = aL\{x[n]\} = ay[n]$$

where $L\{*\}$ es el operador linear que representa el sistema

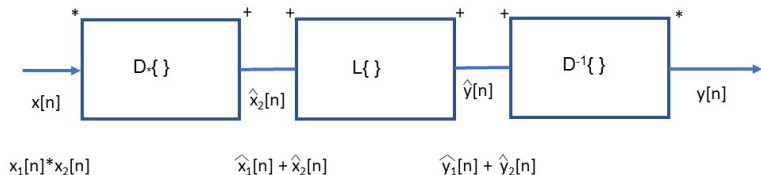
Sistema homomórfico de convolucion

Un sistema homomórfico para la convolucion es un sistema que obedece al principio de sobreposición donde la adición está remplazada por la convolucion

Para una señal $x[n] = x_1[n] * x_2[n]$, la señal de salida $y[n]$ es:

$$\begin{aligned} y[n] &= H[x[n]] = H[x_1[n] * x_2[n]] = \\ &= H[x_1[n]] * H[x_2[n]] = y_1[n] * y_2[n] \end{aligned}$$

Un sistema homomórfico puede ser representado por una secuencia de 3 sistemas homomórficos. El primero convierte la señal de ingreso obtenido por convolucion y lo transforma en una adición. El segundo es un sistema lineal que obedece al principio de superposición. El tercero es el inverso del primer sistema



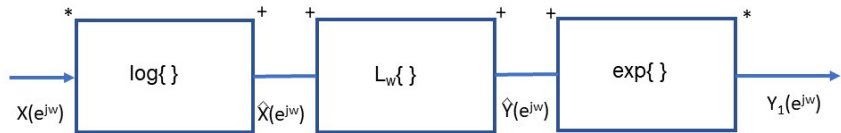
Representación con la transformada discreta de Fourier DFT

Una de las propiedades más importantes de la convolucion es el hecho que el producto de convolucion se puede representar con una multiplicación en el dominio de la frecuencia:

$$x[n] = x_1[n] * x_2[n] \rightarrow X(e^{j\omega}) = X_1(e^{j\omega}) \cdot X_2(e^{j\omega})$$

Entonces el sistema homomórfico se puede representar en el siguiente modo:

$$\begin{aligned}\hat{X}(e^{j\omega}) &= \log\{X(e^{j\omega})\} = \\ &= \log\{X_1(e^{j\omega}) \cdot X_2(e^{j\omega})\} = \log\{X_1(e^{j\omega})\} + \log\{X_2(e^{j\omega})\} = \\ &= \hat{X}_1(e^{j\omega}) + \hat{X}_2(e^{j\omega})\end{aligned}$$



$$X_1(e^{j\omega}) X_2(e^{j\omega})$$

$$\hat{X}_1(e^{j\omega}) + \hat{X}_2(e^{j\omega})$$

$$\hat{Y}_1(e^{j\omega}) + \hat{Y}_2(e^{j\omega})$$

$$Y_1(e^{j\omega}) Y_2(e^{j\omega})$$

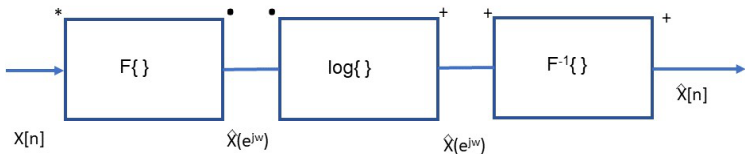
Representación con DFT para representar la señal como una secuencia

Si se quiere representar la señal como una secuencia, se puede utilizar la siguiente configuración de 3 ecuaciones:

$$X(e^{j\omega}) = \sum_{n=-\infty}^{+\infty} x[n]e^{-j\omega n}$$

$$\hat{X}(e^{j\omega}) = \log\{X(e^{j\omega})\}$$

$$\hat{x}[n] = \frac{1}{2\pi} \int_{-\pi}^{+\pi} \hat{X}(e^{j\omega}) e^{jn\omega} d\omega$$



Donde F y F^{-1} son respectivamente la transformada de Fourier y la transformada inversa de Fourier. La señal de salida $\hat{X}[n]$ es el cepstrum complejo, así llamado no porque es complejo, sino porque se basa sobre un logaritmo complejo

Ambigüedad en la representación de un algoritmo complejo

Un ángulo en el plano complejo siempre es ambiguo por un múltiplo de 2π . Este problema no se aplica al exponencial complejo (inverso del logaritmo complejo), porque se representa de la siguiente manera:

$$\begin{aligned} X(e^{j\omega}) &= e^{\log(|X(e^{j\omega})|) + j\arg\{X(e^{j\omega})\}} = \\ &= e^{\log|X(e^{j\omega})|} e^{j\arg\{X(e^{j\omega})\}} = \\ &= |X(e^{j\omega})| e^{j\arg\{X(e^{j\omega})\}} \end{aligned}$$

Dada esta ecuación, se puede concluir que para cada entero múltiplo de 2π , el valor del exponencial complejo no cambia.

Representación con la transformada Z

El sistema homomórfico descrito con la DFT se puede obviamente representar también con la transformada Z:

$$\hat{X}(z) = \log\{X(z)\} =$$

$$\hat{X}(z) = \log\{X_1(z) \cdot X_2(z)\}$$

$$= \log\{X_1(z)\} + \log\{X_2(z)\}$$

Ejemplo: Con una secuencia exponencial

$$x_1 = a^n u[n], \text{ con } |a| < 1$$

Solución

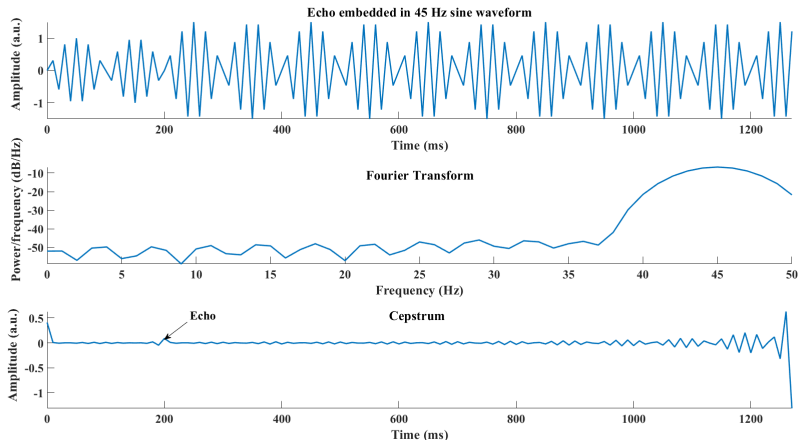
Transformada Z:

$$X_1(z) = \sum_{n=0}^{+\infty} a^n z^{-n} = \frac{1}{1 - az^{-1}}$$

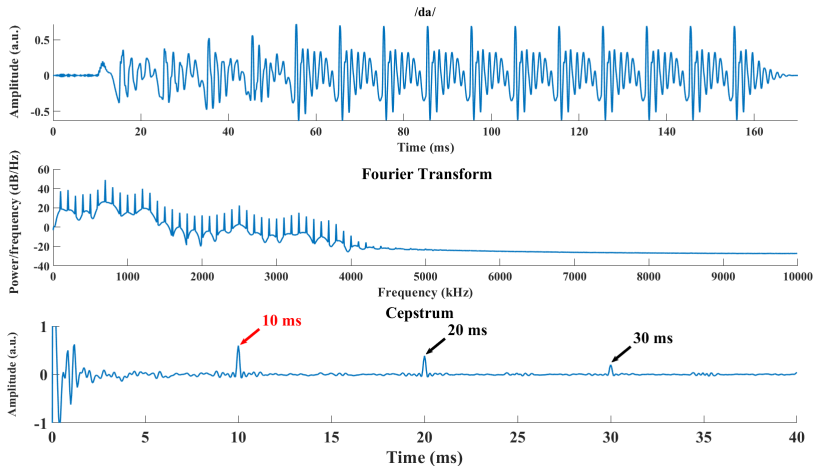
$$\hat{X}_1(z) = \log[X_1(z)] = -\log(1 - az^{-1}) = \sum_{n=1}^{+\infty} \frac{a^n}{n} z^{-n}$$

El cepstrum complejo es: $\hat{x}_1 = \frac{a^n}{n} u[n - 1]$

Ejemplo: Detectar la presencia del eco en una onda sinusoidal con frecuencia de 45 Hz



Ejemplo: Estimar la frecuencia de timbre de la siguiente señal "da"



Cepstrum y coeficientes de Mel

Mel-Frequency Cepstrum (MFC)

MFC es una representación del espectro de potencia a corto plazo de un sonido, basada en una transformada de coseno lineal de un espectro de potencia logarítmica en una escala de frecuencia mel no lineal.

Mel-frequency cepstral coefficients (MFCCs)

MFCCs son coeficientes que forman el MFC. La diferencia entre el cepstrum y el cepstrum de frecuencia mel es que en el MFC, las bandas de frecuencia están igualmente espaciadas en la escala mel.

Ecuaciones para el Mel-Spectrum

El Mel-spectrum de la m^{th} ventana de la DFT se define de la siguiente manera:

$$MF_m[r] = \frac{1}{A_r} \sum_{k=L_r}^{U_r} |V_r[k]X_m[k]|^2, r = 1, 2, \dots, R$$

con $R =$ numero de filtros (40 by default en Matlab).

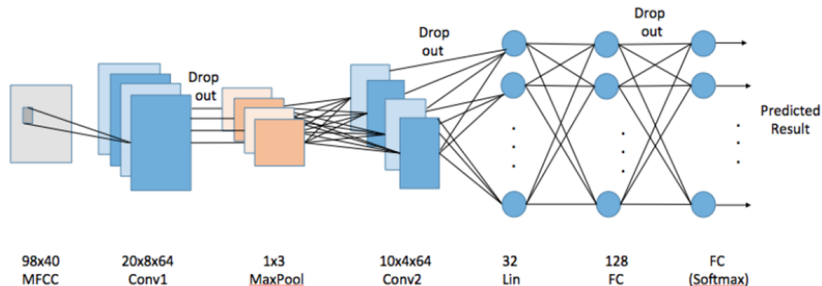
$V_r[k]$ = funcion de peso de la DFT.

$A_r = \sum_{k=L_r}^{U_r} |V_r[k]|^2 =$ factor de normalización para el r^{th} filtro de mel, que se utiliza para obtener uno espectro de mel aplanado

Para cada ventana, m , se calcula una transformada de coseno discreta del logaritmo de la magnitud de las salidas del filtro para formar la función $mfcc_m[n]$:

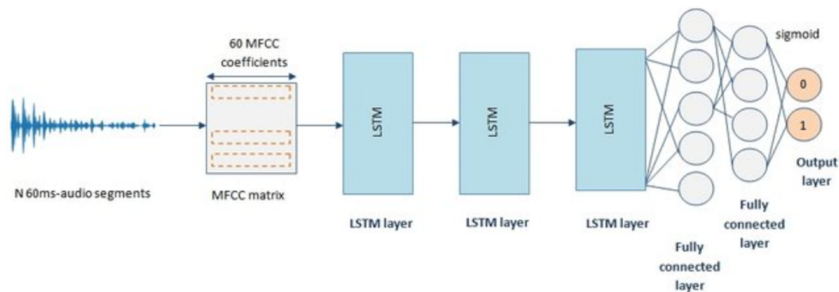
$$mfcc_m[n] = \frac{1}{R} \sum_{r=1}^R \log(MF_m[r]) \cos\left[\frac{2\pi}{R}\left(r + \frac{1}{2}\right)n\right]$$

Ejemplo: MFCC utilizado como señal de ingreso para Convolution Neural Network (CNN)



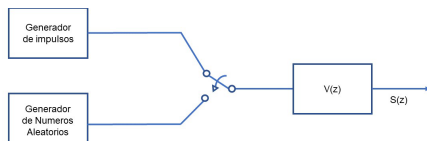
Li et al (2017) "Speech Command Recognition with Convolutional Neural Network" (Computer Science)

Ejemplo: MFCC utilizado como señal de ingreso para a Long short-term memory (LSTM) artificial recurrent neural network (RNN)



Rejaibi et al (2019) "MFCC-based Recurrent Neural Network for Automatic Clinical Depression Recognition and Assessment from Speech" (arxiv.org/abs/1909.07208v1)

Uso del Cepstrum para el procesamiento de Voz.



$$S(n) = v(n) * p(n)$$

El objetivo es separar $p(n)$ de $s(n)$ para así hacer el reconocimiento usando solamente $v(n)$ como patrón

$$S(z) = V(z) \cdot P(z)$$

$$\hat{S}(z) = \log[V(z) \cdot P(z)] = \log[V(z)] + \log[P(z)]$$

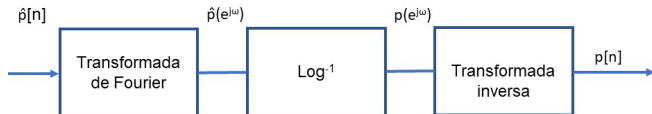
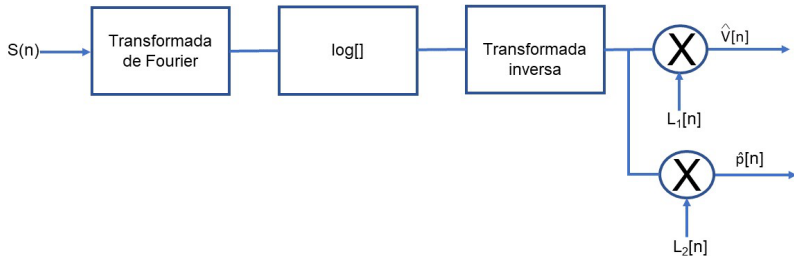
$$\hat{s}(n) = \hat{v}(n) + \hat{p}(n)$$

Si $\hat{s}(n)$ y $\hat{p}(n)$ tienen una region de soporte que no se traslapa, entonces se podrian separar $\hat{s}(n)$ y $\hat{p}(n)$. Por ejemplo, si $\hat{v}[n] = 0$ para $n \geq n_0$ y $\hat{p}(n) = 0$ para $n < n_0$, entonces

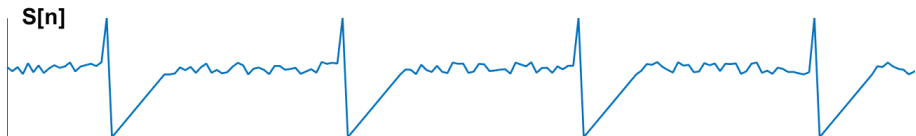
$$\hat{v}[n] = \hat{s}[n]l[n]$$

con

$$l[n] = \begin{cases} 1, & n < n_0 \\ 0, & n \geq n_0 \end{cases}$$



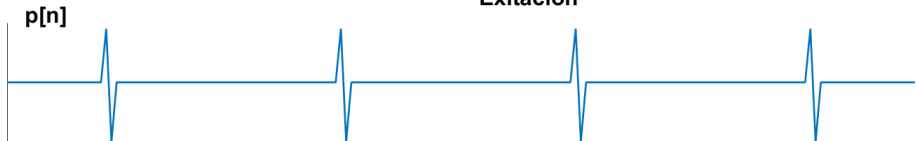
Sonido de salida del tracto vocal



Tracto vocal



Excitacion



Para $V(z) = \frac{G}{A(z)}$ (usando la expansion de Laurent) obtenemos el siguiente resultado

$$\log\left[\frac{G}{A(z)}\right] = \log G + \sum_{n=1}^{\infty} a_n z^{-n}$$

Diferenciando ambos lados con respecto a z^{-1} e igualando los coeficientes para z^{-1} , se puede derivar la siguiente recursion:

$$C_n = -a_n - \frac{1}{n} \sum_{k=1}^{n-1} k c_k a_{n-k}$$

para $a_0 = 1$ y $a_k = 0$ para $k > p$

Usando la serie de Taylor para:

$$\log\left[\frac{G^2}{|A(e^{j\omega})|^2}\right] = \sum_{n=-\infty}^{+\infty} c_n e^{-jn\omega}$$

$$c_0 = \log G^2 \text{ y } C_{-n} = C_n$$

Medidas de distancia usando Cesptrum

$$\begin{aligned}d_2^2(s, s') &= \int_{-\pi}^{+\pi} |\log S(\omega) - \log S'(\omega)|^2 \frac{d\omega}{2\pi} = \\&= \int_{-\pi}^{+\pi} \left| \sum_{n=-\infty}^{+\infty} C_n e^{-jn\omega} - \sum_{n=-\infty}^{+\infty} C'_n e^{-jn\omega} \right|^2 \frac{d\omega}{2\pi} = \\&= \int_{-\pi}^{+\pi} \left| \sum_{n=-\infty}^{+\infty} (C_n - C'_n) e^{-jn\omega} \right|^2 \frac{d\omega}{2\pi} = \\&= \int_{-\pi}^{+\pi} \sum_{n=-\infty}^{+\infty} \sum_{m=-\infty}^{+\infty} (C_n - C'_n)(C_m^* - C_m'^*) e^{-j(n-m)\omega} \frac{d\omega}{2\pi} =\end{aligned}$$

$$\begin{aligned}
&= \sum_{n=-\infty}^{+\infty} \sum_{m=-\infty}^{+\infty} (C_n - C'_n)(C_m^* - C'_m{}^*) \int_{-\pi}^{+\pi} e^{-j(n-m)\omega} \frac{d\omega}{2\pi} = \\
&= \sum_{n=-\infty}^{+\infty} \sum_{m=-\infty}^{+\infty} (C_n - C'_n)(C_m^* - C'_m{}^*) \delta(n - m) = \\
&= \sum_{n=-\infty}^{+\infty} |C_n - C'_n|^2 = \sum_{n=-\infty}^{+\infty} (C_n - C'_n)^2
\end{aligned}$$

S y S' son espectros de potencia, los cuales son funciones pares, por lo tanto sus coeficientes cepstrales son reales