

## Contenido

1. Obtención de señales de voz
  - 1.1 Conversión de archivos de audio tipo .WAV a archivos de datos tipo .PRN
  - 1.2 Conversión de archivos de datos tipo .PRN a archivos de audio tipo .WAV
2. Aspectos básicos de procesamiento digital de señales
  - 2.1 Relación entre el espectro de una señal analógica y su versión muestreada
  - 2.2 Transformada rápida de Fourier, fft
  - 2.3 Diseño de filtros FIR mediante ventanas en el tiempo
  - 2.4 Convolución mediante FFT
3. Modelo de producción de voz
4. Análisis de voz en tiempo corto
  - 4.1 Energía promedio y tasa promedio de cruces por cero. Clasificación sonoro/sordo
  - 4.2 Autocorrelación en tiempo corto
  - 4.3 Estimación de la función de tono mediante la autocorrelación en tiempo corto
  - 4.4 Transformada de Fourier en tiempo corto
5. Análisis predictivo lineal, LPC
6. Análisis homomórfico mediante cepstro real
7. Codificación de señales de voz
  - 7.1 Codificación PCM ley  $\mu$
  - 7.2 Codificación PCM diferencial
  - 7.3 Codificación subbanda mediante banco de filtros QMF FIR de dos bandas
  - 7.4 Codificación CELP
8. Síntesis de voz por concatenación de fonemas
9. Reconocimiento de voz - Alineación dinámica de tiempo

1.1 Conversión de archivos de audio tipo .WAU a archivos de datos tipo .PRN

LISTADO DEL PROGRAMA WAU\_PRN.PAS ESCRITO EN TURBO PASCAL 5.0

{ Para el procesamiento digital de señales de voz lo primero que se necesita es una señal de voz. Con la proliferación de sistemas multimedia para computadores personales, ahora es fácil capturar señales de voz y almacenarlas en archivos tipo .WAU. Como los documentos de MathCAD que se presentan en este reporte requieren señales de voz en el formato .PRN, aquí se presenta un programa en pascal que convierte archivos .WAU MONOFONICOS DE 8 BITS a archivos .PRN. }

```

uses dos,           { Biblioteca de funciones para manejo de archivos }
    crt;           { Biblioteca de funciones para manejo de pantalla }

var
  NombreArchivo : String; { Nombre del archivo. Cambia la extensión WAU-PRN }
  ArchivoIn     : text;   { Archivo de entrada tipo .WAU }
  ArchivoOut    : text;   { Archivo de salida tipo .PRN }
  i,j,n         : integer; { índice auxiliar, valor de c/muestra, # muestras }
  directorio    : SearchRec; { Estructura de pascal para leer directorios }
  c             : char;   { Lectura de cada byte del archivo de entrada }

begin
  clrscr;
  repeat
    write(#7, 'Nombre del Archivo? '); { Acepta un nombre de archivo }
    readln(NombreArchivo);
    NombreArchivo := NombreArchivo + '.wav'; { Añade la extensión .WAU }
    FindFirst(NombreArchivo, archive, directorio); { y verifica su existencia. }
  until DosError = 0; { Repite si no es válido }
  assign(ArchivoIn, NombreArchivo);
  reset(ArchivoIn); { Abre el archivo de entrada para lectura }
  i:=length(NombreArchivo); { Cambia la extensión para conformar el }
  NombreArchivo[i-8]:='n'; { nombre del archivo de salida }
  NombreArchivo[i-1]:='r';
  NombreArchivo[i-2]:='p';
  assign(ArchivoOut, NombreArchivo);
  rewrite(ArchivoOut); { Abre el archivo de salida para escritura }
  writeln('Conversión en progreso...'); { Reporta la actividad de conversión }
  for i:=1 to 48 do read(ArchivoIn, c); { Ignora los primeros 48 bytes del }
  read(ArchivoIn, c); { archivo de entrada. Los siguientes 4 }
  n:=ord(c); { bytes indican el número de muestras }
  read(ArchivoIn, c);
  n:=n + abs(256*ord(c));
  read(ArchivoIn, c);
  read(ArchivoIn, c);
  for i:=1 to n do { Lee cada una de las muestras }
  begin
    read(ArchivoIn, c);
    j:=ord(c)-128; { Elimina el nivel dc }
    writeln(ArchivoOut, j:4);
  end;
  close(ArchivoOut); { Cierra los dos archivos }
  close(ArchivoIn);
end.

```

1.2 Conversión de archivos de datos tipo .PRN a archivos de audio tipo .MAU

LISTADO DEL PROGRAMA PRN\_MAU.PAS ESCRITO EN TURBOPASCAL 5.0

```

uses dos, crt;           ( Bibliotecas de funciones para archivos y pantalla )

var
  NombreArchivo      : String;      ( Nombre de los archivos, .PRN o .MAU )
  ArchivoIn, archivoOut : text;      ( Archivos de entrada y salida )
  ValorTexto         : string;      ( Buffer para lectura de datos )
  i,j,k,n           : integer;      ( Indices Auxiliares )
  directorio        : SearchRec;    ( Estructura de archivos en pascal )
  x,max             : double;       ( Valor de cada muestra, valor maximo )
  c,d               : char;        ( Para escritura de c/byte en .MAU )

begin
  clrscr;
  repeat
    write(#7,'¿Nombre del Archivo? '); ( Acepta un nombre de archivo, le )
    readln(NombreArchivo);           ( añade la extensión .PRN y verifica )
    NombreArchivo:=NombreArchivo+'.prn'; ( su existencia. )
    FindFirst(NombreArchivo,archivo,directorio);
  until DosError=0;                 ( Repite si el archivo no existe )
  assign(ArchivoIn,NombreArchivo);
  reset(ArchivoIn);                 ( Abre el archivo de entrada para lectura )
  n:=0; max:=-1e30;                 ( Numero de muestras, muestra máxima )
  while not(eof(ArchivoIn)) do
  begin
    readln(ArchivoIn,ValorTexto);   ( Lee cada muestra, la convierte a double )
    val(ValorTexto,x,i);            ( y actualize los valores de n y de max )
    if abs(x)>max then max:=abs(x);
    inc(n);
  end;
  close(ArchivoIn);
  i:=length(NombreArchivo);        ( Confirma el nombre del archivo de salida )
  NombreArchivo[i-0]:='v';         ( añadiendo la extensión .MAU )
  NombreArchivo[i-1]:='a';
  NombreArchivo[i-2]:='u';
  assign(ArchivoOut,NombreArchivo); ( Abre el archivo de salida para escritura )
  rewrite(ArchivoOut);              ( y el de entrada para lectura )
  reset(ArchivoIn);
  writeIn('Se detectaron ',n,' muestras en ',NombreArchivo,'. ');
  write('¿Frecuencia de muestreo? '); ( Introduce la tasa de muestreo )
  readln(ValorTexto);              ( y calcula su representación )
  val(ValorTexto,j,k);             ( en bytes para escribirla en )
  i:=j mod 256;                    ( el archivo de salida )
  c:=chr(i);                       ( byte menos significativo )
  i:=j div 256;
  d:=chr(i);                       ( byte más significativo )
  writeIn('Conversión en progreso...');
  write(ArchivoOut,'RIFF');         ( Encabezado constante )
  write(ArchivoOut,chr((n+36) MOD 256)); ( Longitud del archivo )
  write(ArchivoOut,chr(((n+36) DIV 256) MOD 256));
  write(ArchivoOut,#0,#0);
  write(ArchivoOut,'WAUefat ');     ( Formato WAUE )
  write(ArchivoOut,#16,#0,#0,#0,#1,#0,#1,#0); ( Datos constantes... (¿?) )
  write(ArchivoOut,c,d,#0,#0,c,d,#0,#0); ( Frecuencia de muestreo )
  write(ArchivoOut,#1,#0,#8,#0,'data');
  write(ArchivoOut,chr(n MOD 256)); ( Número de muestras )
  write(ArchivoOut,chr((n DIV 256) MOD 256));
  write(ArchivoOut,#0,#0);
  for i:=1 to n do
  begin
    readln(ArchivoIn,ValorTexto);
    val(ValorTexto,x,k);
    j:=trunc(127.0*x/max + 128);
    c:=chr(j);
    write(ArchivoOut,c);
  end;
  close(ArchivoIn); close(ArchivoOut);
end.

```

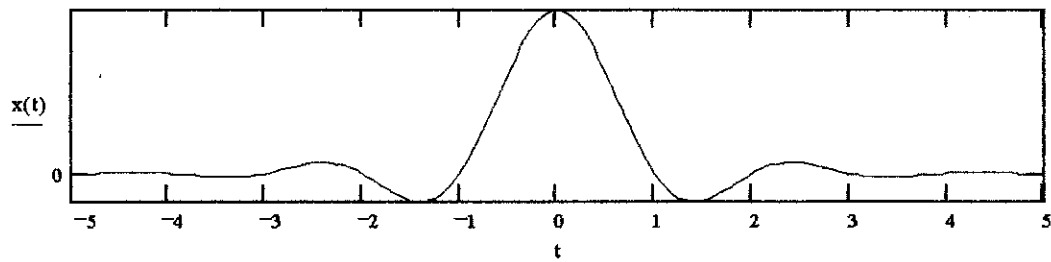
## 2.1 Espectro de $x(t)$ y de su versión muestreada $x[n] = x(nT)$

$$x(t) := \text{if}(t=0, 1, \frac{\sin(\pi \cdot t)}{\pi \cdot t}) \cdot \frac{\cos\left(\frac{\pi \cdot t}{5}\right) + 1}{2}$$

Señal analoga a discretizar

$$t := -5, -4.92 \dots 5$$

Rango de definición de la señal analoga



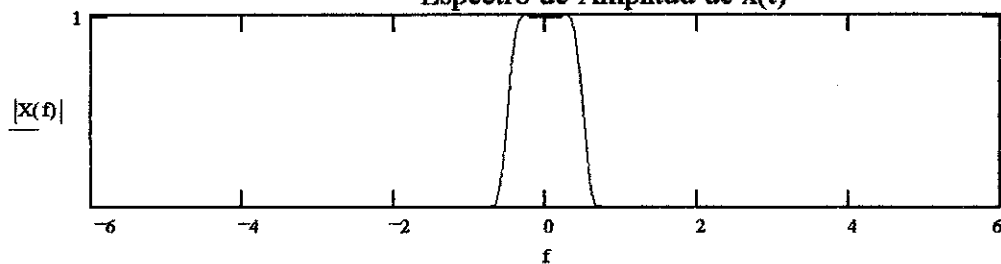
$$X(f) := \int_{-5}^5 x(t) \cdot \exp(-2i \cdot \pi \cdot f \cdot t) dt$$

Composición espectral de  $x(t)$ , calculada mediante la Transformada de Fourier

$$f := -6, -5.95 \dots 6$$

Rango de frecuencias a estudiar

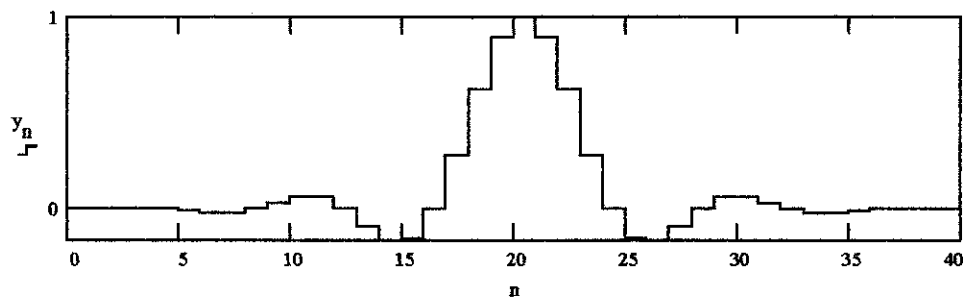
Espectro de Amplitud de  $x(t)$



$$n := 0 \dots 40$$

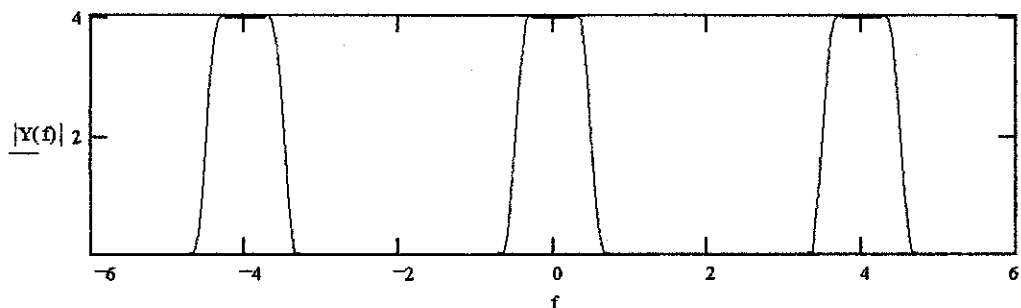
$$y_n := x\left(\frac{n}{4} - 5\right)$$

Se tomarán 41 muestras separadas 0.25 unidades de tiempo entre si para una tasa de muestreo de 4 Hz



$$Y(f) := \sum_n y_n \cdot \exp\left(\frac{-2i \cdot \pi \cdot n \cdot f}{4}\right)$$

Composición espectral de  $y[n]$ , calculada mediante la Transformada de Fourier en tiempo discreto



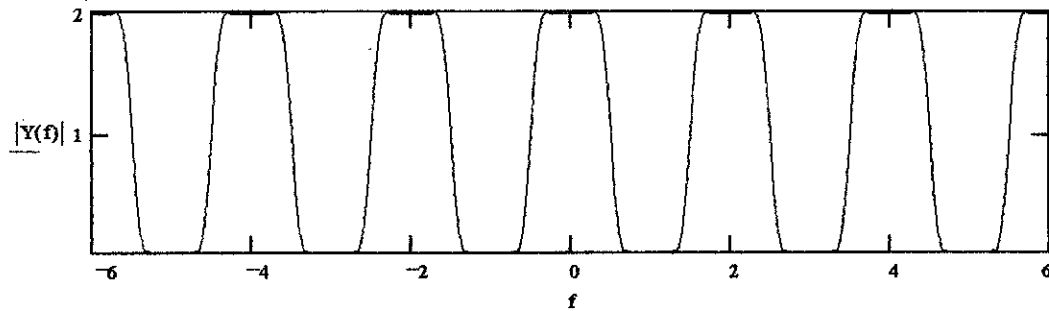
El espectro de la señal analoga se repite cada cuatro Hz en el espectro de la señal muestreada

$n := 0..20$

$$y_n := x\left(\frac{n}{2} - 5\right)$$

$$Y(f) := \sum_n y_n \cdot \exp\left(\frac{-2i \cdot \pi \cdot n \cdot f}{2}\right)$$

Al bajar la tasa de muestreo a 2 Hz las réplicas del espectro de la señal análoga se acercan más, pero aún se pueden separar.

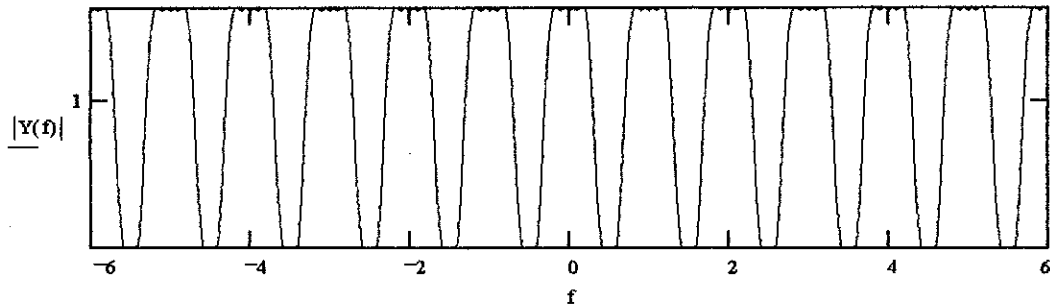


$n := 0..16$

$$y_n := x\left(\frac{n}{1.6} - 5\right)$$

$$Y(f) := \sum_n y_n \cdot \exp(-2i \cdot \pi \cdot n \cdot f)$$

La tasa de muestreo crítica es de 1.6 Hz pues el ancho de banda de la señal análoga original es de cerca de 0.8 Hz. Aunque aún se pueden distinguir las distintas réplicas ya es muy difícil tratar de separarlas.

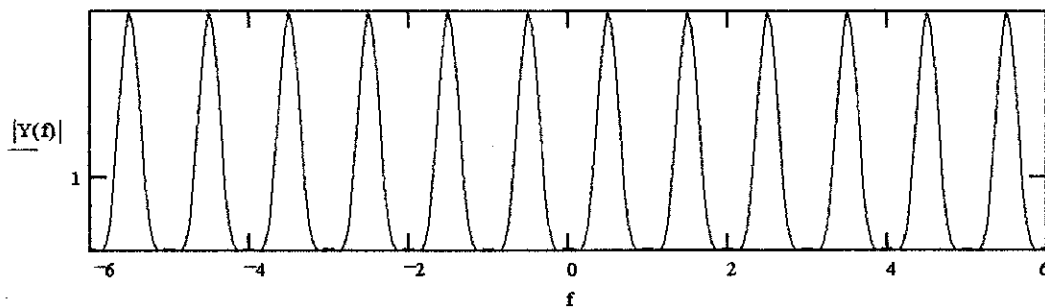


$n := 0..8$

$$y_n := x\left(\frac{n}{0.8} - 5\right)$$

$$Y(f) := \sum_n y_n \cdot \exp(-2i \cdot \pi \cdot n \cdot f)$$

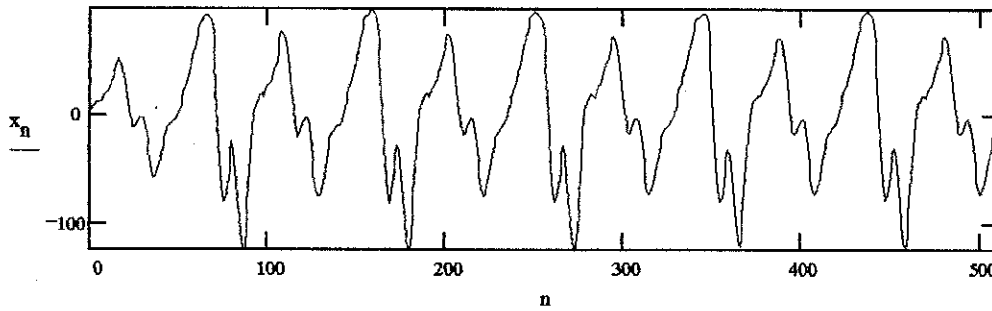
A tasas de muestreo menores que la frecuencia Nyquist el fenómeno del aliasing impide una adecuada representación de la señal original.



## 2.2 Transformada rápida de Fourier, FFT

$x := \text{READPRN}(a)$   
 $N := 512$   
 $n := 0..N-1$

Señal a analizar  
 Número de muestras  
 Índice de tiempo

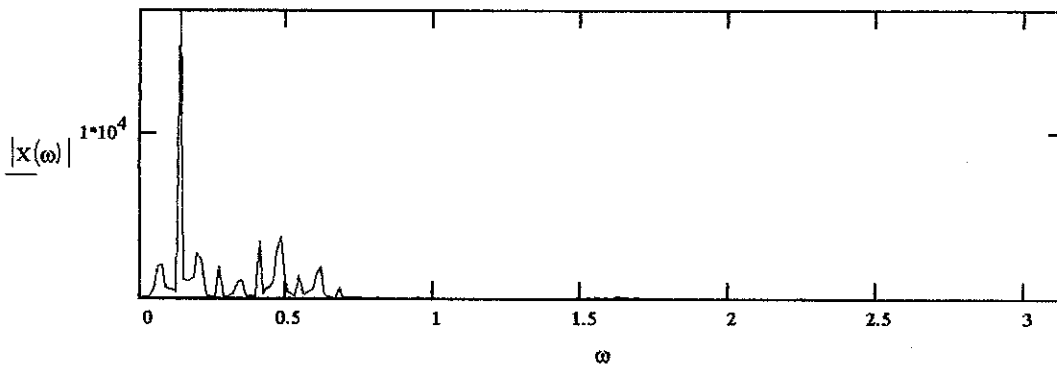


$$X(\omega) := \sum_n x_n \cdot \exp(-i \cdot \omega \cdot n)$$

Transformada de Fourier en tiempo discreto, DTFT, definida en el dominio continuo de la frecuencia normalizada,  $\omega = 2\pi f / f_s$ , donde  $f_s$  es la frecuencia de muestreo.

$$\omega := 0, \frac{2 \cdot \pi}{N} \dots \pi$$

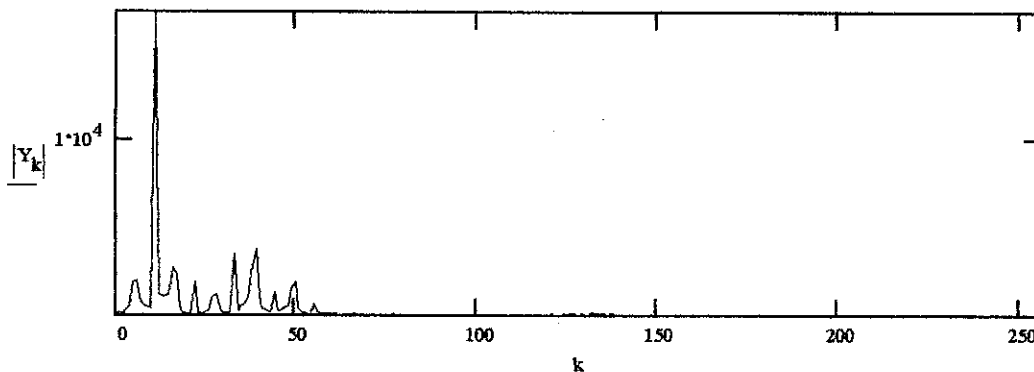
Muestreo en frecuencia de la DTFT para  $f$  de 0 a  $f_s/2$



$$Y := \text{fft}(x) \cdot \sqrt{N}$$

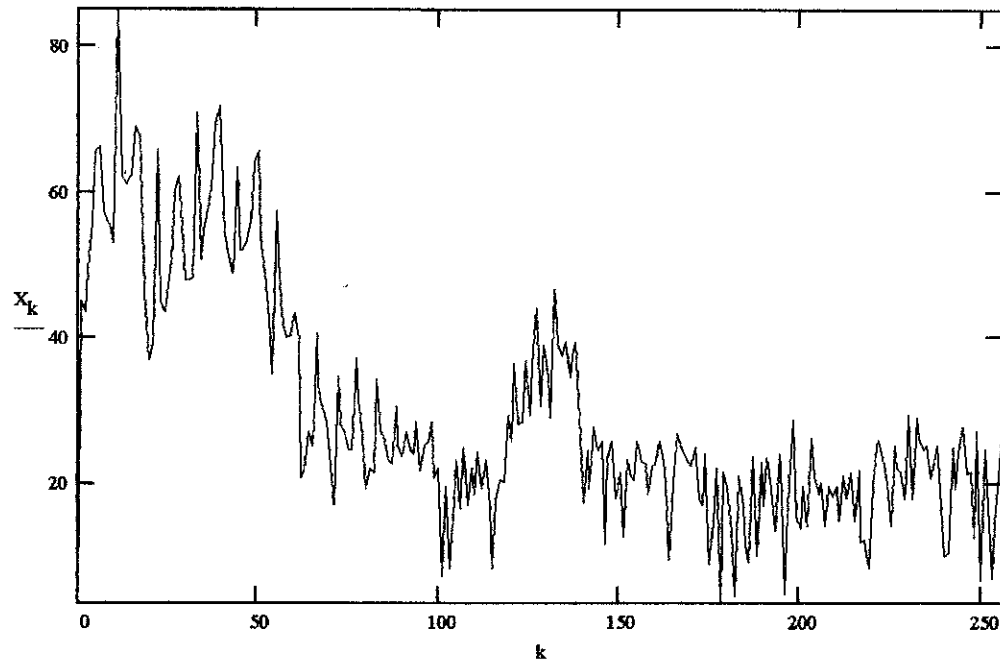
$$k := 0.. \frac{N}{2}$$

La transformada rápida de Fourier, FFT, calcula las muestras de la DTFT directamente con un algoritmo muy eficiente. La FFT particular que utiliza MathCAD requiere que  $N$  sea potencia de dos y retorna  $1+N/2$  muestras en frecuencia, correspondientes a una frecuencia entre 0 y  $f_s/2$ , donde  $f_s$  es la frecuencia de muestreo. Por alguna razón muy particular de MathCAD, la magnitud de la FFT viene normalizada respecto a la raíz del número de muestras en tiempo.



$$X_k := 20 \cdot \log(|Y_k|)$$

Una forma más común de representar el espectro de amplitud es en unidades relativas o decibeles:



Nótese que se trata de una serie de picos separados cada 6.5 muestras en frecuencia, aproximadamente, y multiplicados por una envolvente con máximos a las 13, 40 y 132 muestras. Conociendo que  $f_s=11025$  Hz, podemos determinar las frecuencias correspondientes:

$$\frac{6.5 \cdot 11025}{N} = 139.966$$

Los armónicos se suceden cada 140 Hz, que corresponde al tono de la señal de voz

$$\frac{13 \cdot 11025}{N} = 279.932$$

La primera frecuencia de resonancia de la cavidad oral es de 280 Hz

$$\frac{40 \cdot 11025}{N} = 861.328$$

La segunda frecuencia de resonancia de la cavidad oral es de 861.3 Hz

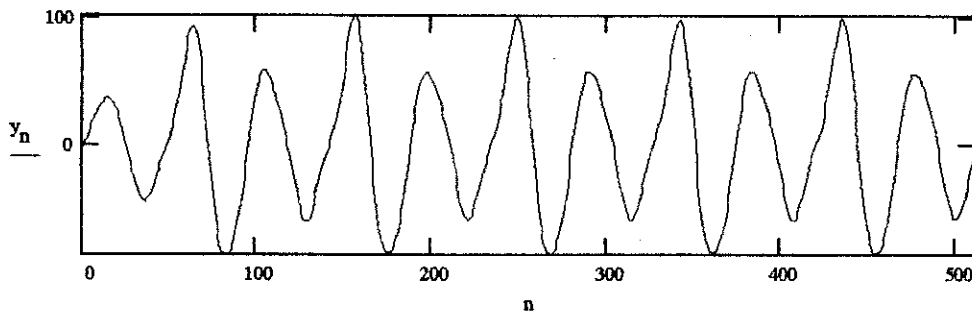
$$\frac{132 \cdot 11025}{N} = 2.842 \cdot 10^3$$

La tercera frecuencia de resonancia de la cavidad oral es de 2.84 kHz

$$k := 30.. \frac{N}{2} \quad Y_k := 0$$

$$y := \frac{\text{ifft}(Y)}{\sqrt{N}}$$

En el lado de la transformada es fácil hacer operaciones como el filtraje mediante la eliminación de las frecuencias indeseadas. En este caso, un filtro pasabajos. Nótese nuevamente la normalización inversa en la IFFT.



### 2.3 Diseño de filtros FIR mediante ventanas en el tiempo

$$\text{sinc}(\alpha) := \text{if} \left( \alpha = 0, 1, \frac{\sin(\pi \cdot \alpha)}{\pi \cdot \alpha} \right)$$

Definición de la función SENO-C

$$N := 20$$

Orden de los filtros

$$k := 0..N-1$$

Índice de tiempo

$$r_k := \frac{\text{sinc}\left(\frac{k-0.5 \cdot N}{2}\right)}{2}$$

Filtro pasabajos de media banda mediante ventana rectangular

$$h_k := \left( 0.54 - 0.46 \cdot \cos\left(\frac{2 \cdot \pi \cdot k}{N}\right) \right) \cdot r_k$$

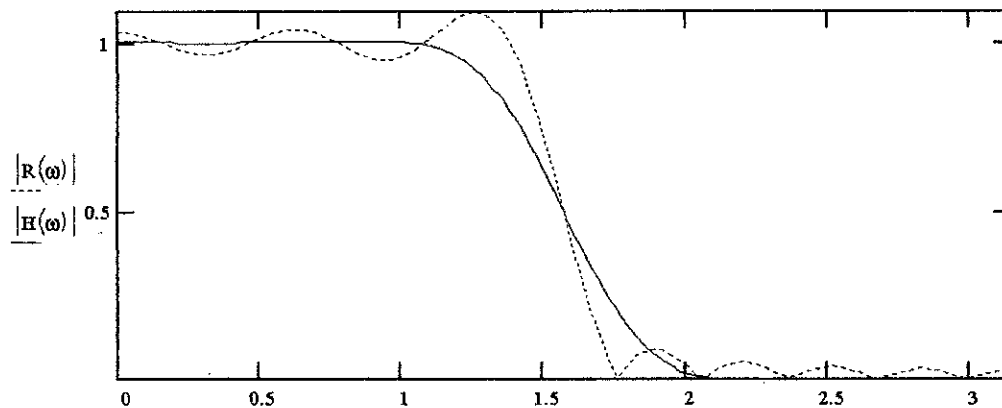
Filtro pasabajos de media banda mediante ventana Hamming

$$R(\omega) := \sum_k r_k \cdot \exp(-i \cdot k \cdot \omega)$$

Respuesta en frecuencia de cada filtro

$$H(\omega) := \sum_k h_k \cdot \exp(-i \cdot k \cdot \omega)$$

$$\omega := 0, 0.01 \cdot \pi .. \pi$$



Aunque el filtro mediante ventana rectangular permite una transición más rápida, sus bandas pasante y no-pasante presentan oscilaciones indeseables. El filtro mediante ventana Hamming tiene una transición más lenta pero, en cambio, sus bandas pasante y no-pasante son más planas.

$$N := 128$$

Orden de los filtros

$$k := 0..N-1$$

Índice de tiempo de las respuestas al impulso

Forma general de la  $k$ -ésima muestra de un filtro FIR pasabanda con frecuencias de corte inferior y superior iguales a  $\theta_1$  y  $\theta_2$ , respectivamente.

$$g(k, \theta_1, \theta_2) := \left[ \frac{\theta_2}{\pi} \cdot \text{sinc}\left[\frac{(k-0.5 \cdot N) \cdot \theta_2}{\pi}\right] - \frac{\theta_1}{\pi} \cdot \text{sinc}\left[\frac{(k-0.5 \cdot N) \cdot \theta_1}{\pi}\right] \right] \cdot \left( 0.54 - 0.46 \cdot \cos\left(\frac{2 \cdot \pi \cdot k}{N}\right) \right)$$

$$\text{lpf}_k := g\left(k, 0, \frac{\pi}{3}\right)$$

Coefficientes de un filtro pasabajos a  $fs/6$

$$\text{bpf}_k := g\left(k, \frac{\pi}{3}, \frac{2 \cdot \pi}{3}\right)$$

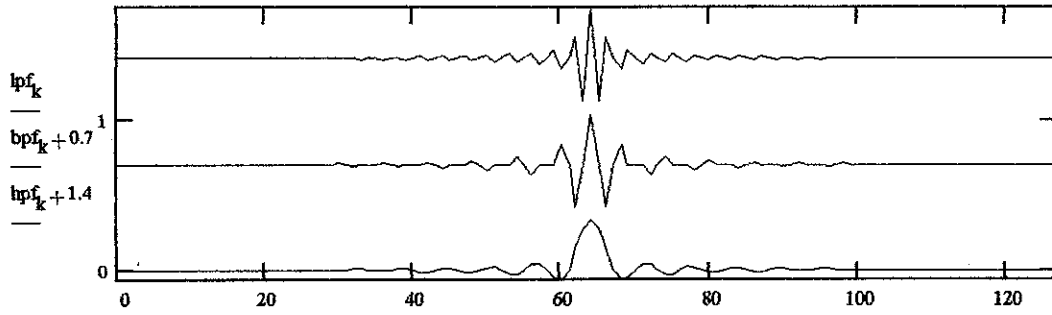
Coefficientes de un filtro pasabanda entre  $fs/6$  y  $fs/3$

$$\text{hpf}_k := g\left(k, \frac{2 \cdot \pi}{3}, \pi\right)$$

Coefficientes de un filtro pasalto desde  $fs/3$

(hasta  $fs/2$ , claro, porque la señal muestreada no debe tener componentes espectrales mayores a  $fs/2$ ).



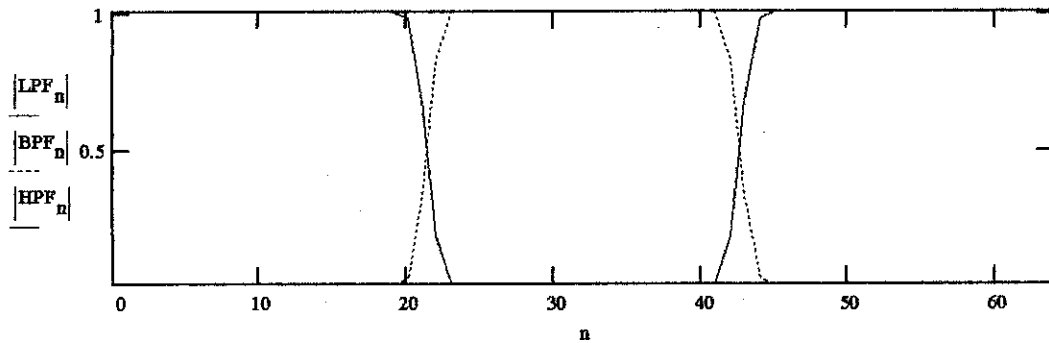


$LPF := \text{fft}(lpf) \cdot \sqrt{N}$   
 $BPF := \text{fft}(bpf) \cdot \sqrt{N}$   
 $HPF := \text{fft}(hpf) \cdot \sqrt{N}$   
 $n := 0..0.5 \cdot N$

Respuesta en frecuencia del filtro pasabajos

Respuesta en frecuencia del filtro pasabanda

Respuesta en frecuencia del filtro pasaaltos



$x := \text{READPRN}(sa)$

Señal de voz a filtrar

$n := 0.. \text{size}(x) + N$

Número de muestras en las señales filtradas

$xl_n := \sum_k lpf_k \cdot \text{if}(n \geq k, \text{if}(n - k \leq \text{size}(x), x_{n-k}, 0), 0)$

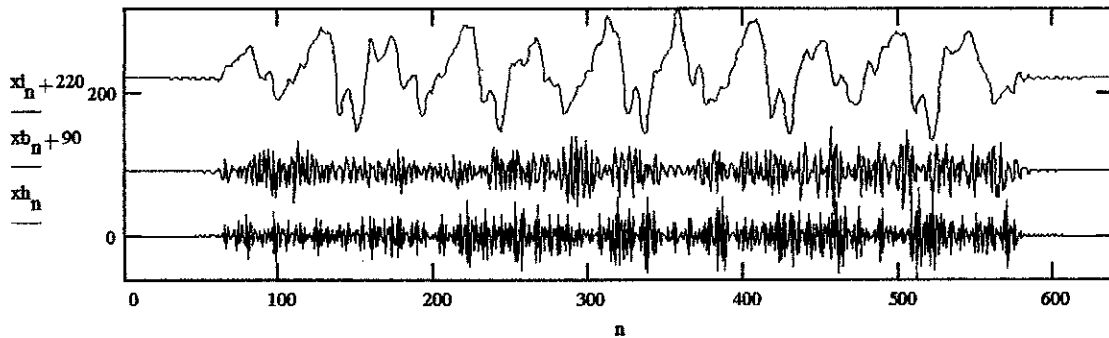
Salida del filtro pasabajos

$xb_n := \sum_k bpf_k \cdot \text{if}(n \geq k, \text{if}(n - k \leq \text{size}(x), x_{n-k}, 0), 0)$

Salida del filtro pasabanda

$xh_n := \sum_k hpf_k \cdot \text{if}(n \geq k, \text{if}(n - k \leq \text{size}(x), x_{n-k}, 0), 0)$

Salida del filtro pasaaltos

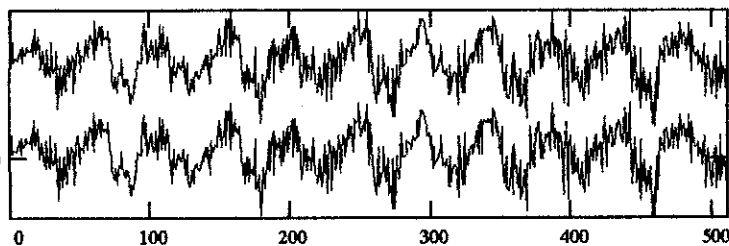


$n := 0.. \text{size}(x)$

La señal original es aproximadamente igual a la suma de las tres bandas

$x_n + 200$

$xl_n + 0.5 \cdot N + xb_n + 0.5 \cdot N + xh_n + 0.5 \cdot N$



## 2.4 Convolución mediante FFT

$N := 512$

Orden del filtro

$k := 0..N - 1$

Índice de tiempo para la respuesta al impulso

$$h_k := \text{if} \left[ k = 0.5 \cdot N, 0.2, \frac{\sin \left[ \frac{(k - 0.5 \cdot N) \cdot \pi}{5} \right]}{(k - 0.5 \cdot N) \cdot \pi} \right] \cdot \left( 0.54 - 0.46 \cdot \cos \left( \frac{2 \cdot \pi \cdot k}{N} \right) \right) \quad \text{Respuesta al impulso del Filtro pasabajos}$$

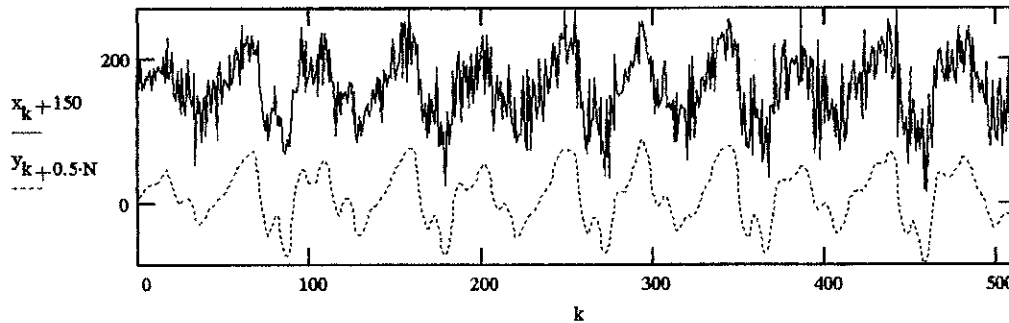
$x := \text{READPRN}(sa)$

Señal a filtrar

$n := 0..2 \cdot N - 1$

Índice de tiempo de la señal filtrada

$$y_n := \sum_k h_k \cdot \text{if}(n \geq k, \text{if}(n - k < N, x_{n-k}, 0), 0) \quad \text{Convolución lineal}$$



¿Se debe esperar obtener el mismo resultado con la IFFT del producto de las FFT?

$$X := \text{fft}(x) \cdot \sqrt{N}$$

Espectro de la señal a filtrar

$$H := \text{fft}(h) \cdot \sqrt{N}$$

Espectro de la respuesta al impulso del filtro

$$n := 0..0.5 \cdot N$$

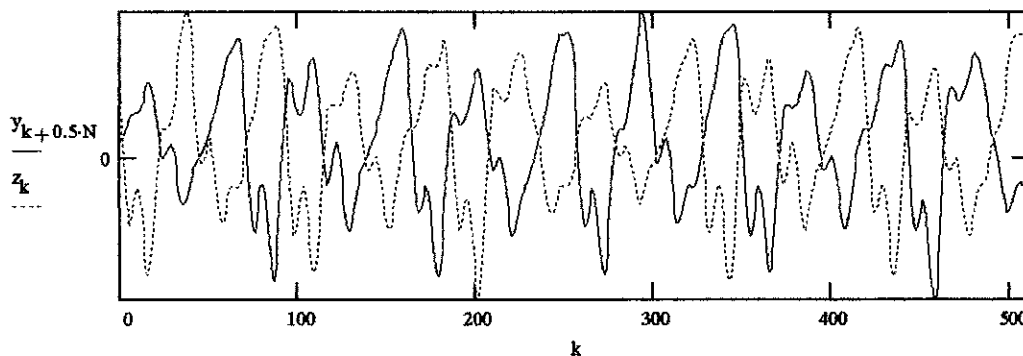
Índice de frecuencia

$$Z_n := X_n \cdot H_n$$

Producto de las transformadas

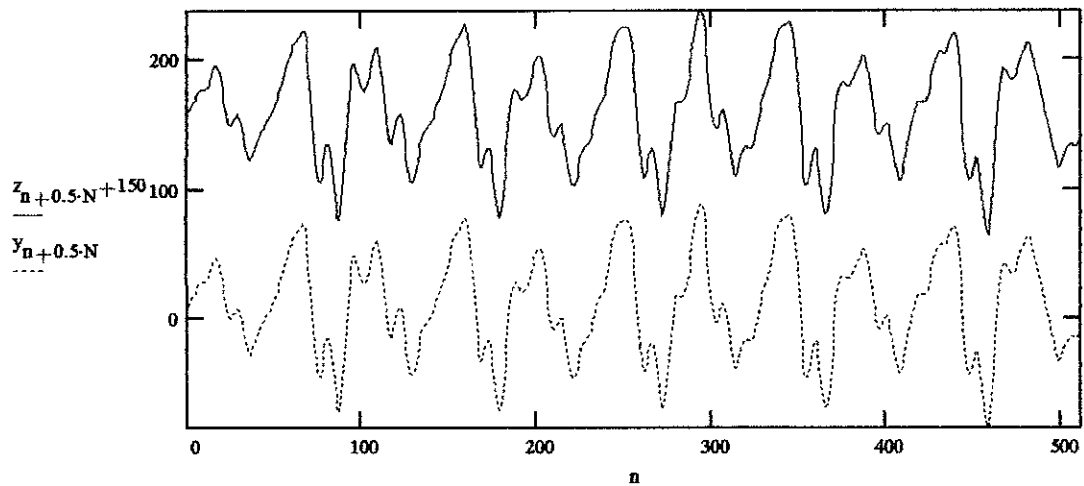
$$z := \frac{\text{ifft}(Z)}{\sqrt{N}}$$

Transformada inversa del producto de las transformadas:  
¡No es la convolución lineal!



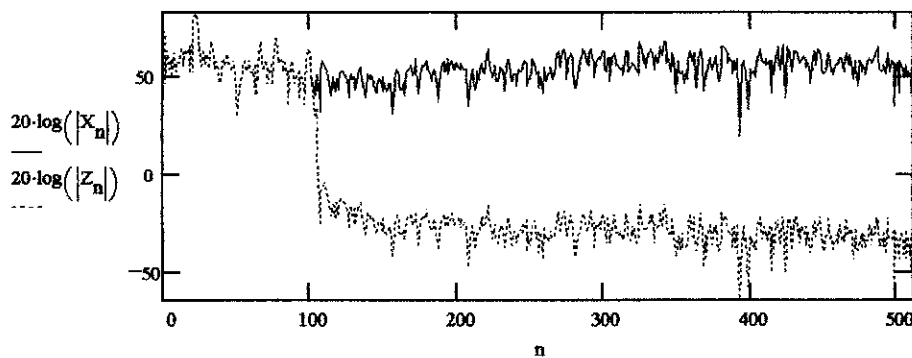
Se obtuvo la "Convolución Circular" debido a que el espectro discreto de señales muestreadas equivale a la serie de Fourier de señales periódicas muestreadas. Para calcular la convolución lineal mediante la FFT, se debe hacer un "relleno" de la señal y la respuesta al impulso del filtro para evitar la superposición de los supuestos periodos:

$k := N..2 \cdot N - 1$	Muestras adicionales
$x_k := 0$	Relleno de la señal
$h_k := 0$	Relleno de la respuesta al impulso
$X := \text{fft}(x) \cdot \sqrt{2 \cdot N}$	Nuevo espectro de amplitud de la señal
$H := \text{fft}(h) \cdot \sqrt{2 \cdot N}$	Nuevo espectro de la respuesta al impulso
$n := 0..N$	Indice de frecuencias
$Z_n := X_n \cdot H_n$	Espectro de la convolución lineal
$z := \frac{\text{ifft}(Z)}{\sqrt{2 \cdot N}}$	Salida del filtro
$n := 0..N - 1$	Indice de tiempo



Ahora si se obtuvo la verdadera salida del filtro. El cálculo mediante la FFT se hizo muchas veces más rápido que mediante la evaluación de la suma de convolución.

Comparación espectral a la entrada y la salida del filtro:



### 3. Modelo de producción de voz

Coefficientes del filtro inverso para el fonema /A/. Sonoro, con período de 85 muestras

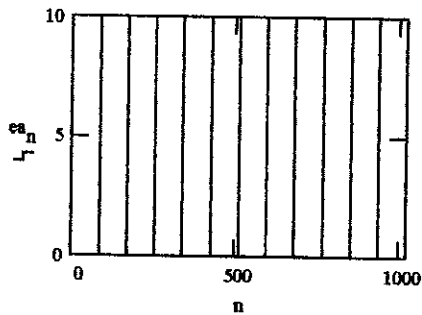
a := READPRN(c\_lpc\_a)

$$a = \begin{bmatrix} 1 \\ -2.326 \\ 1.652 \\ -0.3113 \\ -0.08789 \\ 0.6424 \\ -0.5074 \\ -0.5582 \\ 0.3833 \\ 0.3652 \\ -0.2332 \end{bmatrix}$$

n := 0..1023

Fuente de excitación sonora

$$ea_n := \text{if}(n \div 85 \cdot \text{floor}\left(\frac{n}{85}\right), 10, 0)$$



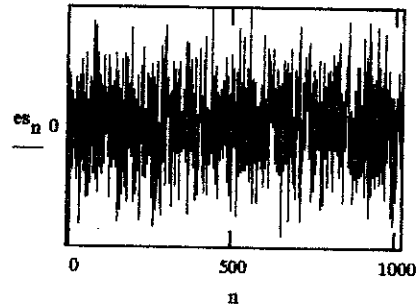
Coefficientes del filtro inverso para el fonema /S/. Sordo, con excitación en forma de ruido

s := READPRN(c\_lpc\_s)

$$s = \begin{bmatrix} 1 \\ 0.5311 \\ 0.2112 \\ -0.2393 \\ -0.3038 \\ -0.3328 \\ -0.2541 \\ -0.246 \\ -0.04053 \\ 0.01189 \\ 0.02777 \end{bmatrix}$$

Fuente de excitación sorda

$$es_n := \sqrt{-40 \cdot \ln(\text{rnd}(1))} \cdot \cos(2 \cdot \pi \cdot \text{rnd}(1))$$



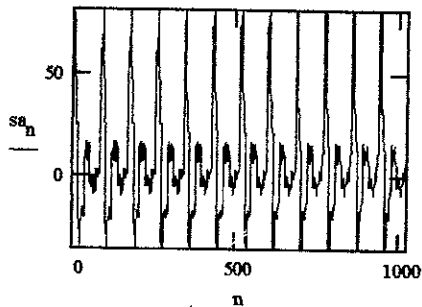
Síntesis del fonema /A/

sa\_0 := ea\_0

n := 1..1023      k := 1..10

$$sa_n := ea_n - \sum_k a_k \cdot \text{if}(n - k < 0, 0, sa_{n-k})$$

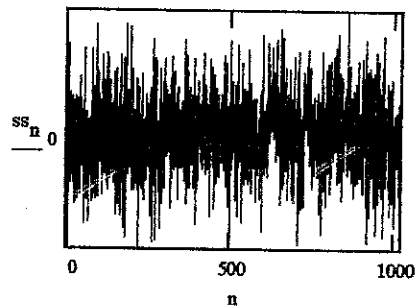
n := 0..1023



Síntesis del fonema /S/

ss\_0 := es\_0

$$ss_n := es_n - \sum_k s_k \cdot \text{if}(n - k < 0, 0, ss_{n-k})$$



## Síntesis de la palabra /ASA/

$$x_n := sa_n \cdot \left( 0.55 - 0.45 \cdot \cos\left(\frac{\pi \cdot n}{512}\right) \right)$$

Fonema /A/

$$n := 800..1023$$

$$x_n := x_n + \frac{n-800}{224} \cdot ss_{n-800}$$

Transición de /A/ a /S/

$$n := 224..1023$$

$$x_{n+800} := ss_n$$

Fonema /S/

$$n := 1600..1823$$

$$x_n := x_n \cdot \left( 1 - \frac{n-1600}{224} \right) + x_{n-1600}$$

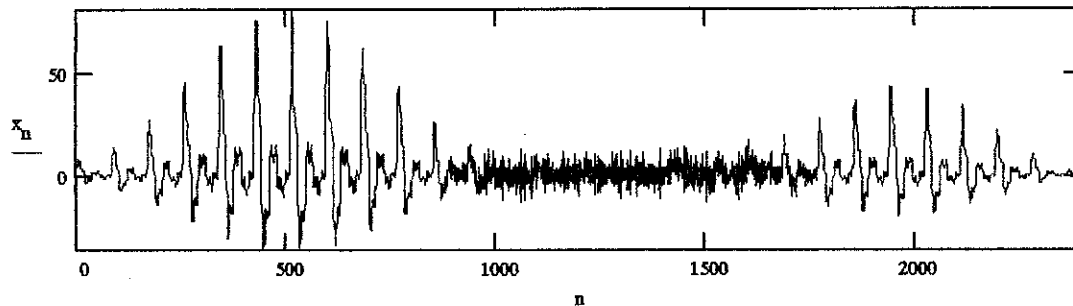
Transición de /S/ a /A/

$$n := 224..800$$

$$x_{n+1600} := x_n \cdot \frac{800-n}{676}$$

Fonema /A/ final

$$n := 0..2400$$

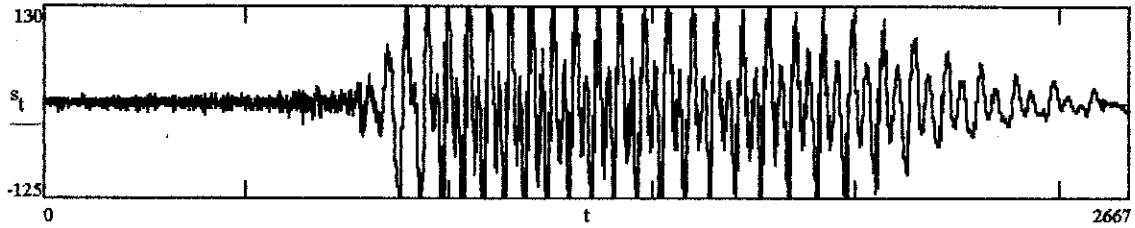


WRITEPRN(as) := x    Se almacena en un archivo .PRN para convertir a .WAV  
y escuchar la palabra sintetizada (ver sección 1.2)

#### 4.1 Energía promedio y Tasa promedio de cruces por cero

$s := \text{READPRN}(si)$   
 $t := 0.. \text{size}(s)$

Señal de voz a analizar: /SI/  
 Índice de tiempo



Se escogen 51 tramas de 100 muestras cada una, separadas entre si por 50 muestras:

$n := 0.. 99$

Índice dentro de cada Trama

$k := 0.. 50$

Número de tramas

$\text{signo}(x) := \text{if}(x < 0, -1, 1)$

Función 'Signo de x'

$$E_k := \frac{\sum_{n=0}^{99} (s_{n+50 \cdot k})^2}{100}$$

$$Z_k := \frac{\sum_{n=0}^{99} |\text{signo}(s_{n+50 \cdot k}) - \text{signo}(s_{n+50 \cdot k+1})|}{200}$$

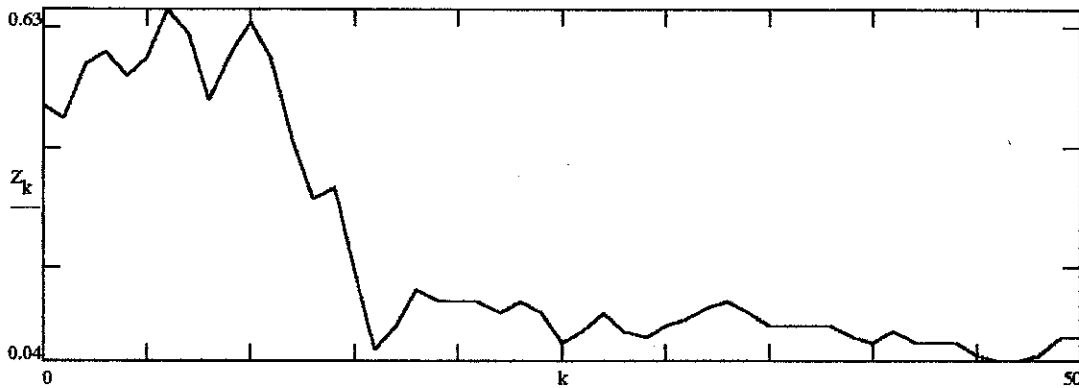
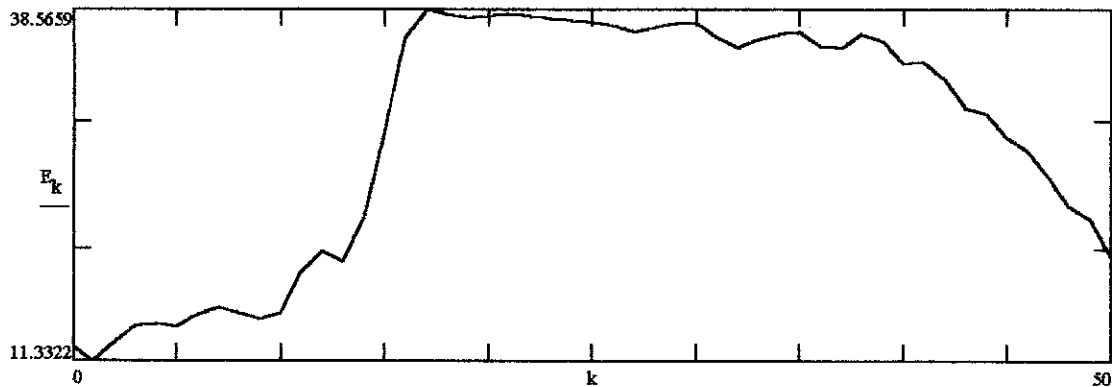
Función de Energía en  
 Tiempo Corto

Función de Tasa Promedio de Cruces por  
 Cero en Tiempo Corto

Debido al rango dinámico de la Energía promedio en tiempo corto, es conveniente visualizarla en una escala logarítmica:

$E_k := 10 \cdot \log(E_k)$

Energía en decibeles



Obsérvese cómo, durante el sonido sordo, la energía es baja y la tasa de cruces por cero es alta mientras que, durante el sonido sonoro, la energía es alta y la tasa promedio de cruces por cero es baja. Esto sugiere un algoritmo para determinar la condición sonoro/sordo:

Primero se normalizan las funciones de energía y cruces por cero:

$$e_m := \min(E)$$

$$z_m := \min(Z)$$

$$e_x := \max(E)$$

$$z_x := \max(Z)$$

$$E_k := \frac{E_k - e_m}{e_x - e_m}$$

$$Z_k := \frac{Z_k - z_m}{z_x - z_m}$$

Y ahora se comparan sus valores respecto a dos niveles:

$$s_k := -1$$

La clasificación de cada ventana se deja incierta

$$s_k := \text{if}(E_k > 0.66, \text{if}(Z_k < 0.66, 1, s_k), s_k)$$

Ventana sonora

$$s_k := \text{if}(Z_k > 0.66, \text{if}(E_k < 0.66, 0, s_k), s_k)$$

Ventana sorda

$$s_k := \text{if}(E_k > 0.33, \text{if}(Z_k < 0.33, 1, s_k), s_k)$$

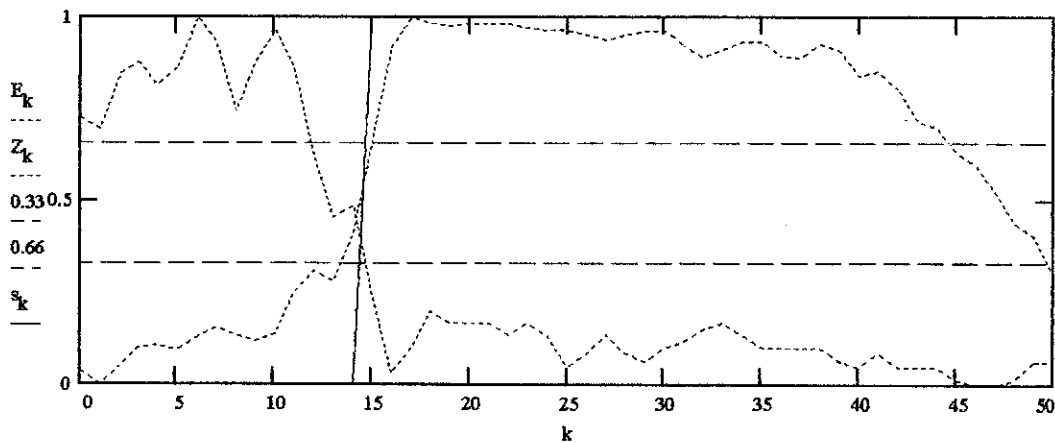
Ventana sonora

$$s_k := \text{if}(Z_k > 0.33, \text{if}(E_k < 0.33, 0, s_k), s_k)$$

Ventana sorda

$$s_k := \text{if}(s_k < 0, s_{k-1}, s_k)$$

Se deja un ciclo de histéresis entre distintas clasificaciones

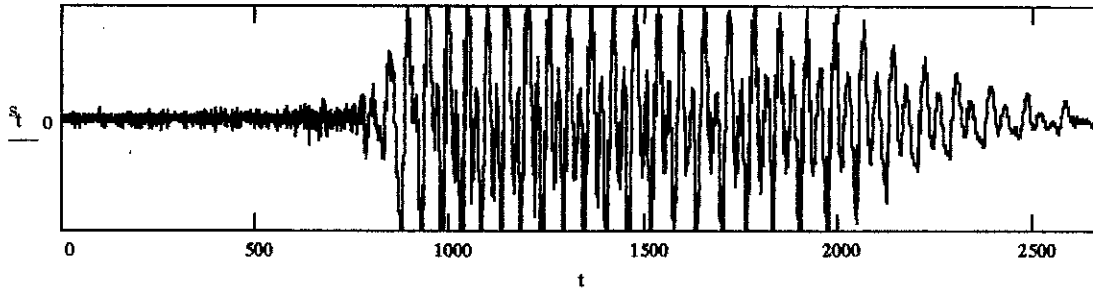


Hasta la ventana 14 la señal se clasifica como sorda y, desde la ventana 15, como sonora.

## 4.2 Autocorrelación en tiempo corto

$s := \text{READPRN}(si)$   
 $t := 0.. \text{size}(s)$

Señal de voz a analizar  
 Índice de tiempo



Ahora se calculará la función de autocorrelación en un segmento sordo y en un segmento sonoro de la señal de voz:

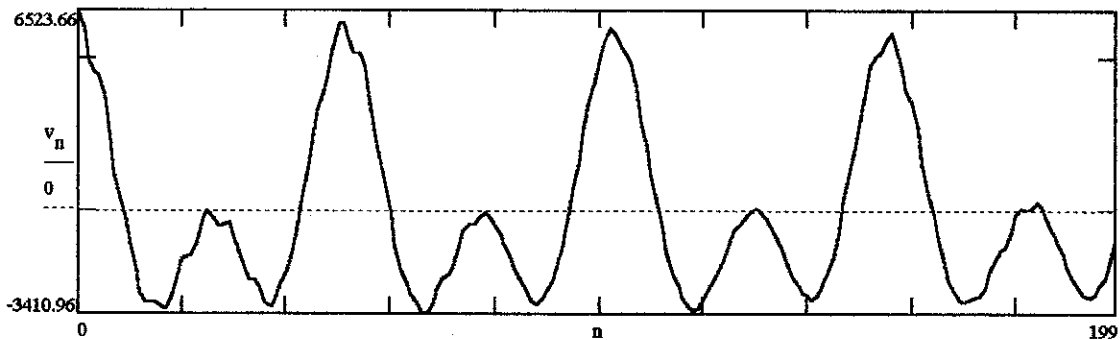
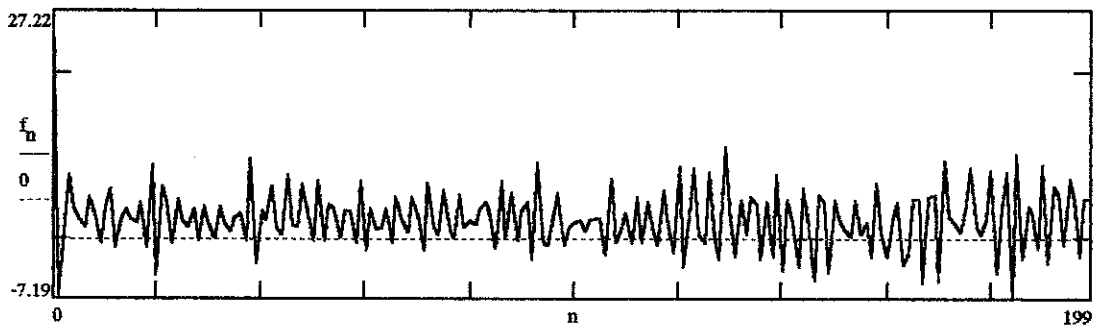
$k := 0..99$       $n := 0..199$

$$f_n := \sum_k \frac{s_{200+k} s_{200+k+n}}{100}$$

$$v_n := \sum_k \frac{s_{1000+k} s_{1000+k+n}}{100}$$

Autocorrelación del segmento sordo  
 (fricativo) entre las muestras 200 y 400

Autocorrelación del segmento sonoro  
 (vocalizado) entre las muestras 1000 y 1200

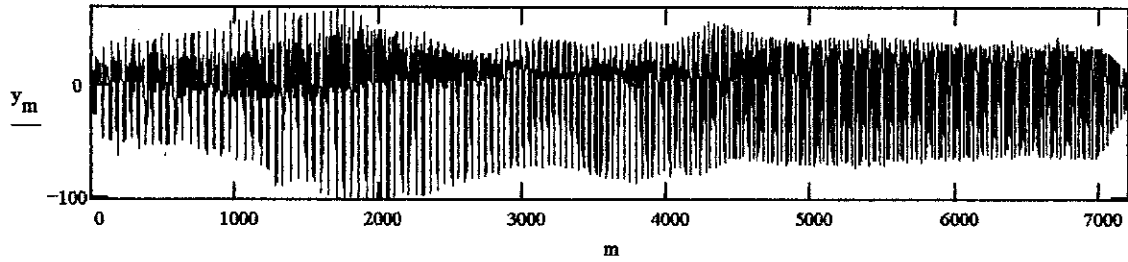


En el primer caso no hay evidencia de periodicidad, mientras en el segundo hay un pico claro cada período de la señal de voz. Esta es una indicación de sonoridad pero, principalmente, sugiere un método para detección de tono, como se explora en la sección 4.3



### 4.3 Estimación de tono

$y := \text{READPRN}(\text{do\_mayor})$  Señal de voz a analizar  
 $m := 0.. \text{size}(y)$  Índice de tiempo



Primero se pasa la señal por un filtro pasabajos a 375 Hz. Por tratarse de una frecuencia de muestreo de 6 kHz, el filtro es de 1/8 de banda:

$$\text{sinc}(\alpha) := \text{if}(\alpha=0, 1, \frac{\sin(\pi \cdot \alpha)}{\pi \cdot \alpha})$$

Definición de la función SENO-C

$k := 0.. 19$

Índice de tiempo

$$h_k := (0.54 - 0.46 \cdot \cos(0.1 \cdot \pi \cdot k)) \cdot \text{sinc}((k - 10) \cdot 0.125) \cdot 0.125$$

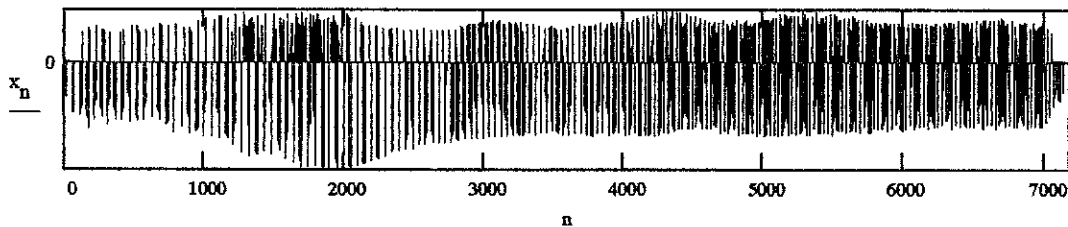
Filtro pasabajos de 1/4 banda

$$n := 0.. \text{size}(y) - 20 \quad x_n := \sum_k h_k \cdot y_{n+20-k}$$

Se filtra la señal

$$x_n := \text{if}(|x_n| < 20, 0, x_n)$$

Y se recorta para "aplanamiento" espectral



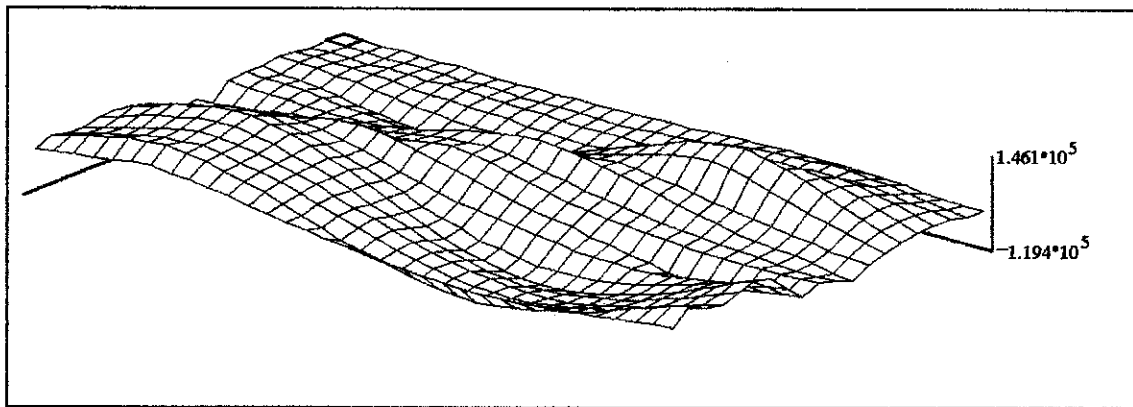
$k := 0.. 19$   
 $i := 0.. 159$   
 $n := 26.. 59$

Se tomarán 20 ventanas de 160 muestras cada una

y se considerarán posibles tonos entre 100 y 230 Hz ( $6000/59..6000/26$ )

$$T_{k,n-26} := \sum_i x_{365 \cdot k + i} \cdot x_{365 \cdot k + i + n}$$

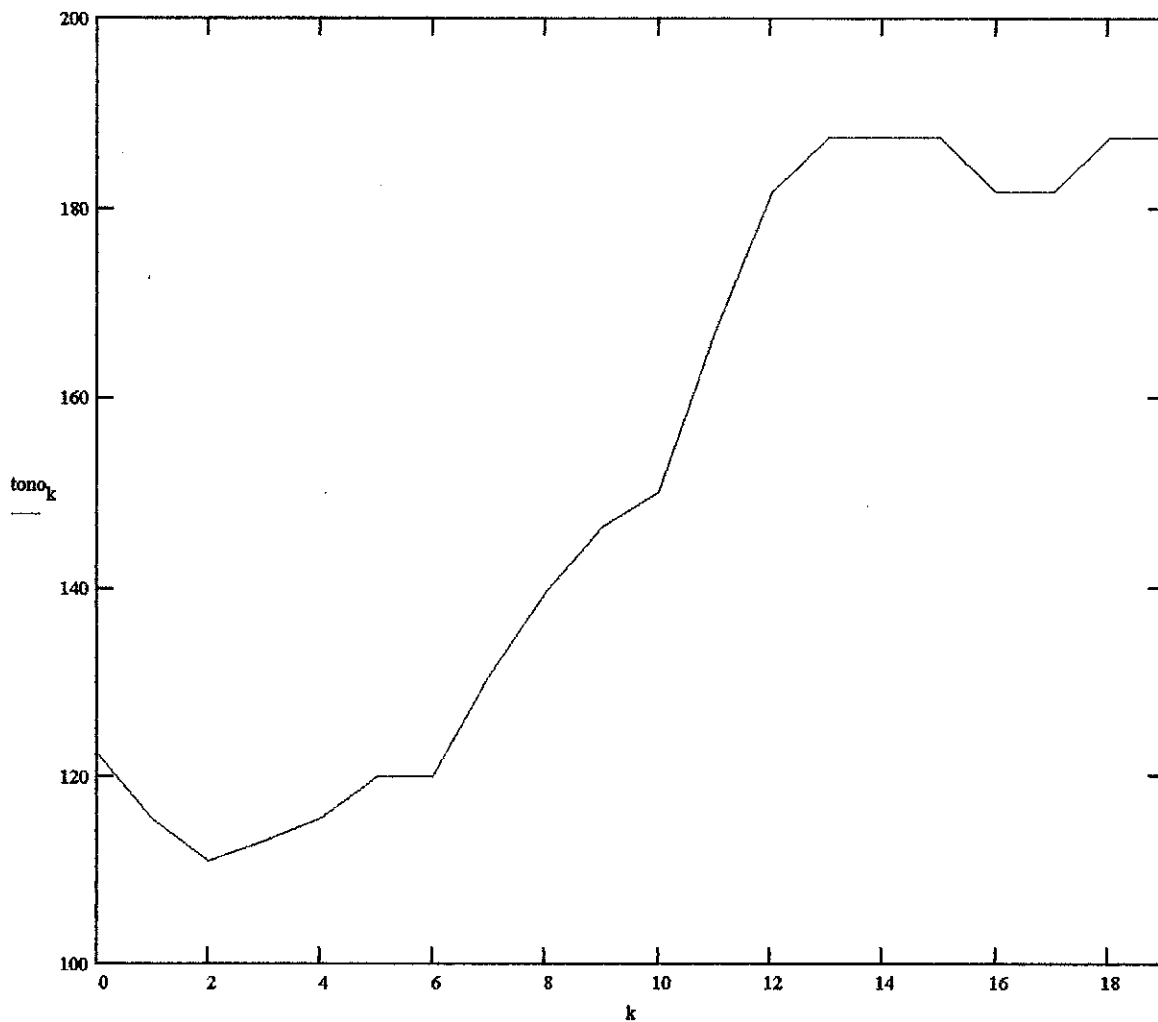
Se espera que la función de autocorrelación presente máximos cada período.



T

$n := 0..33$				
$i := 0$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	
$i := 1$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 1
$i := 2$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 2
$i := 3$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 3
$i := 4$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 4
$i := 5$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 5
$i := 6$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 6
$i := 7$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 7
$i := 8$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 8
$i := 9$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 9
$i := 10$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 10
$i := 11$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 11
$i := 12$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 12
$i := 13$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 13
$i := 14$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 14
$i := 15$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 15
$i := 16$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 16
$i := 17$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 17
$i := 18$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 18
$i := 19$	$t_n := T_{i,n}$	$m := \max(t)$	$\text{tono}_i := \frac{n \cdot 6000}{26 + \sum \text{if}(m=t_n, n, 0)}$	Ventana 19

El instante en que se presente el máximo de cada ventana indica su período. El tono está dado por su inverso.



Se trató de cantar la secuencia de la tríada correspondiente a Do Mayor (do - mi - sol). Así, la función de tono debió ser una escalera con escalones en 131, 165 y 196 Hz. El programa muestra que sólo se alcanzó a cantar la secuencia de Si Mayor ligeramente desafinada...

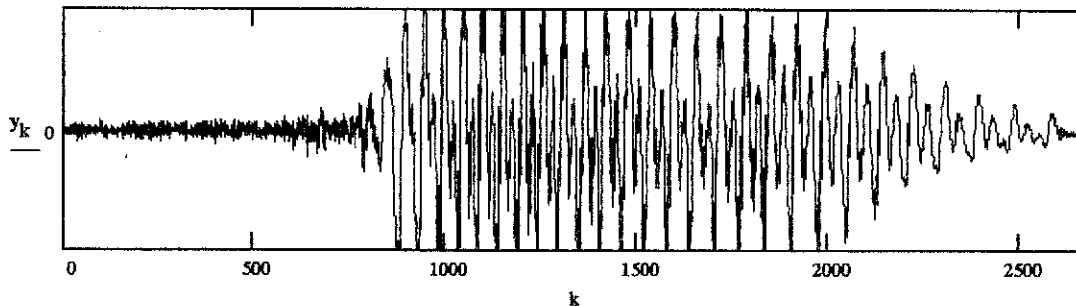
#### 4.4 Transformada de Fourier en tiempo corto

y := READPRN(si)

Señal a analizar

k := 0..size(y)

Índice de tiempo



Se construirán 36 ventanas de 256 muestras cada una, separadas 64 muestras entre si.

j := 0..35

Número de ventana

i := 0..90

Índice de frecuencia (de 0 a 3900 Hz)

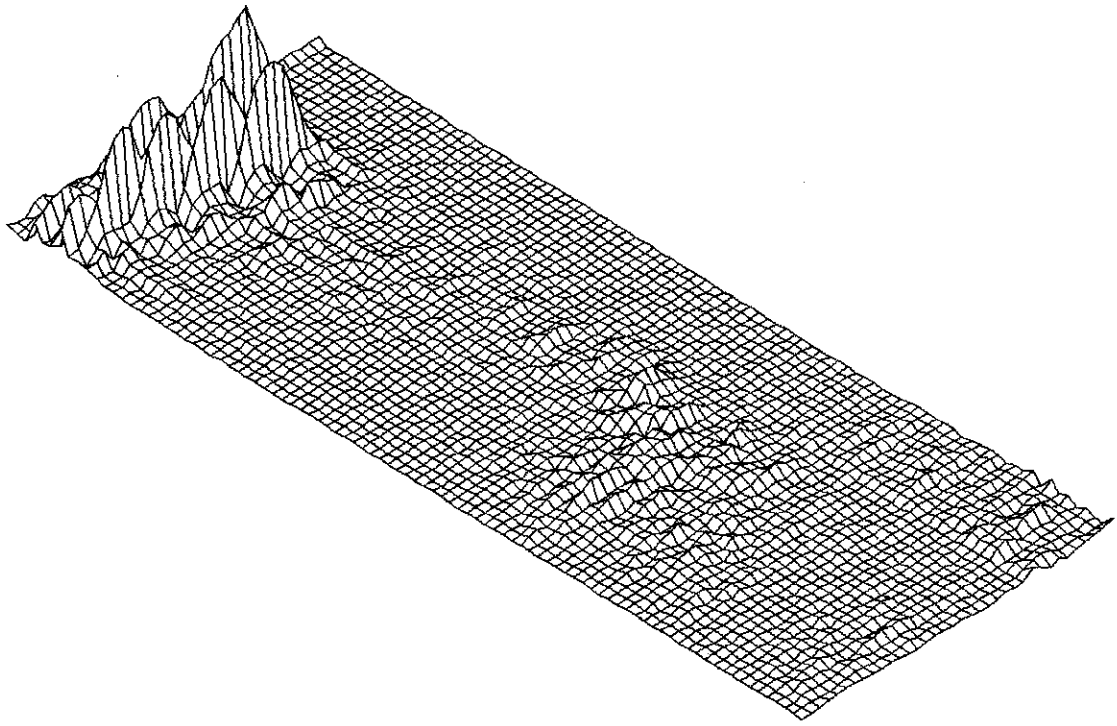
n := 0..255

Índice de tiempo en cada ventana

Construye cada ventana, calcula su FFT y almacena la amplitud en un arreglo bidimensional

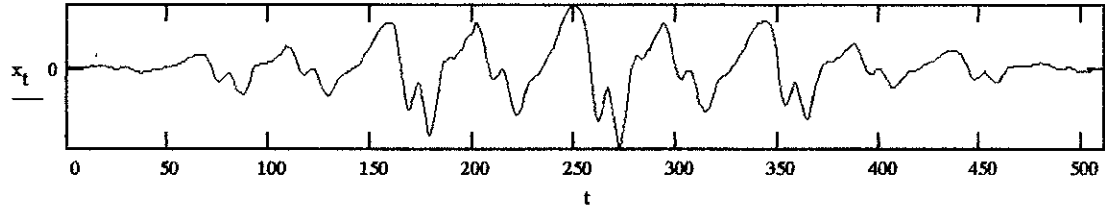
$x_n := y_{64 \cdot 0 + n}$	$X := \text{fft}(x)$	$S_{1,0} :=  X_i $	$x_n := y_{64 \cdot 1 + n}$	$X := \text{fft}(x)$	$S_{1,1} :=  X_i $
$x_n := y_{64 \cdot 2 + n}$	$X := \text{fft}(x)$	$S_{1,2} :=  X_i $	$x_n := y_{64 \cdot 3 + n}$	$X := \text{fft}(x)$	$S_{1,3} :=  X_i $
$x_n := y_{64 \cdot 4 + n}$	$X := \text{fft}(x)$	$S_{1,4} :=  X_i $	$x_n := y_{64 \cdot 5 + n}$	$X := \text{fft}(x)$	$S_{1,5} :=  X_i $
$x_n := y_{64 \cdot 6 + n}$	$X := \text{fft}(x)$	$S_{1,6} :=  X_i $	$x_n := y_{64 \cdot 7 + n}$	$X := \text{fft}(x)$	$S_{1,7} :=  X_i $
$x_n := y_{64 \cdot 8 + n}$	$X := \text{fft}(x)$	$S_{1,8} :=  X_i $	$x_n := y_{64 \cdot 9 + n}$	$X := \text{fft}(x)$	$S_{1,9} :=  X_i $
$x_n := y_{64 \cdot 10 + n}$	$X := \text{fft}(x)$	$S_{1,10} :=  X_i $	$x_n := y_{64 \cdot 11 + n}$	$X := \text{fft}(x)$	$S_{1,11} :=  X_i $
$x_n := y_{64 \cdot 12 + n}$	$X := \text{fft}(x)$	$S_{1,12} :=  X_i $	$x_n := y_{64 \cdot 13 + n}$	$X := \text{fft}(x)$	$S_{1,13} :=  X_i $
$x_n := y_{64 \cdot 14 + n}$	$X := \text{fft}(x)$	$S_{1,14} :=  X_i $	$x_n := y_{64 \cdot 15 + n}$	$X := \text{fft}(x)$	$S_{1,15} :=  X_i $
$x_n := y_{64 \cdot 16 + n}$	$X := \text{fft}(x)$	$S_{1,16} :=  X_i $	$x_n := y_{64 \cdot 17 + n}$	$X := \text{fft}(x)$	$S_{1,17} :=  X_i $
$x_n := y_{64 \cdot 18 + n}$	$X := \text{fft}(x)$	$S_{1,18} :=  X_i $	$x_n := y_{64 \cdot 19 + n}$	$X := \text{fft}(x)$	$S_{1,19} :=  X_i $
$x_n := y_{64 \cdot 20 + n}$	$X := \text{fft}(x)$	$S_{1,20} :=  X_i $	$x_n := y_{64 \cdot 21 + n}$	$X := \text{fft}(x)$	$S_{1,21} :=  X_i $
$x_n := y_{64 \cdot 22 + n}$	$X := \text{fft}(x)$	$S_{1,22} :=  X_i $	$x_n := y_{64 \cdot 23 + n}$	$X := \text{fft}(x)$	$S_{1,23} :=  X_i $
$x_n := y_{64 \cdot 24 + n}$	$X := \text{fft}(x)$	$S_{1,24} :=  X_i $	$x_n := y_{64 \cdot 25 + n}$	$X := \text{fft}(x)$	$S_{1,25} :=  X_i $
$x_n := y_{64 \cdot 26 + n}$	$X := \text{fft}(x)$	$S_{1,26} :=  X_i $	$x_n := y_{64 \cdot 27 + n}$	$X := \text{fft}(x)$	$S_{1,27} :=  X_i $
$x_n := y_{64 \cdot 28 + n}$	$X := \text{fft}(x)$	$S_{1,28} :=  X_i $	$x_n := y_{64 \cdot 29 + n}$	$X := \text{fft}(x)$	$S_{1,29} :=  X_i $
$x_n := y_{64 \cdot 30 + n}$	$X := \text{fft}(x)$	$S_{1,30} :=  X_i $	$x_n := y_{64 \cdot 31 + n}$	$X := \text{fft}(x)$	$S_{1,31} :=  X_i $
$x_n := y_{64 \cdot 32 + n}$	$X := \text{fft}(x)$	$S_{1,32} :=  X_i $	$x_n := y_{64 \cdot 33 + n}$	$X := \text{fft}(x)$	$S_{1,33} :=  X_i $
$x_n := y_{64 \cdot 34 + n}$	$X := \text{fft}(x)$	$S_{1,34} :=  X_i $	$x_n := y_{64 \cdot 35 + n}$	$X := \text{fft}(x)$	$S_{1,35} :=  X_i $

La gráfica del arreglo S se conoce como espectrograma y en él se nota cómo, durante el fonema /S/, prevalecen las altas frecuencias mientras que, durante el fonema /I/, prevalecen las bajas frecuencias con mucho mayor energía.

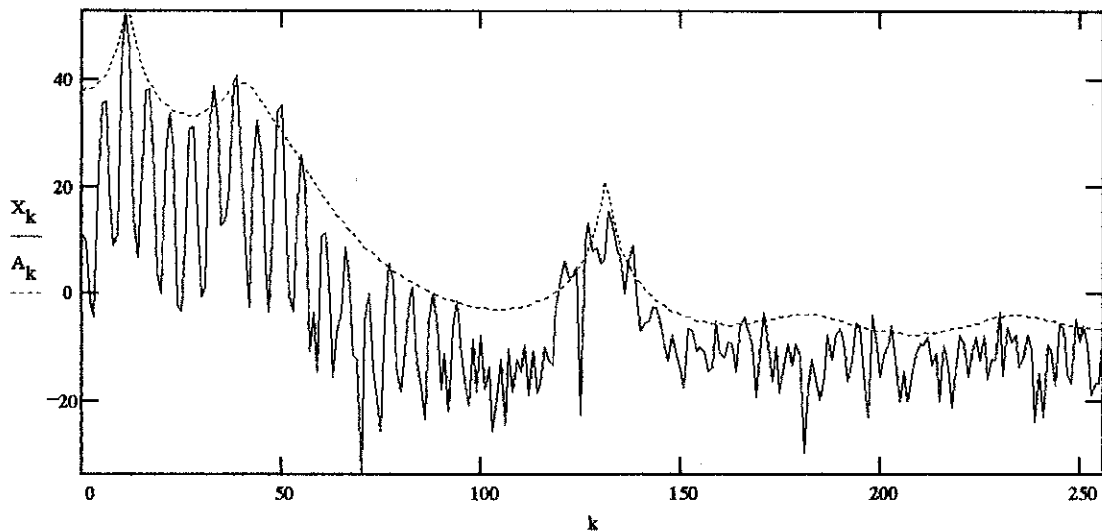


## 5. Análisis predictivo lineal, LPC

$x := \text{READPRN}(a)$  Muestras de la señal de voz  
 $N := (\text{size}(x) + 1)$  Número de muestras  
 $t := 0..N - 1$  Índice de tiempo  
 $x_t := (x_t) \cdot \left[ 0.54 + 0.46 \cdot \cos \left[ 2 \cdot \pi \cdot \left( \frac{t}{N} - 0.5 \right) \right] \right]$  Ventana Hamming



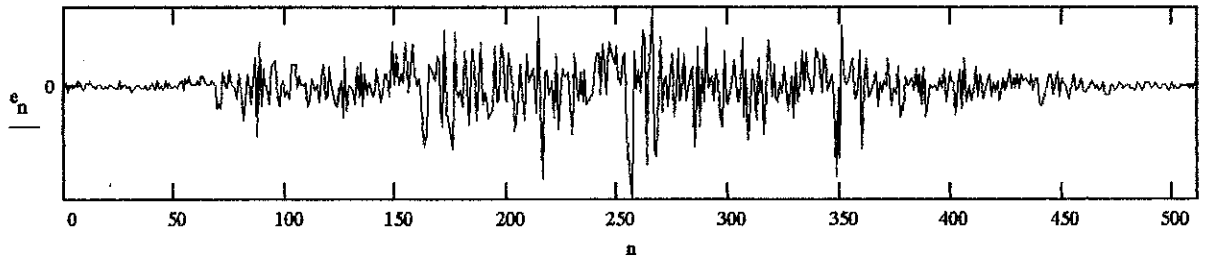
$p := 12$  Orden del filtro predictor  
 $\tau := 0..p$  Índice de la función de autocorrelación  
 $r_\tau := \sum_t x_t \cdot \text{if}(t + \tau < N, x_{t+\tau}, 0)$  Función de autocorrelación  
 $i := 0..p - 1$        $j := 0..p - 1$  Índices de la matriz de autocorrelación  
 $R_{i,j} := r_{|i-j|}$  Matriz de autocorrelación  
 $\rho := \bar{r}_{i+1}$  Vector de autocorrelación  
 $R := R^{-1}$  Inversa de la matriz de autocorrelación  
 $a := R \cdot \rho$   
 $a_{i+1} := -a_i$        $a_0 := 1$  Vector de Coeficientes predictores  
 $k := p + 1..511$  Rellena con ceros para calcular la FFT  
 $\alpha_k := 0$   
 $A := \text{fft}(\alpha)$  Respuesta en frecuencia del filtro inverso  
 $k := 0..256$  Índice en la frecuencia  
 $A_k := 20 \cdot \log \left( \left| \frac{1}{A_k} \right| \right) - 20$  Respuesta en frecuencia del filtro directo  
 $X := \text{fft}(x)$  Composición espectral de la señal original  
 $X_k := 20 \cdot \log(|X_k|)$



Los formantos (frecuencias de resonancia de la cavidad oral) se encuentran en  $k=11, 40$  y  $131$ , correspondiendo a  $11025 \cdot 11/512 = 237$  Hz,  $11025 \cdot 40/512 = 861$  Hz y  $11025 \cdot 131/512 = 2821$  Hz.

$$n := 0..511 \quad i := 0..p \quad e_n := \sum_i \alpha_i \cdot \text{if}(n-i \geq 0, x_{n-i}, 0)$$

Error de predicción



$$k := 0..40$$

$$h_k := \text{if} \left[ k=20, 0.125, 0.125 \cdot \frac{\sin(0.125 \cdot \pi \cdot (k-20))}{0.125 \cdot \pi \cdot (k-20)} \right] \cdot \left( 0.54 - 0.46 \cdot \cos\left(k \cdot \frac{\pi}{20}\right) \right)$$

Filtro pasabajos de 1/8 banda

$$n := 0..511$$

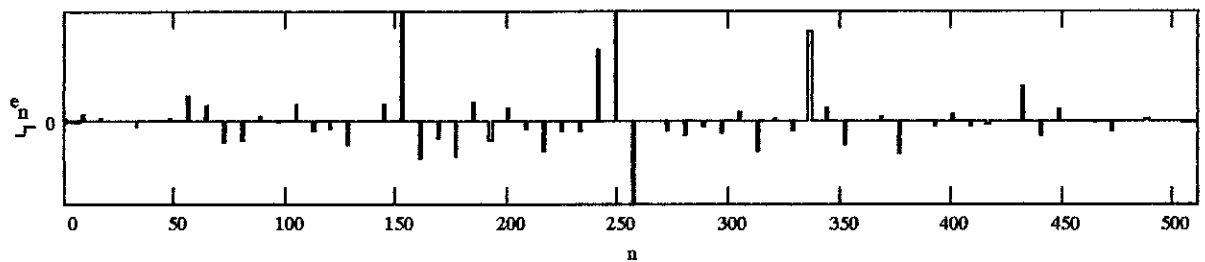
$$y_n := \sum_k h_k \cdot \text{if}(n-k < 0, 0, \text{if}(n-k > 511, 0, e_{n-k}))$$

Error filtrado

$$n := 0..511$$

$$e_n := \text{if}(n \neq 8 \cdot \text{floor}\left(\frac{n}{8}\right), y_{n+20} \cdot 8, 0)$$

Error submuestreado

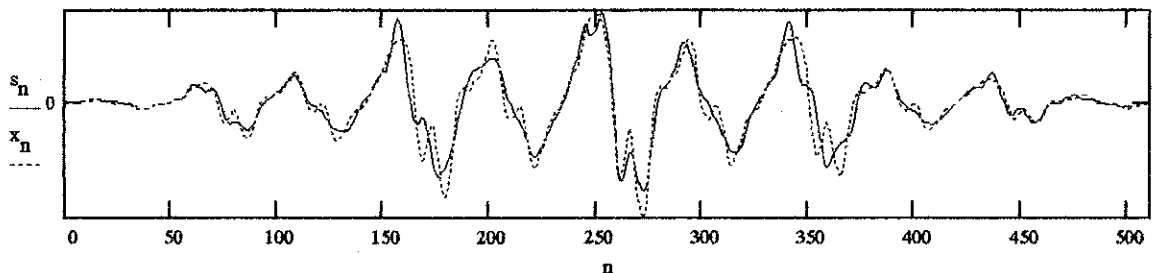


$$i := 1..p \quad n := 1..511$$

$$s_0 := e_0 \quad s_n := e_n - \left( \sum_i \alpha_i \cdot \text{if}(n-i \geq 0, s_{n-i}, 0) \right)$$

Síntesis de la señal con el error filtrado y submuestreado por 8

$$n := 0..511$$



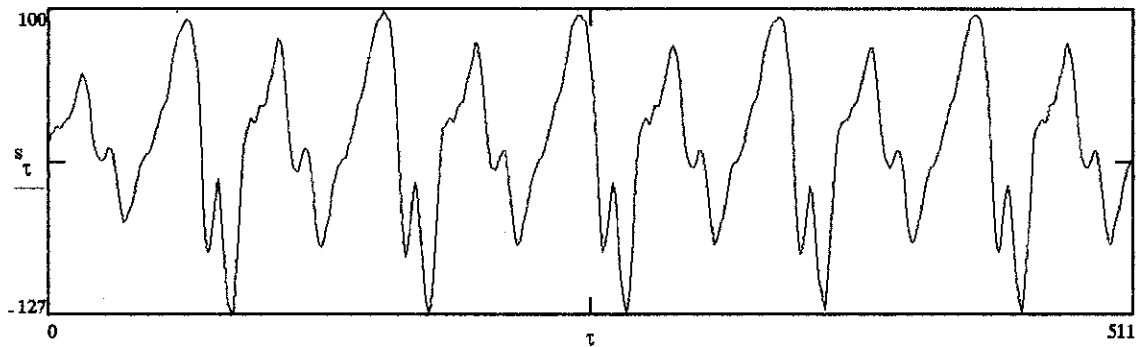
Obsérvese cómo, de las 512 muestras originales, se puede reproducir el segmento con suficiente exactitud a partir de  $64 + 12 = 76$  muestras

## 6. Análisis cepstral - Cepstro real

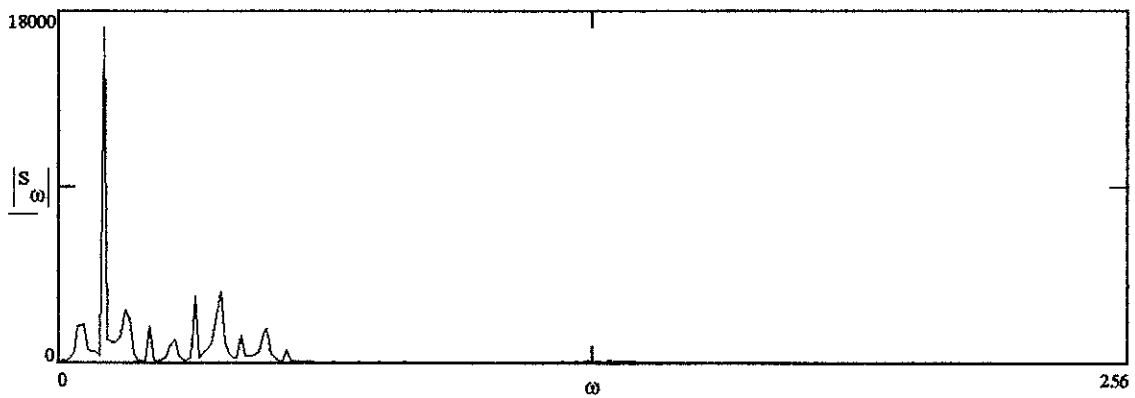
$s := \text{READPRN}(a)$  Muestras de la señal de voz

$L := \text{size}(s) + 1$  Número de muestras

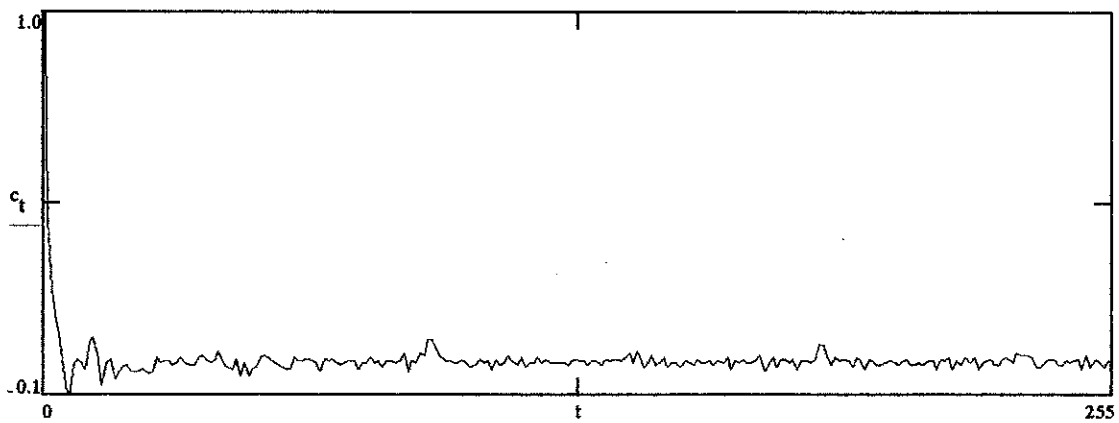
$\tau := 0..L - 1$  Índice de tiempo



$S := \text{fft}(s) \cdot \sqrt{L}$        $\omega := 0.. \frac{L}{2}$       Espectro de la señal original



$S_\omega := \log(|S_\omega|)$        $c := \frac{\text{ifft}(S)}{\sqrt{L}}$       Cepstro real de la señal



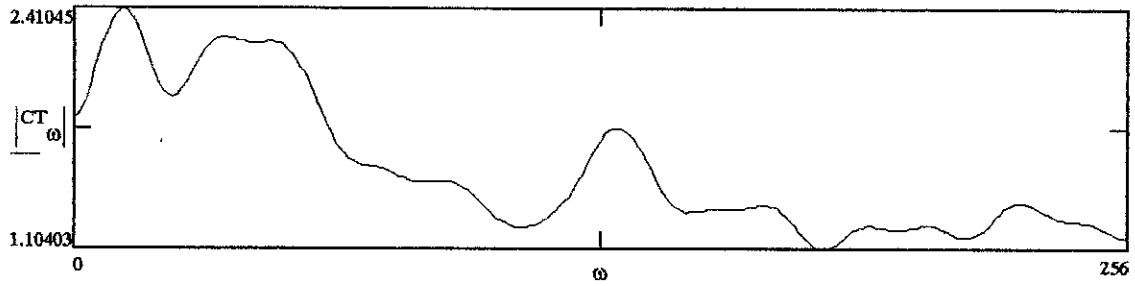


ct := c    i := 30..511    ct<sub>i</sub> := 0

"Liftraje" para separar la envolvente espectral

CT := fft(ct) · √L

Respuesta frecuencia del tracto oral. Compárese con la obtenida mediante LPC en la sección 5.



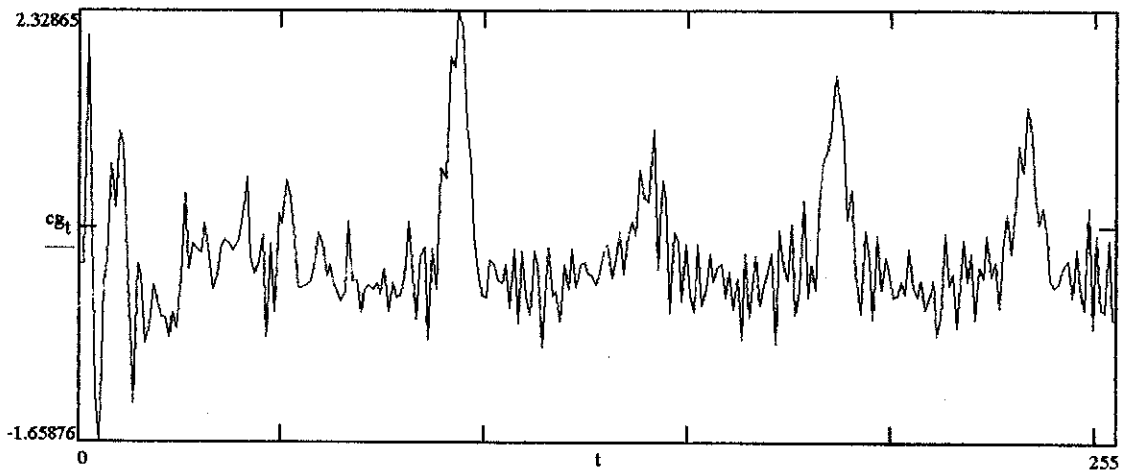
cg := c    i := 0..72    cg<sub>i</sub> := 0

"Liftraje" para separar la excitación glotal

CG := fft(cg) · √L

CG<sub>omega</sub> := exp(|CG<sub>omega</sub>|)    cg := ifft(CG)

Estimación de la excitación glotal



El pulso de periodicidad aparece cada 93 muestras, indicando una frecuencia fundamental de  $11025/93 = 118.5$  Hz

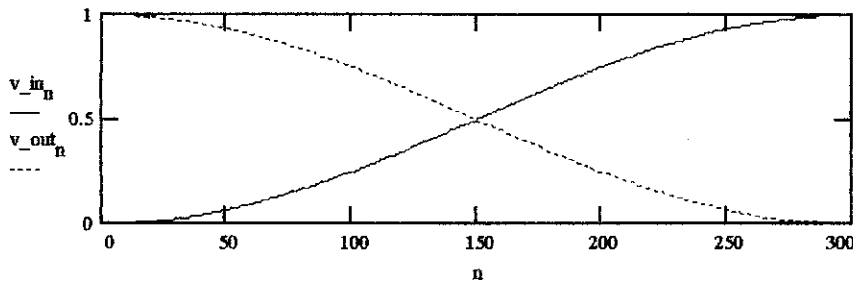
## 8. Síntesis de voz mediante concatenación de fonemas

$n := 0..300$

$$v_{out_n} := 0.5 + 0.5 \cdot \cos\left(\frac{\pi \cdot n}{300}\right)$$

$$v_{in_n} := 1 - v_{out_n}$$

La transición entre fonemas se hará mediante la superposición de los extremos de los fonemas, ponderados por estas ventanas.



$y := \text{READPRN}(ss)$

Lee el fonema /S/

$x := y$

Lo "concatena" con el anterior

$y := \text{READPRN}(i)$

Lee el siguiente fonema

$n := 900..1200$

Región de transición entre fonemas

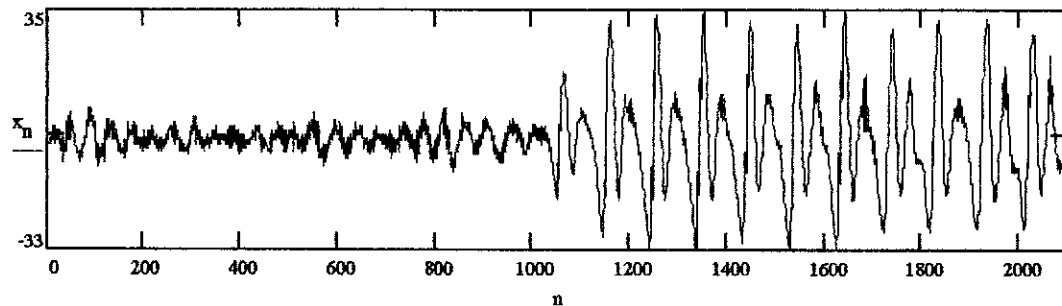
$$x_n := x_n \cdot v_{out_{n-900}} + y_{n-900} \cdot v_{in_{n-900}} \quad \text{Paso de un fonema a otro}$$

$n := 1200..2100$

Resto del fonema "entrante"

$$x_n := y_{n-900}$$

$n := 0..2100$



Repite el proceso para los demás fonemas:

$y := \text{READPRN}(nn) \quad n := 1800..2100 \quad x_n := x_n \cdot v_{out_{n-1800}} + y_{n-1800} \cdot v_{in_{n-1800}}$

$n := 2100..3000 \quad x_n := y_{n-2100}$

$y := \text{READPRN}(tt) \quad n := 2700..3000 \quad x_n := x_n \cdot v_{out_{n-2700}} + y_{n-2700} \cdot v_{in_{n-2700}}$

$n := 3000..3900 \quad x_n := y_{n-3000}$

$y := \text{READPRN}(ee) \quad n := 3600..3900 \quad x_n := x_n \cdot v_{out_{n-3600}} + y_{n-3600} \cdot v_{in_{n-3600}}$

$n := 3900..4800 \quad x_n := y_{n-3900}$

$y := \text{READPRN}(ss) \quad n := 4500..4800 \quad x_n := x_n \cdot v_{out_{n-4500}} + y_{n-4500} \cdot v_{in_{n-4500}}$

$n := 4800..5700 \quad x_n := y_{n-4800}$

$y := \text{READPRN}(ii) \quad n := 5400..5700 \quad x_n := x_n \cdot v\_out_{n-5400} + y_{n-5400} \cdot v\_in_{n-5400}$

$n := 5700..6600 \quad x_n := y_{n-5700}$

$y := \text{READPRN}(ss) \quad n := 6300..6600 \quad x_n := x_n \cdot v\_out_{n-6300} + y_{n-6300} \cdot v\_in_{n-6300}$

$n := 6300..7500 \quad x_n := y_{n-6300}$

Ahora hace la "concatenación" con supuestos silencios al comienzo y al final de la palabra

$n := 0..300 \quad x_n := x_n \cdot v\_in_n$

$n := 7200..7500 \quad x_n := x_n \cdot v\_out_{n-7200}$

Y, por último, pone el acento en la primera sílaba

$n := 0..1000$

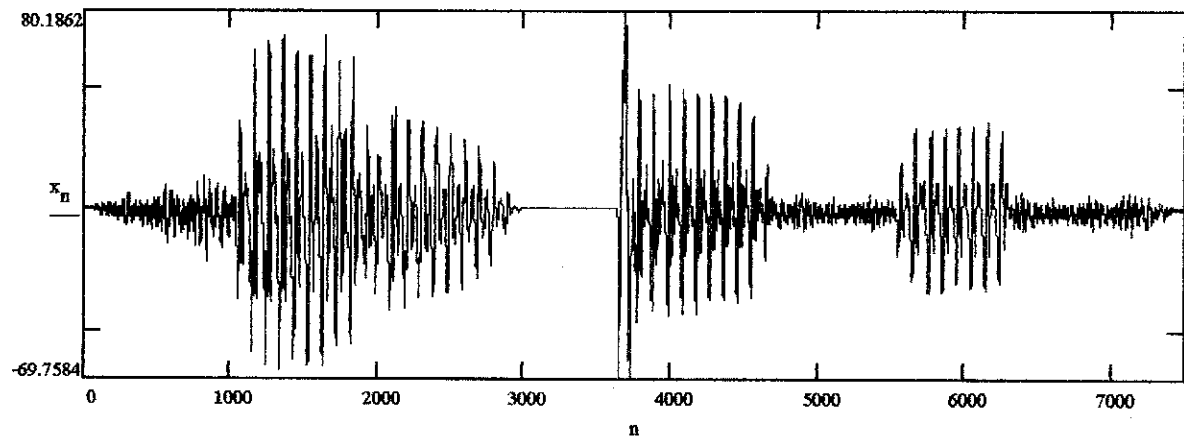
$$x_n := x_n \cdot \left( 1.5 - 0.5 \cdot \cos\left(\frac{\pi \cdot n}{1000}\right) \right)$$

$$x_{n+1000} := 2 \cdot x_{n+1000}$$

$$x_{n+2000} := x_{n+2000} \cdot \left( 1.5 + 0.5 \cdot \cos\left(\frac{\pi \cdot n}{1000}\right) \right)$$

$n := 0..7500$

Forma de onda de la palabra /SINTESIS/ sintetizada



$\text{WRITEPRN}(\text{sinthesis}) := x$

Se almacena en un archivo .PRN para convertirlo a .WAV  
y escucharlo mediante una tarjeta de audio

## 9. Alineación Dinámica de Tiempo -DTW-

$a := \text{READPRN}(\text{cero1})$  Señal "patrón"

$b := \text{READPRN}(\text{cero2})$  Señal de prueba

$I := \text{floor}\left(\frac{\text{size}(a) + 1}{128} - 1\right)$  Número de tramas en la señal patrón  $I = 49$

$J := \text{floor}\left(\frac{\text{size}(b) + 1}{128} - 1\right)$  Número de tramas en la señal de prueba  $J = 36$

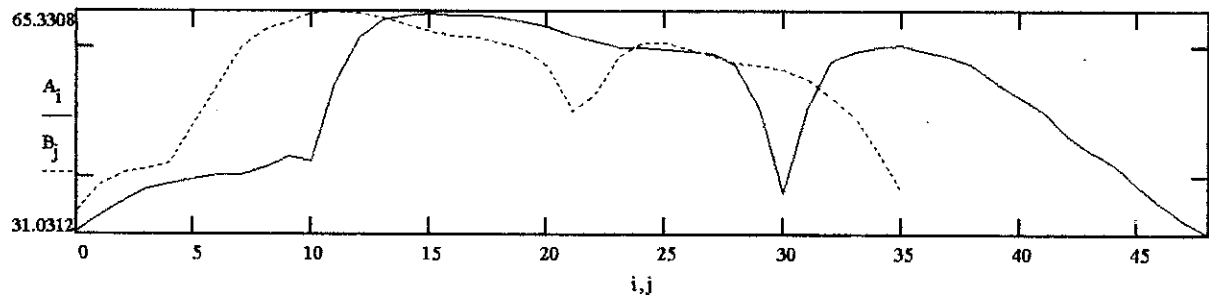
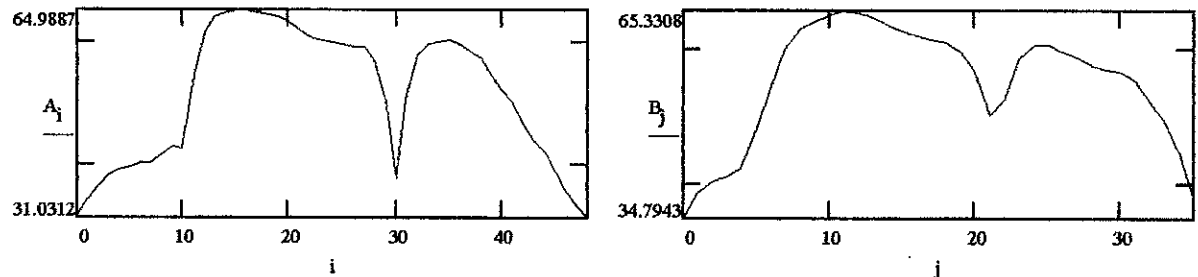
$i := 0..I - 1$  Índice de la Función de Energía en Tiempo Corto (patrón)

$j := 0..J - 1$  Índice de la Función de Energía en Tiempo Corto (prueba)

$n := 0..255$  Índice de las muestras de cada trama

$A_i := \sum_n (a_{128i+n})^2$   $A_i := 10 \cdot \log(A_i)$  Función de Energía en Tiempo Corto en dB (patrón)

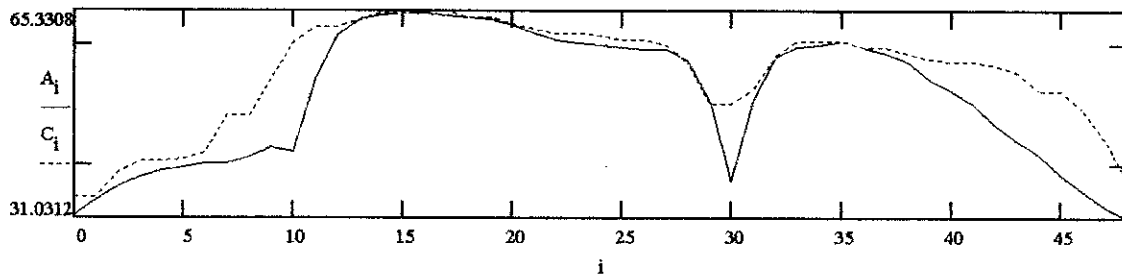
$B_j := \sum_n (b_{128j+n})^2$   $B_j := 10 \cdot \log(B_j)$  Función de Energía en Tiempo Corto en dB (prueba)



Aunque ambas corresponden a la función de energía de la palabra /cero/, sería absurdo calcular el error entre ellas, muestra a muestra, sin hacer antes alguna alineación:

$x_i := \text{floor}\left(i \cdot \frac{J-1}{I-1}\right)$  Primero se intenta una alineación lineal y se observa como ahora los patrones se hacen más "comparables".

$C_i := B_{(x_i)}$  Sin embargo, en algunos puntos, como  $i=10$  e  $i=45$ , el error sigue siendo innecesariamente grande.



$$E2 := \frac{1}{I} \sum_i (A_i - C_i)^2 \quad E2 = 36.34302$$

El método de alineación escogido es, entonces, aquel que optimiza el error promedio entre el patrón y la prueba. La mejor técnica de optimización es la Programación Dinámica:

$d_{i,j} := (A_i - B_j)^2$  Se asigna una "distancia" (la diferencia de Energías) a cada punto (i,j) y del estudio de esta matriz se encuentra la siguiente "ruta óptima":

$$\text{min2}(a, b) := \text{if}(a < b, a, b)$$

$$\text{min4}(a, b, c, d) := \text{min2}(\text{min2}(a, b), \text{min2}(c, d))$$

$$\text{min5}(a, b, c, d, e) := \text{min2}(\text{min4}(a, b, c, d), e)$$

$$f(a, b, c, d, e) := \text{if}(\text{min5}(a, b, c, d, e) = a, 0, \text{if}(\text{min5}(a, b, c, d, e) = b, 0, \text{if}(\text{min5}(a, b, c, d, e) = d, 1, \text{if}(\text{min5}(a, b, c, d, e) = c, 1, 2))))$$

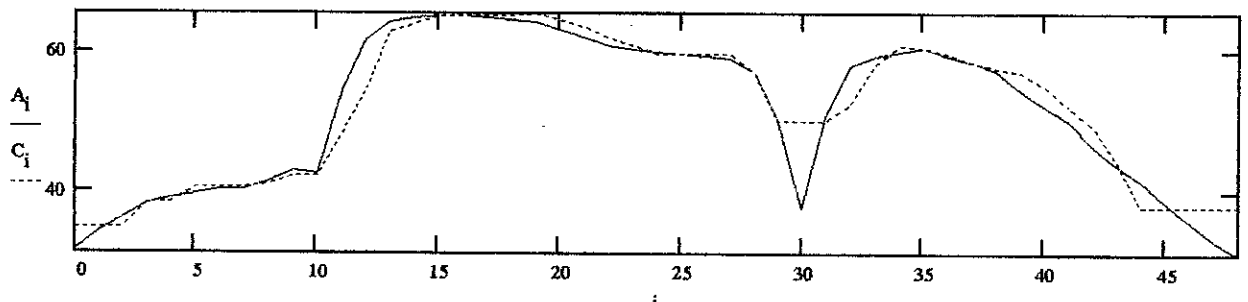
$$x_0 := 0 \quad n := 1..I - 6$$

$$x_n := x_{n-1} + f(d_{n, x_{n-1}}, d_{n+1, x_{n-1}}, d_{n-1, x_{n-1}+1}, d_{n, x_{n-1}+1}, d_{n-1, x_{n-1}+2})$$

$$j := I - 5..I - 1$$

$$x_j := J - 1$$

$$C_i := B(x_i)$$

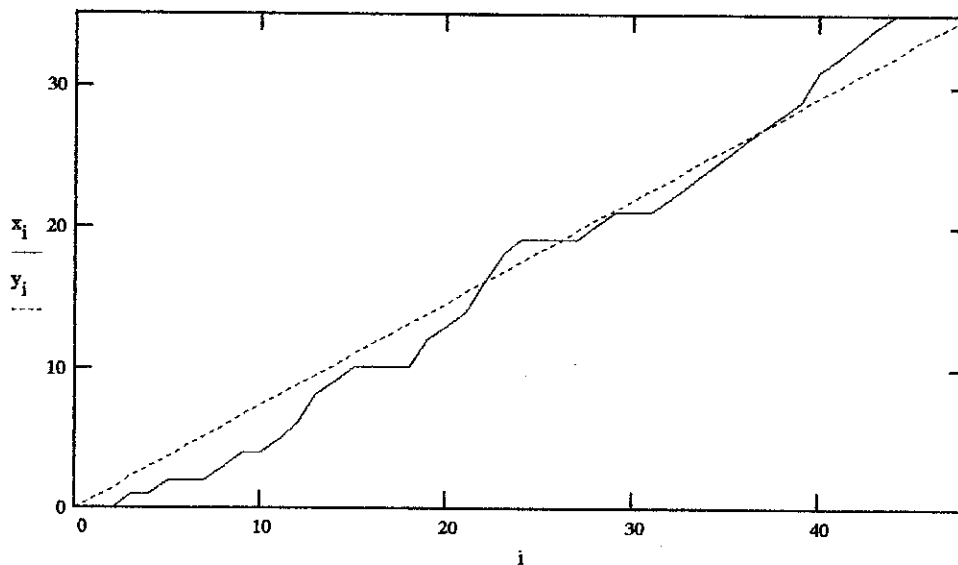


$$E2 := \frac{1}{I} \sum_i (A_i - C_i)^2$$

$$E2 = 8.60524$$

Como se esperaba, la técnica DTW permite un error promedio mucho menor que la compresión lineal

$$y_i := i \cdot \frac{J-1}{I-1}$$



Función de alineación óptima para estas dos funciones de energía de la palabra /CERO/