

1
29j



UNIVERSIDAD NACIONAL AUTÓNOMA
DE MÉXICO

FACULTAD DE INGENIERÍA

RECONOCIMIENTO AUTOMÁTICO DE VOZ

T E S I S
QUE PARA OBTENER EL TÍTULO DE
INGENIERO EN TELECOMUNICACIONES
P R E S E N T A:

ALEXANDRE FREDÉRIC BOUCHET LOPEZ

DIRECTOR DE TESIS M. C. ABEL HERRERA C.



CIDAD UNIVERSITARIA

TESIS CON
FALLA DE ORIGEN
1997

TESIS CON
FALLA DE ORIGEN



Universidad Nacional
Autónoma de México

Dirección General de Bibliotecas de la UNAM

Biblioteca Central



UNAM – Dirección General de Bibliotecas
Tesis Digitales
Restricciones de uso

DERECHOS RESERVADOS ©
PROHIBIDA SU REPRODUCCIÓN TOTAL O PARCIAL

Todo el material contenido en esta tesis esta protegido por la Ley Federal del Derecho de Autor (LFDA) de los Estados Unidos Mexicanos (México).

El uso de imágenes, fragmentos de videos, y demás material que sea objeto de protección de los derechos de autor, será exclusivamente para fines educativos e informativos y deberá citar la fuente donde la obtuvo mencionando el autor o autores. Cualquier uso distinto como el lucro, reproducción, edición o modificación, será perseguido y sancionado por el respectivo titular de los Derechos de Autor.

**Dedico esta tesis a mis padres, quienes me han dado
mas de lo que he querido creer durante muchos años.**

Agradecimientos

Agradezco a Sonia por tantas cosas...

Al M. en Ingeniería Abel Herrera Camacho por haberme iniciado al fascinante mundo de la investigación.

A mi familia por haber soportado días de humor variable y por su apoyo constante a lo largo de la elaboración de este trabajo escrito.

A la familia Cheng por su cariño y apoyo constante.

A Ricardo Ibarra con quien descubrí tantas veces el hilo negro.

A Valerie, Guanna, El Chavo, Federico ("porque sino se enoja") y a todos los que me apoyaron de alguna u otra forma en la elaboración de esta tesis.

A todos los compañeros de la Facultad, a la familia Telleria, y en especial a Javier, Fabricio y Alfredo por su amistad y su apoyo constante, discreto e incondicional.

Índice

INTRODUCCIÓN	1
CAPÍTULO I. ELEMENTOS DE ACÚSTICA, DIGITALIZACIÓN Y CARACTERÍSTICAS DE LA VOZ	3
I.1. Elementos de acústica	3
I.2. Elementos de digitalización	5
I.3. Elementos de acústica vocal y características de la voz	7
I.3.1. El modelo fuente-filtro	7
I.3.2. Sonidos sonoros	9
I.3.3. Sonidos sordos	10
I.3.4. Energía de una señal de voz	11
I.3.5. Tasa de cruce por cero de una señal de voz	14
CAPÍTULO II. EL MODELO LPC	16
II.1. La hipótesis autoregresiva	16
II.2. Obtención de los coeficientes LPC	18
II.3. El método recursivo de Ljung y Quenben	21
II.4. La distancia de Itakura	22
II.5. Importancia del presentador	22
CAPÍTULO III. CUANTIFICACIÓN DINÁMICA EN EL TIEMPO	24
III.1. Palabras acústicas	25
III.1.1. Definición de las palabras en la práctica	25
III.1.2. Definición de las palabras acústicas por procedimiento LPC de las palabras	27
III.1.3. Método de cuantificación adaptativa de las palabras acústicas	30

III.2. Palabras conectadas	33
III.2.1. Concatenación de patrones de referencia	34
III.2.2. Algoritmo de una sola etapa para reconocimiento de palabras conectadas	34
CAPITULO IV. DESCRIPCIÓN DE LOS EXPERIMENTOS Y RESULTADOS	37
IV.1. Generalidades acerca de los experimentos	37
IV.2. Experimento I. Verificación experimental de la validez del modelo fuente-filtro LPC	37
IV.3. Experimento II. Reconocimiento de palabras aisladas	42
IV.4. Experimento III. Reconocimiento de palabras conectadas	47
CAPITULO V. ANÁLISIS Y DISCUSIÓN DE LOS RESULTADOS	48
V.1. Experimento I	48
V.2. Experimento II	49
V.3. Experimento III	50
CONCLUSIONES	53
BIBLIOGRAFÍA	55

Introducción

La meta de la investigación en el reconocimiento automático de voz es desarrollar técnicas y sistemas que permitan que la voz sea aceptada por las computadoras como señal de entrada. Este problema ha sido estudiado ampliamente desde años atrás, y sin embargo existe hoy, a nivel mundial, un número sumamente limitado de sistemas comerciales que resuelven este problema de forma satisfactoria. En todos los casos, estos sistemas funcionan bajo una serie de restricciones que permiten limitar la dificultad del problema. La experiencia de años de investigación en esta área ha logrado definir los siguientes factores determinantes en la dificultad que involucra este reconocimiento:

- Grado de conexión de las palabras
- Tamaño del vocabulario de trabajo
- Restricciones impuestas por la tarea específica y el lenguaje
- Dependencia e independencia del parlante
- Grado de ambigüedad acústica
- Ruido presente en el medio ambiente

El grado de conexión de las palabras determina si se desea reconocer palabras aisladas o cadenas de palabras pronunciadas sin pausas intermedias. En este último caso, que llamaremos palabras conectadas, la dificultad es mayor ya que no se tiene ninguna consecuencia previa acerca de la terminación de cada palabra dentro de la fraseación. Además, en palabras conectadas, una ambigüedad de semejanza de pronunciación son palabras que en el caso de palabras aisladas, no ocurren. Por ejemplo, en un vocabulario de trabajo, distinguirse la pronunciación en el siguiente ambiente: "ca que machucan las perforaciones de ventilación entre palabras" por palabras que suenan muy similares como "ca que machucan las perforaciones de ventilación" requiere, en el gobierno de palabras, una técnica o una computadora por la tarea específica del problema. En otros casos, se requiere de técnicas especiales para la dificultad del problema, como el uso de un diccionario, o el uso de palabras que son aisladas en

cualquier instante. Los sistemas dependientes de un solo parlante funcionan en general mejor que los sistemas independientes de la identidad del parlante. Esto se debe a que las características espectrales de la voz cambian de forma muy marcada de parlante a parlante. El nivel de ruido puede afectar gravemente a sistemas de reconocimiento ya que sus características se mezclan con las características de interés, distorsionando así la información útil para el reconocimiento.

El ajuste dinámico del tiempo es una técnica de comparación no lineal que ha demostrado ser indispensable cuando las palabras son representadas como secuencias de parámetros extraídos de la misma. En el caso de reconocimiento de voz, la representación de las señales por secuencias de coeficientes LPC ha comprobado ser una valiosa herramienta de compresión y extracción de información.

Los objetivos de este trabajo fueron los siguientes:

- Realizar una investigación exhaustiva y un estudio completo del modelo digital fuente-filtro de la voz basado en coeficientes LPC.
- Realizar una investigación exhaustiva y un estudio completo de las diferentes variantes de ajuste dinámico del tiempo en sus aplicaciones a palabras aisladas y palabras conectadas.
- Desarrollar para ambos casos un algoritmo adaptativo y determinar experimentalmente sus parámetros óptimos.

El primer capítulo del trabajo describe brevemente en forma elemental los conceptos básicos comprendidos en el desarrollo del resto del trabajo. Los capítulos II a III presentan en forma sintética los resultados de algunas investigaciones realizadas, mostrando que el cuarto capítulo describe la metodología empleada en los experimentos. El quinto capítulo describe una metodología alternativa de los resultados de los experimentos, en forma sencilla y directa de interpretación.

Capítulo I

Elementos de Acústica, Digitalización y Características de la Voz

1.1. Elementos de Acústica

El sonido proviene de alguna perturbación en un medio elástico, generalmente el aire, que causa una alteración de presión y un desplazamiento de sus partículas. La sensación auditiva se genera en un intervalo de frecuencias de 20 Hz a 20 KHz aproximadamente. Mediciones efectuadas indican que las máximas variaciones toleradas por el oído humano son del orden de 280 dinas/cm², mientras que variaciones menores de $3 \cdot 10^{-8}$ dinas/cm² no producen sensación auditiva. La propagación del sonido es esencialmente longitudinal, es decir que la dirección del movimiento de las partículas es la misma que la de la onda. La velocidad de propagación del sonido en el aire es de alrededor de 340 m/s, valor que varía según la presión atmosférica y la temperatura ambiente.

La intensidad de una onda sonora se define como la potencia media transportada por unidad de superficie y esta equivale a

$$I = \frac{P}{S} \quad (1.1)$$

donde P es la amplitud de variación de presión y v la velocidad de propagación. En general se emplea una escala logarítmica para expresar la intensidad a través del nivel de intensidad β , con

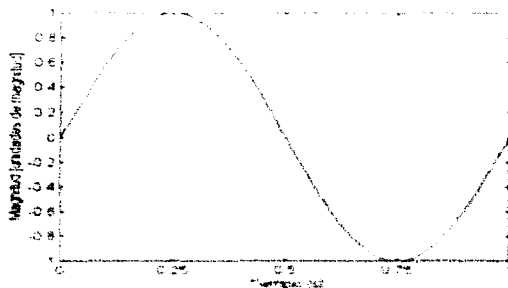
$$\beta = 10 \log \frac{I}{I_0} \quad (1.2)$$

donde I_0 es el nivel de referencia empleado, que en general es de 10^{-16} a 10^{-12} W/cm², el rango del umbral mínimo audible. La siguiente tabla indica algunos niveles de intensidad representativos calculados con una referencia $I_0 = 10^{-12}$.

Decibeles	Fuente Acustica
200	Transbordador espacial despegando
190	
180	
170	Avion de propulsion a chorro con post-combustion
160	
150	Avion de helicos
140	
130	Concierto amplificado de Rock
120	Orquesta de 75 instrumentos (pico maximo)
110	Piano (pico maximo)
100	Automovil en area residencial
90	Voz gritando
80	
70	Voz de conversacion
60	
50	
40	
30	Susurros
20	Furgoneta de granada en granada silenciosamente
10	
0	Nivel de referencia (por la definicion de decibeles)

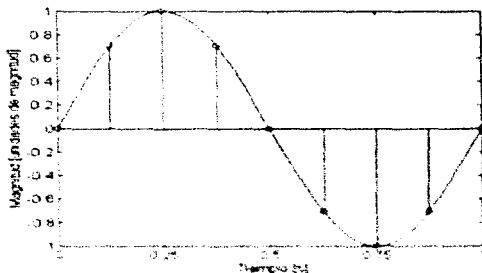
1.2. Elementos de Digitalización

El sonido, como muchas señales naturales, es una señal continua en el tiempo y en amplitud. Esto quiere decir que la señal puede tomar cualquier valor dentro de un intervalo continuo (que contiene una infinidad de valores posibles) y que puede medirse para cualquier tiempo dentro de un intervalo de tiempo continuo (que contiene una infinidad de valores posibles), como lo ejemplifica la siguiente figura.

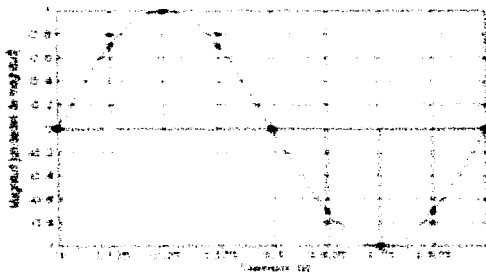


Para poder procesar esta señal en una computadora digital, es necesario limitar el número de valores que pueden tomar tanto la amplitud como el tiempo por que solo existe número finito de bits para representar estos números finitos de valores. Se requiere entonces el proceso de digitalización que consiste de dos partes: el muestreo en el tiempo y el de cuantización de la amplitud. Se está procesando en palabras una representación llamada PCM que puede verse ilustrada que es la representación digital más simple y fácil de entender. Existen más sofisticadas de estas representaciones, más sofisticadas en su análisis, pero de que más detalladamente se presentará en que se profundiza en el curso de estos dos semestres de voz.

El primer paso consiste en tomar una muestra del valor de la señal cada cierto tiempo, como lo ejemplifica la siguiente figura.



El segundo paso consiste en aproximar el valor de la amplitud de la señal en cada muestra con el valor más cercano o inmediatamente inferior de un conjunto finito de valores. Este proceso se conoce como **corte de cuantización**.



Es común que ocurran ambos pasos de forma simultánea en sistemas reales. En este ejemplo se tomó una muestra cada 0.125s y el número de niveles de cuantización es 11. El número de valores posibles, niveles de cuantización, queda determinado por el número de bits empleados para cada muestra mientras que el número de muestras tomadas por unidad de tiempo depende de la frecuencia de muestreo. La relación entre el número de niveles existentes (N) y el número de bits empleados para representarlos (n) es

$$N = 2^n \quad (1.4)$$

y la frecuencia de muestreo (f_m) debe cumplir con el teorema de Nyquist

$$f_m \geq 2f_{max} \quad (1.5)$$

donde f_{max} es la frecuencia máxima que contiene la señal por muestrear.

1.3. Elementos de Acústica Vocal y Características Importantes de la Voz

1.3.1. Elementos fuente-filtro

El tracto vocal es un tubo resonante de un extremo por la laringe y en el otro por los labios. Otro tubo, el tracto nasal, puede ser considerado si deseamos, como una rama resonante al ramamiento del tubo. A lo largo del tracto pueden darse fenómenos como el tracto vocal puntual por la laringe. En este sistema el sonido se genera de tres formas diferentes:

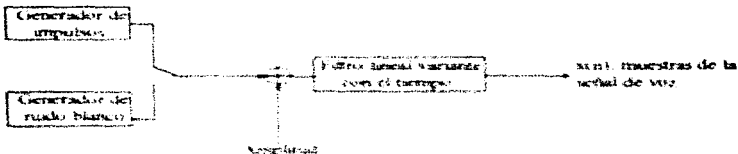
Los sonidos vocales se producen principalmente en los labios, en la cavidad nasofaríngea de la boca y en la cavidad oral en el flujo de aire en vibración por la vibración de las cuerdas vocales y el desplazamiento de la lengua y paladar.

-Los sonidos fricativos se producen al formar una constricción en algún lugar del tracto, creando así turbulencias que producen una fuente de sonido de banda ancha que excita el tracto vocal. Los sonidos fricativos pueden producirse con o sin fonación.

-Los susurros se generan en la laringe al producirse una turbulencia por el paso del aire por una pequeña apertura triangular entre los cartílagos epiglótico y aritenoides.

Todas estas fuentes crean una excitación de banda relativamente ancha del tracto vocal que a su vez se comporta como un filtro lineal, variante en el tiempo, que impone sus propiedades de transmisión sobre el espectro de la fuente. Se ha observado que la variación de este filtro es lenta y que por lo tanto se puede considerar invariante en periodos de alrededor de 10 ms.

Porque la fuente de sonido y la forma del tracto son relativamente independientes, es razonable modelarlas en forma separada, como en la siguiente figura

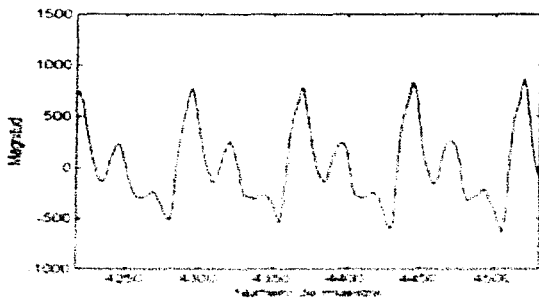


En este modelo digital, las muestras de la señal de voz son la salida de un filtro digital variante en el tiempo que procesa los desplazamientos del tracto vocal a través de un filtro de variación de tiempo de los generadores de ruido blanco. El ruido blanco que se introduce al sistema proviene de una fuente variable que a su vez es filtrada por el filtro. Este modelo debe considerarse como una aproximación de un modelo más complejo que debería considerar los efectos de los cambios de longitud del tracto vocal y los desplazamientos de los segmentos del tracto vocal. Este modelo podría ser mejorado por el uso de un filtro de variación de tiempo que procesa los desplazamientos del tracto vocal a su vez.

realidad representa las propiedades de radiación, del tracto vocal y de los pulsos de fonación.

13.2. Sonidos sonoros

La característica principal de los sonidos sonoros es que son cuasiperiódicos como se ve claramente en la siguiente figura:



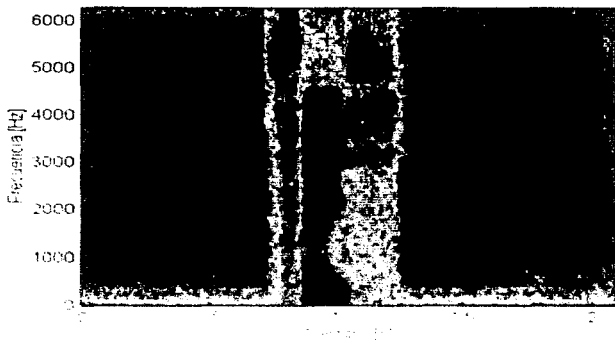
La vibración de las cuerdas vocales que es el punto del que parten una excitación cuasiperiódica del tracto vocal que modifica el espectro de dicha excitación según su respuesta en frecuencia. Diferentes sonidos cuasiperiódicos pueden tener la misma frecuencia fundamental y sus armónicos ser distribuidos de diferente forma. Una característica de estos sonidos es que su potencia varía entre armónicos que son armónicos simples. Las formantes determinan el timbre del sonido o cómo se percibe sobre su intensidad dependiente de la respuesta en frecuencia del tracto vocal que corresponde a sus matices. Entre matices corresponden a una excitación cuasiperiódica del sonido. Anexo 510p. 4

continuación presentamos el espectrograma (aproximación gráfica del contenido espectral instantáneo de una señal) de las vocales /a/, /e/, /i/ y /u/:



Fig. 1. Señal de las vocales

Una vez obtenida la información de la señal de las vocales, se puede observar que el tiempo que se tarda en producir la /a/ de las vocales es el mayor y el menor tiempo que se tarda en producir la /i/. Por otro lado, cuando se producen las vocales /e/ y /i/ se produce un mayor número de armónicos que cuando se producen las vocales /a/ y /u/. La representación de los espectrogramas de las palabras "arroz" y "café" se muestra en las figuras 2 y 3.



1.3.4. Energía de una señal discreta

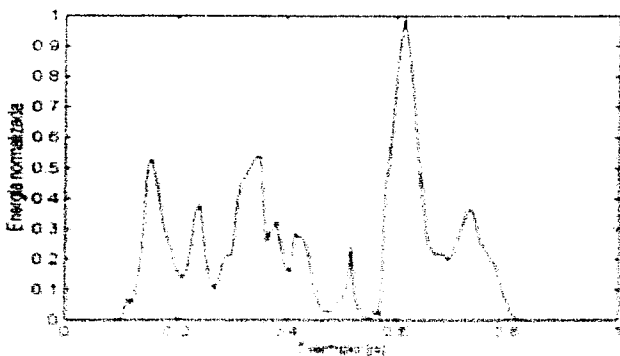
Una de las propiedades más importantes de una señal es su energía. En el caso de una señal discreta esta se define como la suma de los cuadrados de sus valores.

$$E = \sum_{n=-\infty}^{\infty} |x[n]|^2 \quad (1.3.4)$$

Para señales con datos reales esta fórmula se simplifica sustituyendo en el numerador el cuadrado absoluto de los tiempos de muestreo.

$$E = \sum_{n=-\infty}^{\infty} x[n]^2 \quad (1.3.5)$$

donde $w(m)$ es una ventana que selecciona un segmento de $x(n)$ y N es el número de muestras de la ventana. Si se desea dar el mismo peso a todo el intervalo analizado, se emplea una ventana rectangular, mientras que si se desea dar mayor peso a la sección central del intervalo se puede emplear una ventana de pesos variables, como una ventana de Hamming. La selección de la longitud de la ventana es muy importante. Si se selecciona una ventana menor al periodo (en el caso de voz sonora), $E(m)$ fluctuará muy rápido dependiendo de los detalles precisos de la forma de onda, en cambio si se selecciona una ventana de longitud mayor a varios periodos, $E(m)$ tendrá muy poca variación y no reflejará las propiedades variantes de la señal de voz. Una longitud de 10-20 ms ofrece buenos resultados para señales de voz. La siguiente figura muestra el perfil de energía normalizada para la palabra "normalizada".

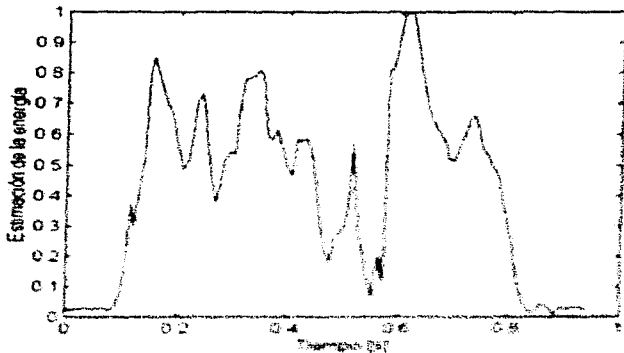


Una dificultad que puede verse al analizar la energía de una voz es la de que los estrechos bandos de energía a menudo están en series parciales entre sí, como se ilustra en la siguiente figura.

Se puede evitar este problema de una simple forma, usando como estimación de la energía la función

$$E(n) = \frac{1}{M} \sum_{m=0}^{M-1} |w(m)x(n-m)|. \quad (1.5)$$

La siguiente figure muestra esta estimación normalizada de la energía de la misma palabra /normalizada/.

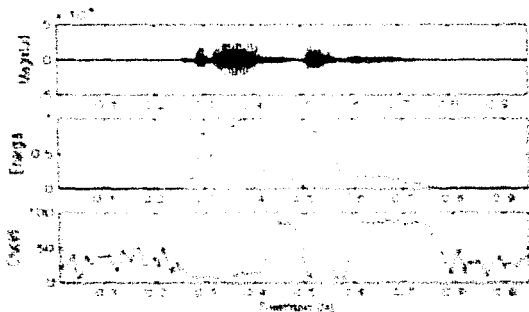


1.3.5. Tasa de cruces por cero de una señal de voz

Otra característica interesante de una señal de voz es su tasa de cruces por cero. En una señal digital, un cruce por cero ocurre entre dos instantes de muestreo n y $n-1$ si

$$\text{signo}\{x(n)\} \neq \text{signo}\{x(n-1)\} \quad (1.6)$$

Esta característica permite por ejemplo estimar las posiciones de sonidos sonoros y sonidos sordos. En efecto, la energía de los sonidos sonoros está concentrada debajo de 3 kHz, mientras que la energía de los sonidos sordos lo está arriba de 3 kHz, y por lo tanto la tasa de cruces por cero sería más alta para sonidos sordos que para sonidos sonoros. Además, los sonidos sonoros generalmente tienen mayor energía que los sordos, lo que ayuda también en la localización de estos tipos de sonidos. La siguiente figura muestra dichas diferencias de tasas y de energía para una grabación de la palabra "crucero".



Los numerosos cruces por cero observables antes y después de la palabra se deben al alto nivel de ruido en la grabación.

Para que la estimación de la tasa de cruces por cero sea correcta, es necesario que la señal no contenga componentes de frecuencias bajas como ruido de 60Hz o niveles de DC, ya que su efecto es desastroso en esta medición. Realizando una medición de cruces por la media de la ventana en cuestión se puede solucionar fácilmente este inconveniente. En dicho caso, es aconsejable restar primero la media del contenido de la ventana de análisis a las muestras de dicha ventana y posteriormente realizar un análisis de tasa de cruces por cero convencional.

Capítulo II

El Modelo LPC

II.1. La hipótesis auto-regresiva

Los métodos más exitosos hasta la fecha para realizar análisis de la voz son los que se basan en el concepto de predicción lineal, porque son precisos y requieren de poco cómputo.

La idea básica detrás de la codificación lineal predictiva (LPC por Linear Predictive Coding) es que una muestra de voz puede ser aproximada por una combinación lineal de las p muestras anteriores, donde p es el orden del modelo. Minimizando el error cuadrático medio entre las muestras reales y las predichas linealmente se pueden determinar los coeficientes del predictor, es decir los pesos de la combinación lineal. El uso del análisis lineal predictivo fue inspirado por el modelo digital de la voz mencionado.

Suponiendo que las muestras de una señal de voz son producidas por el modelo de la siguiente figura, un modelo fuente filtro de la voz:



donde el sistema tiene la siguiente función de transferencia.

$$H(z) = \frac{A}{1 - \sum_{k=1}^L a_k z^{-k}} \quad (2.1)$$

Para sonidos sonoros el sistema es excitado por un tren de impulsos cuasiperiódicos, y para sonidos sordos es excitado por ruido blanco. El análisis por predicción lineal se basa en la observación que para un tal sistema, las muestras $x(n)$ están relacionadas con la excitación $\delta(n)$ por la siguiente ecuación en diferencias.

$$x(n) = \sum_{k=1}^L a_k x(n-k) + \delta(n) \quad (2.2)$$

Supongamos que procesamos esta señal con un predictor lineal

$$\hat{x}(n) = \sum_{k=1}^L \alpha_k x(n-k) \quad (2.3)$$

Entonces el error de predicción se define como

$$e(n) = x(n) - \hat{x}(n) = x(n) - \sum_{k=1}^L \alpha_k x(n-k) \quad (2.4)$$

Comparando (2.2) y (2.4) vemos que si $\alpha_k = a_k$, y si la señal realmente obedece al modelo de (2.2), entonces $e(n) = \delta(n)$. Por lo tanto, entre los modelos de excitación de la voz sonora, el error de predicción debería ser más pequeño si los coeficientes del predictor α_k son iguales a los parámetros a_k de la ecuación de transferencia del trazo vocal. Para sonidos sordos, el error será pequeño a lo largo de la señal. Por lo tanto el problema del predictor

$$\text{Problema: } \sum_{k=1}^L \alpha_k x^2(n-k) \quad (2.5)$$

es una buena aproximación del denominador de la función de transferencia del tracto vocal.

11.2 Obtención de los coeficientes LPC

Una forma de obtener los coeficientes del predictor se basa en la minimización del error cuadrático medio de predicción dentro de un periodo corto; se buscan los valores α_k que minimicen

$$E = \sum_{n=1}^N [x_1(n) - \hat{x}(n)]^2 \quad (2.6)$$

$$E = \sum_{n=1}^N \left[x_1(n) - \sum_{k=1}^p \alpha_k x_1(n-k) \right]^2 \quad (2.7)$$

donde $x_1(n)$ es un segmento de voz seleccionado en la vecindad de la muestra 1,

$$x_1(n) = x(n-1) \quad (2.8)$$

y al resolver $\partial E / \partial \alpha_k = 0$ (queremos satisfacer las condiciones para que E sea mínimo) para $k=1, 2, \dots, p$ obtenemos las ecuaciones

$$\sum_{n=1}^N x_1(n-k) [x_1(n) - \sum_{l=1}^p \alpha_l x_1(n-l)] = \sum_{n=1}^N \alpha_l x_1(n-k) x_1(n-l) = \delta_{kl} \sum_{n=1}^N x_1^2(n) \quad \text{para } 1 \leq k \leq p \quad (2.9)$$

No definamos

$$\varphi_i(i,k) = \sum_{n=0}^N x_i(n-i)x_i(n-k) \quad (2.10)$$

entonces podemos reescribir (2.9) como

$$\sum_{k=0}^p \alpha_k \varphi_i(i,k) = \varphi_i(i,0), \quad i=1,2,\dots,p. \quad (2.11)$$

Este sistema de p ecuaciones con p incógnitas puede ser resuelto calculando $\varphi_i(i,k)$ para $1 \leq i \leq p$ y $1 \leq k \leq p$. Sustituyendo variables reescribimos (2.10) como

$$\begin{aligned} \varphi_i(i,k) &= \sum_{n=0}^N x_i(n)x_i(n+i-k) \\ \varphi_i(i,k) &= \sum_{n=0}^{N-k} x_i(n)x_i(n+k-i) \end{aligned} \quad (2.12)$$

de donde claramente vemos que $\varphi_i(i,k) = \varphi_i(i,k-i)$. También observamos de (2.12) que valores de $x_i(n)$ son requeridos fuera del intervalo $0 \leq n \leq N-1$. Para no proveer los valores fuera de este intervalo es necesario emplear alguna ventana finita de duración finita para reducir los efectos de límite que surgen al tratar de producir las p primeras muestras basándose en muestras cuyo valor es nulo por estar fuera del intervalo $0 \leq n \leq N-1$. Existe para estos casos límite un error de predicción muy grande - es necesario introducir entonces una ventana que reduce sustancialmente la señal hasta cero en los extremos del intervalo $0 \leq n \leq N-1$. Obtenemos entonces:

$$x_i(n) = \begin{cases} x_i(n) & 0 \leq n \leq N-1 \\ 0 & \text{de otra forma} \end{cases} \quad (2.13)$$

Utilizando esta definición de $x_i(n)$, $2 \leq i \leq p$, se puede escribir:

$$\varphi_i(i,k) = \sum_{n=0}^{N-(i+k)} x_i(n)x_i(n+i-k)$$

$$\varphi_i(i,k) = \sum_{n=0}^{N-(k-i)} x_i(n)x_i(n+k-i)$$

$$\varphi_i(i,k) = r_i(i-k) = r_i(k-i) \quad (2.14)$$

En este caso (2.11) se transforma en

$$\sum_{k=0}^p \alpha_k r_i(j-k) = r_i(j), \quad j=1,2,\dots,p. \quad (2.15)$$

Tomando en cuenta que la definición de la autocorrelación en tiempo corto es

$$R_i(m) = \frac{1}{N} \sum_{n=0}^{N-m} x_i(n)w(n)x_i(n-m)w(n+m) \quad (2.16)$$

queda claro que $r_i(n) \approx NR_i(n)$. Esto implica que para obtener los coeficientes LPC es necesario calcular primero los coeficientes de autocorrelación en tiempo corto de la ventana de interés y luego resolver un sistema de ecuaciones (2.15) o (2.16) con p ecuaciones y p incógnitas que puede ser expresado de la siguiente forma

$$\Phi \mathbf{a} = \mathbf{r} \quad (2.17)$$

Este sistema puede ser resuelto considerando cualquiera de los métodos propuestos para resolver sistemas de ecuaciones lineales. Sin embargo, es importante señalar que al resolver $\Phi \mathbf{a} = \mathbf{r}$ en este caso una matriz simétrica y con autocorrelación como elementos en la parte diagonal. Φ es una matriz de Toeplitz y el vector \mathbf{r} puede considerarse un vector de tipo triangular hacia arriba. El método más adecuado para resolver este sistema de ecuaciones es el método de Levinson, el cual puede ser expresado de la siguiente forma

II.3. Método recursivo de Levinson-Durbin

De (2.15) vemos que el sistema por resolver es el siguiente:

$$\begin{bmatrix} r(0) & r(1) & r(2) & \dots & r(p-1) \\ r(1) & r(0) & r(1) & \dots & r(p-2) \\ r(2) & r(1) & r(0) & \dots & r(p-3) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r(p-1) & r(p-2) & r(p-3) & \dots & r(0) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \vdots \\ \alpha_p \end{bmatrix} = \begin{bmatrix} r(1) \\ r(2) \\ r(3) \\ \vdots \\ r(p) \end{bmatrix} \quad (2.18)$$

Los pasos del método recursivo de Levinson-Durbin son los siguientes:

$$\begin{aligned} E^{(0)} &= r(0) \\ k_p &= \frac{r(p) - \sum_{j=0}^{p-1} \alpha_j^{(p-1)} r(p-1-j)}{E^{(p-1)}} \quad \text{para } 1 \leq p \\ \alpha_j^{(p)} &= k_p \\ \alpha_j^{(p)} &= \alpha_j^{(p-1)} - k_p \alpha_{p-j}^{(p-1)} \quad \text{para } 1 \leq j \leq p-1 \\ E^{(p)} &= (1 - k_p^2) E^{(p-1)} \end{aligned} \quad (2.19)$$

Los coeficientes LPS son los valores α obtenidos después de realizar los pasos descritos para todos los valores válidos de i y j . Estos valores definen totalmente las características del filtro digital que representa al trazo social en el modelo fuente-filtro descrito. Estos parámetros caracterizan un argumento de escala, de alrededor de 10ms. Una señal completa de voz puede por lo tanto ser representada como una sucesión de estos parámetros formando así una matriz de coeficientes LPS.

Una característica interesante del modelado por análisis LPC es que una medida de distancia o de disimilitud, entre dos modelos LPC ha sido desarrollada por Fumitada Itakura. Esta medida ha demostrado ser extremadamente eficiente.

11.4. La distancia de Itakura

Supongamos que tenemos dos señales T y R cuyos coeficientes a_1 y a_0 LPC han sido determinados. V es la matriz de autocorrelaciones de T . La distancia entre los modelos LPC definidos para T y R es la siguiente:

$$D(T,R) = \log \frac{a_1 V a_1^T + a_0^2}{a_1 V a_1^T} \quad (2.20)$$

Esta distancia relaciona al error cuadrático medio de predicción de la señal T por el modelo a_1 con el error cuadrático medio de predicción de la señal T por el modelo a_0 . Si la señal T , cuyo modelo LPC es a_1 , responde al modelo a_0 de la señal R , entonces ambos errores de predicción serán similares, su cociente será cercano a 1 y el logaritmo del cociente será cercano a cero. Por otra parte, si T no es similar a R , el error de predicción de T mediante el modelo de R será mayor al error de predicción de T mediante su propio modelo y el cociente será mayor que 1, dando un logaritmo mayor que cero.

11.5. Importancia del procesamiento

La importancia esencial que aquí que pone el acento de atención al momento, es importante tener una tendencia general de cómo por ciertos condiciones aumentan la frecuencia. Cando que en este caso la importancia de frecuencia es fundamental en otros aspectos que poseen en sus diferentes valores algunos momentos de atención. Estas se refieren a la importancia para que la información que generamos en las distancias entre una señal y otra. Este fenómeno puede efectuarse con un filtro para el caso de un filtro de paso bajo, de paso alto o de paso banda de corte

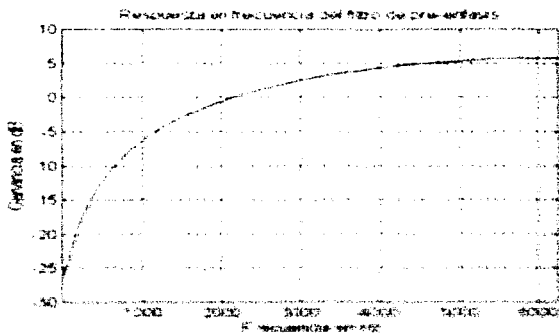
no es crítico ya que se busca principalmente realzar las frecuencias altas. Un tal filtro puede ser implementado digitalmente empleando la ecuación en diferencias

$$y(n) = x(n) - a x(n-1] \quad (2.21)$$

donde $y(n)$ es la salida presente del filtro, $x(n)$ su entrada presente y $x(n-1)$ su entrada anterior. El valor de a debe ser tal que $0.95 \leq a \leq 1$. Otra vez, el valor de a no es muy crítico. Efectuando una transformación discreta Z obtenemos

$$Y(z) = X(z) - a z^{-1} X(z) = (1 - a z^{-1}) X(z) \quad (2.22)$$

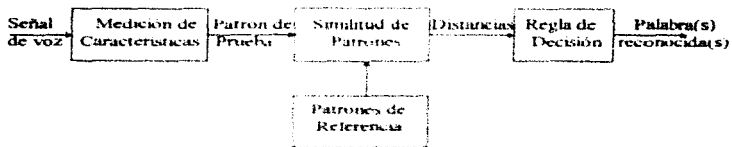
A continuación presentamos la respuesta en frecuencia de dicho filtro cuando $a=0.95$ y la frecuencia de muestreo es de 12.5 kHz



Capítulo III

El ajuste dinámico en el tiempo.

La siguiente figura muestra el modelo de reconocimiento de patrones utilizado en la mayoría de los sistemas de reconocimiento de palabras.



Este modelo consta de tres pasos:

- Medición de características.
- Determinación de similitud con los diferentes patrones de referencia.
- Decisión de la palabra reconocida, basada en alguna regla.

El primer paso es básicamente alguna técnica de extracción de datos. Las características medidas pueden ser de una gran variedad de tipos, como momentos de rectángulos, tasas de cruce por ceros o representaciones como mapas de similitud y similitud de forma en espacios en tiempo cortos o los coeficientes del modelo LPN, es un ejemplo. El segundo paso consiste en la medida de similitud o distancia (según el caso) entre los patrones de prueba que está dependiente de las características analizadas. Existen un número de métodos comunes para esto, tales como representaciones de la señal con coeficientes LPN. Una representación común de la señal de palabras y sus patrones de referencia es un vectorialmente formado de momentos (segunda, tercera y cuarta) calculados en varias

representación por LPC la señal quedará representada por una secuencia de grupos de coeficientes (patrón de coeficientes) donde cada grupo corresponde a un periodo corto (alrededor de 10 ms). El número de grupos de coeficientes que representen a cierta señal dependerá por lo tanto de la longitud de dicha señal. La regla empleada para determinar a qué palabra corresponde la señal analizada requiere de las distancias obtenidas entre la señal de prueba y todos los patrones de referencia. Se escoge la referencia cuya distancia con la señal de prueba es menor, o se escoge la clase (porque varias referencias pueden representar una misma palabra, y por lo tanto pertenecer a la misma clase) cuya presencia sea mayor dentro de los n primeros candidatos. Por ejemplo si $n=3$ y el mejor candidato es cierta palabra a y los segundo y terceros candidatos son otra palabra b , se reconoce la señal analizada como b .

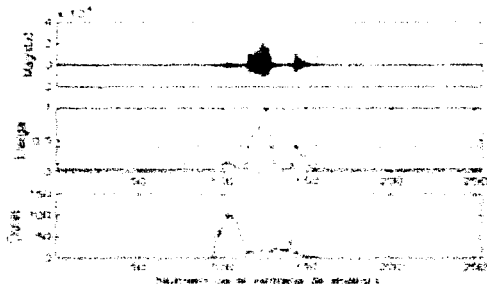
Nos interesaremos primero al caso más simple donde se desea computar una señal de prueba con una sola señal o patrón de referencia a la vez; caso general en el reconocimiento de palabras aisladas.

III.1. Palabras aisladas

III.1.1. Detección de inicio y fin de palabra

En un sistema de reconocimientos de palabras aisladas el primer paso siempre es el de detectar la palabra que se quiere reconocer dentro de el espacio de grabación. Esto permite reducir el tiempo de procesamiento ya que los análisis correspondientes al inicio y al fin de grabación serán descartados. Por otra parte, en la mayoría de los sistemas es necesario realizar dicha detección para poder asegurar que realmente se está escuchando una palabra con otra, y no una mezcla de sonidos o palabras con otra mezcla de sonidos y palabras. El método que presentaremos a continuación demostrando que E. Martinez ha demostrado ser el más preciso y versátil.

Este método de detección de inicio y fin de palabra se basa en dos mediciones en el dominio del tiempo: la energía de la señal y su tasa de cruces por cero. Se efectúa un análisis de estos dos parámetros por ventanas de alrededor de 10ms obteniendo así dos patrones de variación. Es necesario para este algoritmo que los primeros 100ms no contengan voz para que se pueda obtener así una caracterización estadística del ruido de fondo. Usando esta caracterización y el nivel máximo de energía se calculan umbrales de tasa de cruces por cero y de energía. Después se analiza el patrón de energía para determinar el intervalo en el que la señal siempre excede un umbral de energía muy conservador (EE). Se hace entonces la suposición que los puntos de inicio y fin se encuentran fuera de este intervalo y se busca en qué puntos la energía de la señal cae debajo de otro umbral, menor a EE, que llamaremos ET. Estos puntos provisionales de inicio y fin se denominan N1 y N2. El último paso consiste en analizar las 25 ventanas (250ms) anteriores y posteriores a los puntos N1 y N2 respectivamente para determinar cuántas veces el umbral de cruces por cero, TCZ , es superado en ese intervalo. Si el umbral es superado 7 o más veces, los puntos de inicio y/o fin se desplazan al primer (último) punto dentro de estos 250ms donde el umbral fue superado. La siguiente figura ejemplifica el método para la palabra "siete", que representa uno de los casos más difíciles puesto que contiene un silencio dentro de la palabra.



III.1.2. Definición de distancia entre los patrones LPC de dos palabras

Una vez efectuado el análisis de inicio y fin de palabra, se procede a un análisis LPC determinando así para cada palabra de referencia el patrón de coeficientes LPC correspondiente. Cada vez que se desea reconocer alguna palabra, que llamaremos palabra prueba, se debe detectar su inicio y su fin y luego calcular su patrón de coeficientes LPC.

Una vez obtenidos los patrones de características, el siguiente paso es el de determinar la similitud o disimilitud entre patrones de prueba y patrones de referencia. Porque las variaciones de velocidades del habla afectan no solamente la longitud de las palabras, sino también la localización interna de eventos para una misma palabra, la determinación de similitud involucra no solamente el cálculo de distancias, sino también algún tipo de ajuste dinámico de los ejes temporales. Por ajuste dinámico de los ejes temporales entendemos o bien referimos a algún mecanismo que pueda alinear las características de las palabras para que puedan así ser comparadas características correspondientes.

Supongamos que tenemos un patrón de prueba $T(t) = \{T(t_1), T(t_2), \dots, T(t_N)\}$ y un patrón de referencia $R(t) = \{R(t_1), R(t_2), \dots, R(t_N)\}$. El problema de ajuste dinámico del tiempo es de comparar las $R(t_i)$ adecuadas con las $T(t_j)$ adecuadas de tal forma que se minimize la distancia total entre los patrones T y R . Sea C la función que relaciona a los elementos de T con los elementos de R .

$$C = \{c(k_1), c(k_2), \dots, c(k_N)\} \quad (3.1)$$

donde cada elemento de C es un par de asignaciones que empuja los elementos de R y T por compararse

$$c(k) = \{(k), j(k)\} \quad (3.2)$$

El problema es entonces el de encontrar la función C que minimice la distancia entre T y R. Para cada c(k) tenemos una función de distancia o costo d[c(k)]. La distancia entre dos elementos en comparación es, en el caso de análisis LPC, la distancia de Itakura entre los dos modelos LPC. (En este caso los elementos de R y T son vectores de coeficientes LPC). La función de ajuste dinámico del tiempo debe entonces minimizar el costo o la distancia total:

$$D(C) = \sum d[c(k)] \quad (3.3)$$

Existen ciertos requerimientos impuestos por las características de las señales de voz sobre la función c(k). Estos son:

- La función debe ser monótona creciente:

$$\begin{aligned} r(k) &\geq r(k-1) \\ j(k) &\geq j(k-1) \end{aligned} \quad (3.4)$$

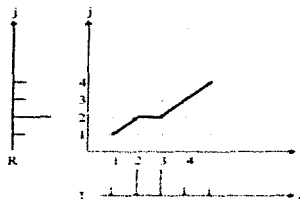
- La función debe cumplir las siguientes condiciones en los límites:

$$\begin{aligned} r(1) &= k(1) = 1 \\ r(K) &= NT \\ j(k) &= NR \end{aligned} \quad (3.5)$$

- La función debe conservar todos los elementos de los patrones introducidos:

$$\begin{aligned} r(k) &= r(k-1) + 1 \\ j(k) &= j(k-1) + 1 \end{aligned} \quad (3.6)$$

Encontrar la función de ajuste temporal es equivalente a encontrar el camino de mínimo costo a través de una rejilla de puntos como lo indica la siguiente figura.



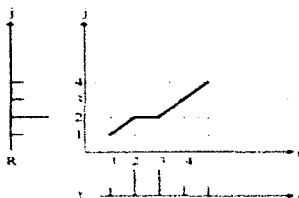
Los elementos 2 y 3 de T se asocian y comparan con el elemento 2 de R ya que representan una misma característica (magnitud grande). En el caso de análisis LPC, la medida de similitud no es la magnitud, sino la distancia de Itakura entre los coeficientes LPC involucrados.

A primera vista puede parecer que DPA (1) tenga que ser evaluada para un número demasiado grande o poco práctico de caminos posibles. Sin embargo, la programación dinámica controla este problema al tomar en cuenta que el mejor camino desde el punto (1,1) hasta cualquier otro punto es independiente de lo que ocurra después de ese punto. Por lo tanto, el costo total del camino que llega hasta el punto (i,j,k) es el costo de ese punto más el costo del camino hasta su precedente óptimo:

$$DPA(i,j,k) = \min_{(i',j',k')} [DPA(i',j',k') + C(i,j,k)] \quad (3.7)$$

Los DPA se calculan por separado para cada punto del espacio que compare con todas las combinaciones más cercanas.

Encontrar la función de ajuste temporal es equivalente a encontrar el camino de mínimo costo a través de una rejilla de puntos como lo indica la siguiente figura.



Los elementos 2 y 3 de I se asocian y comparan con el elemento 2 de R ya que representan una misma característica (magnitud grande). En el caso de análisis LPC, la medida de similitud no es la magnitud, sino la distancia de Itakura entre los coeficientes LPC involucrados.

A primera vista puede parecer que DCTC tenga que ser evaluada para un número demasiado grande e poco práctico de caminos posibles. Sin embargo, la programación dinámica controla este problema al tomar en cuenta que el mejor camino desde el punto (1,1) hasta cualquier otro punto es independiente de lo que ocurre después de ese punto. Por lo tanto, el costo total del camino que llega hasta el punto (k,k) es el costo de ese punto más el costo del camino hacia su predecesor inmediato.

$$DCTC(k,k) = \min_{j \in I} \{ DCTC(k,j) + C(k,j) \}$$

$$(3.7)$$

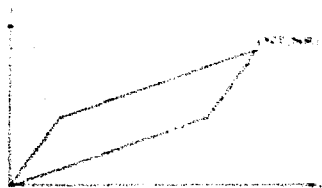
Las ecuaciones anteriores muestran cómo encontrar el camino de mínimo costo a través de una rejilla de puntos.

La programación dinámica funciona por etapas. En su aplicación al ajuste dinámico del tiempo, cada etapa corresponde a una columna del enrejado. Para cada columna, empezando por $i=1$, se calculan las distancias óptimas acumuladas para los puntos de dicha columna. Dado que el cálculo de la siguiente columna está basado en las distancias óptimas para la columna anterior, podemos asegurar que a su vez, las distancias obtenidas para esta nueva columna son las óptimas.

A continuación presentamos un resumen de una serie de restricciones adicionales en la trayectoria de C con las que diferentes investigadores han trabajado. Algunas permiten no solamente mejorar el éxito de reconocimiento, sino también disminuir el número de cálculos necesarios.

III.1.3. Restricciones comúnmente empleadas en la función de ajuste dinámico del tiempo.

- L. Rabiner propone definir la siguiente región de validez o de libertad en forma de paralelogramo cuyos bordes sean rectas con pendientes de ± 0.5 .

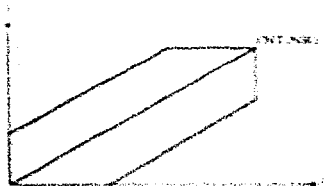


La función de ajuste dinámico debe necesariamente encontrarse dentro de esta región. Esta restricción asegura que se produzcan solamente $g(t)$ y $c(t)$ como un desplazamiento horizontal, pero

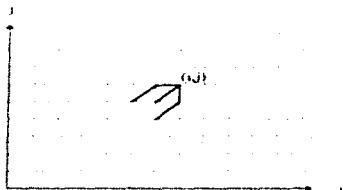
cualquier par de palabras es 2 y que la mínima es $\frac{1}{2}$. Esto implica a su vez que patrones cuya relación de distancias sea menor a $\frac{1}{2}$ o mayor a 2 no podrán ser comparados. Esta restricción es válida porque las variaciones de velocidad del habla son pequeñas, y por lo tanto grabaciones de longitudes considerablemente distintas muy probablemente correspondan a palabras distintas. El número de puntos calculados disminuye ya que únicamente se calculan los puntos internos a la región definida. Los puntos predecesores válidos para algún punto (i, j) son $(i-1, j-1)$, $(i-1, j)$, $(i, j-1)$.

- El Itakura propone una función de ajuste que depende linealmente del índice i . Esto trae como consecuencia que transiciones verticales de la función de ajuste no son permitidas. Por lo tanto los predecesores válidos para un punto (i, j) son $(i-1, j-1)$ y $(i-1, j)$. En este caso se emplea también la región de validez descrita por El Kabiner.

- El estudio más completo realizado sobre restricciones útiles en el ajuste dinámico del tiempo fue realizado por H. Sakoe. Primero define una región de validez, o de libertad, limitada por dos rectas paralelas a la recta de ajuste lineal. La región de validez queda por lo tanto definida por un parámetro ϵ que indica el número de ventanas válidas arriba y abajo de la recta de ajuste lineal.



En su estudio determinó de manera experimental que la mejor restricción que hay que imponer sobre la pendiente local de la función de ajuste es que no sea mayor a 2 ni menor a $\frac{1}{2}$. Esta restricción difiere de las anteriores porque se refiere a la pendiente local y no a la global como en el caso del paralelogramo. La definición de predecesores difiere en este caso. Los caminos válidos que nos lleven hasta un punto (i, j) son:



y de esta forma la distancia mínima acumulada en (i, j) es:

$$g(i, j) = \text{MIN} \begin{cases} g(i-1, j-2) + 2d(i, j-1) + d(i, j) \\ g(i-1, j-1) + 2d(i, j) \\ g(i-2, j-1) + 2d(i-1, j) + d(i, j) \end{cases} \quad (3.8)$$

donde $g(a, b)$ es la distancia mínima acumulada hasta el punto (a, b) y $d(a, b)$ es la distancia local del punto (a, b)

En todos los casos la distancia local sobre el punto P = el patrón R es la distancia acumulada en el punto (N, N, R) que corresponde a la simulación finita especificada en (3.5).

III.2. Palabras conectadas

Sorprendentemente, una de las aplicaciones más importantes de las técnicas para reconocimiento de palabras aisladas ha sido en el área de reconocimiento de palabras conectadas. En el reconocimiento de palabras conectadas, la señal de entrada es una secuencia de palabras pertenecientes a cierto vocabulario definido y el reconocimiento se basa en la comparación del patrón resultante de entrada con patrones de palabras aisladas del mismo vocabulario. La mayor dificultad en este tipo de sistemas reside en un fenómeno conocido como co-articulación. En efecto, en el caso de palabras conectadas las características del principio de cada palabra se ven afectadas por las características del final de la palabra anterior. Esto trae una distorsión de las palabras que al ser comparadas con palabras aisladas de referencia acumularán distancias mayores, disminuyendo así la diferencia de distancia entre una comparación entre dos señales equivalentes, que representan la misma palabra, y dos señales no equivalentes. Ejemplos típicos de reconocimiento de palabras conectadas incluyen el reconocimiento de dígitos conectados donde el vocabulario es el conjunto de los 10 dígitos (0-9) y el reconocimiento de letras conectadas donde el vocabulario es el alfabeto.

El sistema general de reconocimiento de palabras conectadas empleando patrones es muy similar al sistema de reconocimiento de palabras aisladas. El único paso que varía es el de comparación del patrón de entrada con los patrones de referencia, porque en este caso el patrón de entrada contiene varias palabras cuyas localizaciones dentro de la señal no son conocidas. Las técnicas de programación dinámica son especialmente indicadas para este tipo de problemas ya que pueden resolver de forma automática los problemas de detección de límites entre palabras, automáticamente a partir de los datos de tiempos y reconocimiento. Esto descarta la posibilidad de errores causados por una segmentación errónea de la señal de prueba. Existen dos enfoques principales para diseñar algoritmos de comparación: la

concatenación de patrones de referencia para crear patrones de referencia de palabras conectadas y la generalización de las reglas de transición descritas por (3.4), (3.5) y (3.6) para generar un algoritmo de una sola etapa.

III.2.1. Concatenación de patrones de referencia.

Este primer enfoque considera que la cadena de palabras por reconocerse debe ser comparada con las diferentes cadenas posibles formadas por la concatenación de patrones de referencia. La dificultad que esto involucra se debe a la cantidad de comparaciones necesarias para poder realizar una búsqueda exhaustiva entre todas las posibilidades de concatenación. Diferentes autores (El Rabiner, H. Sakoe) encuentran soluciones aproximadas que minimizan, siguiendo ciertos criterios, el número de cálculos requeridos. Dichas optimizaciones se basan tanto en las propiedades del algoritmo dinámico como en la definición de la zona de libertad. Es necesario dividir el problema de optimización en diferentes etapas, por lo que son algoritmos muy complejos.

El problema principal que hemos notado que coexiste con estos métodos es que, si bien si se compara el patrón de prueba con todas las concatenaciones posibles de patrones, no se analizan, explícita o implícitamente, todos los caminos, o funciones de ajuste dinámico que relacionen al patrón de prueba con las cadenas de prueba. Por otra parte, para asegurar su buen funcionamiento, estos métodos requieren de parámetros determinados experimentalmente, y su eficiencia depende en función del número de palabras que contengan las señales de prueba.

III.2.2. Algoritmo de una sola etapa para reconocimiento de palabras conectadas.

Este algoritmo fue originalmente desarrollado por Sakoe, aunque no se popularizó hasta que H. Ney lo reformuló de una manera más sencilla, aunque igualmente. Este algoritmo es una generalización del algoritmo implementado para palabras aisladas. En el caso de palabras aisladas, la función de ajuste dinámico tiene los límites $m = 0$, donde el patrón

de entrada está compuesto por $i=1, \dots, N$ ventanas de análisis y cada patrón de referencia está compuesto por $j=1, \dots, J(k)$ ventanas de análisis. k es el identificador del patrón de referencia. En el caso de reconocimiento de palabras conectadas permitiremos que la función de ajuste tenga tres dimensiones, respetando ciertas características impuestas por la naturaleza de las señales involucradas:

$$W = (w(1), w(2), \dots, w(L)) \quad \text{donde } w(i) = (d(i), j(i), k(i)) \quad (3.9)$$

Tomando en cuenta que las distancias locales ahora pueden tener tres argumentos, $d(i, j, k)$, entonces el problema de reconocimiento puede ser tratado como el siguiente problema de minimización:

$$\text{MIN } \sum d(w(i)) \quad (3.10)$$

Este problema es el de minimizar la distancia global con respecto a todos los caminos posibles. Como ya vimos, este problema puede resolverse mediante programación dinámica. Falta sin embargo definir las reglas de transición descabies e impuestas por las características de las señales de voz.

- Regla de transición en el interior de una referencia:

Si $w(i) = (i, j, k) \neq 1$, entonces:

$$w(i) = (i, j, k), (i, j, k), (i, j, k), (i, j, k) \quad (3.11)$$

Esta regla de transición es idéntica al caso de palabras aisladas.

- Regla de transición entre referencias

si $w(l)=(i,l,k)$, entonces,

$$w(l-1) \in \{(i-1,l,k), (i-1,j(k'),k') \mid k'=1,\dots,K\} \quad (3.12)$$

Esta regla implica que se puede llegar a la referencia k desde la última ventana de análisis de cualquier referencia, incluyendo la misma referencia k .

Mediante programación dinámica se puede determinar fácilmente la distancia mínima acumulada. En este caso nos interesa tanto la distancia acumulada como la trayectoria que sigue la función de ajuste dinámico. En efecto, una vez encontrada esta distancia, se puede determinar la trayectoria que la generó, y determinar así por qué referencias pasó y en qué instantes ocurrió transición de una referencia a otra, resolviendo así simultáneamente los problemas de segmentación y de reconocimiento.

Capítulo IV

Descripción de los experimentos y resultados

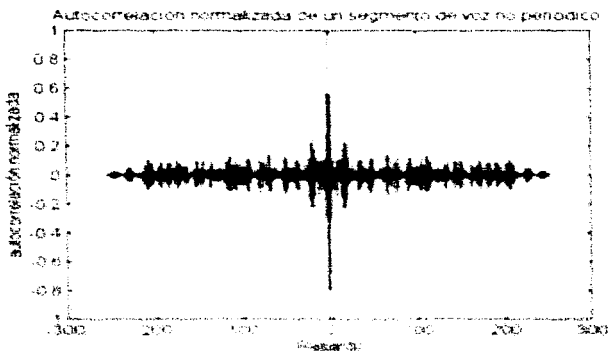
IV.1. Generalidades acerca de los experimentos

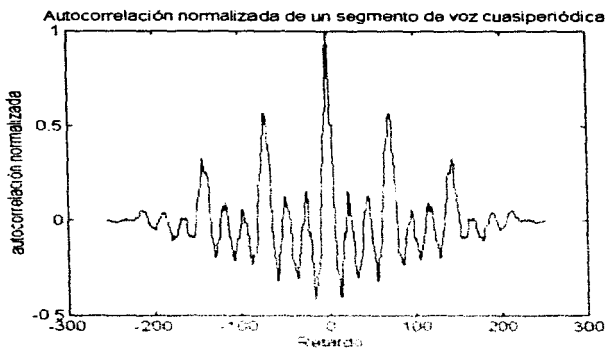
Todos los experimentos que a continuación se presentaran involucraron el uso de algún lenguaje de programación. Se optó por emplear el ambiente de programación MATLAB porque ofrece un gran diversidad de opciones de visualización y flexibilidad en el manejo de todo tipo de variables. La contraparte de estas cualidades es que los tiempos de ejecución fueron mayores a los que hubiésemos observado al emplear el lenguaje de programación C. Se optó por MATLAB tomando en cuenta que el objetivo de este trabajo no es producir un sistema de reconocimiento de palabras que funcione en tiempo real, sino analizar exhaustivamente los métodos de ajuste dinámico del tiempo y la representación por modelado de tipo LPC. Las herramientas gráficas del ambiente mencionado resultaron fundamentales en el estudio de los conceptos analizados. Se empleó la versión 1.0 del compilador para MATLAB que permite compilar rutinas escritas en MATLAB (extensión .m) en rutinas más rápidas llamadas ejecutables de MATLAB (extensión .mex) mediante el compilador ANSI WATCOM versión 1.0.

IV.2. Experimento 1. Verificación experimental de la validez del modelo fuente: filtro LPC

En este experimento se buscó determinar si se analiza de la suficiente manera el contenido como para poder efectuar una reconstrucción exitosa destinada de modular alguna señal de voz. Se escribió un programa que analiza cualquier señal de voz muestreada a 12.800Hz por ventanas de 128 muestras cada 64 muestras. Para cada ventana el programa retorna la energía y calcula los coeficientes LPC del orden que desea el usuario. Por otro lado, el programa determina si existe o no periodicidad en el caso de ser necesario en voz. Para lograr

esta última estimación es necesario emplear ventanas más largas que contengan más de un periodo y el programa analiza por lo tanto ventanas de 256 muestra cada 64 muestras. El análisis de periodicidad esta basado en la autocorrelacion de la señal. Existen metodos complejos para determinar la periodicidad de una señal, pero en este caso se determinó el procedimiento a seguir basandonos en observaciones experimentales ya que el rango de frecuencias fundamentales posibles es pequeño. Primero se limita el analisis al rango de tonos fundamentales que tiene una voz humana convencional (100-300Hz). Dentro del rango de periodos equivalentes se determina si alguna autocorrelacion tiene como magnitud por lo menos el 40% de la magnitud maxima. Si existen varios picos mayores a este umbral se escoge el pico maximo. Si no existen picos mayores a este umbral se determina que el segmento no es periodico. En efecto se observo que en el caso de no periodicidad las autocorrelaciones son menores a este umbral mientras que en el caso de periodicidad la autocorrelacion correspondiente al periodo es mayor a este umbral, como se ve a continuacion

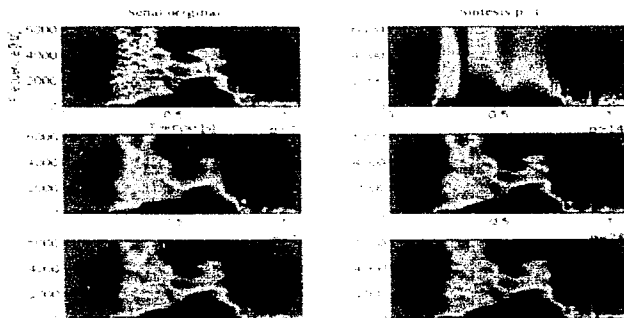




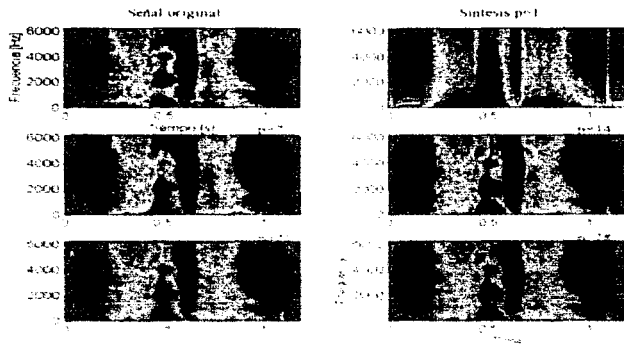
El programa resintetiza después una señal, basándose en los patrones de parámetros obtenidos en el análisis. Para efectos de continuidad en la fase de la frecuencia fundamental se genera primero una señal excitadora completa, lo que facilita preservar dicha continuidad cuando se llega a un fin de ventana y cambios de frecuencia fundamental. Después se procede a filtrar esta señal mediante los coeficientes LPC correspondientes. Se obtiene entonces una señal sintética cuya energía se estimara (se analiza por ventanas de 128 muestras cada 64 muestras para determinar la energía y esta se normaliza a 1). El paso final consiste en multiplicar los segmentos de señal por los valores adecuados para que el perfil final de energía de la señal sintetizada correspondiera al perfil de energías de la señal original.

Mediante técnicas analíticas se observó que señales sintetizadas con modelos de orden mayor a 7 son inteligibles y que señales resintetizadas con modelos de orden mayor a 21 demuestran gran similitud a la original, hasta el punto de ser en ocasiones difíciles de distinguir. La mayor precisión se alcanza cuando se resintetizan (comparativamente a la señal original) las señales sintetizadas. Como se verá en los próximos capítulos, presentamos a continuación los

espectrogramas de una señal original y de las señales resintetizadas con diferentes ordenes de modelado para la palabra /zero/. La distribución de los espectrogramas es la siguiente: de derecha a izquierda y de arriba hacia abajo, señal original, señal sintética de orden 1, señal sintética de orden 7, señal sintética de orden 14, señal sintética de orden 21 y señal sintética de orden 28.



Esto es la misma distribución de espectrogramas para la palabra /six/.



En algunas partes se observan mejor que la palabra sintetizada conserva algo similar a un formante que no existía en la grabación original.

IV.3. Experimento II: Reconocimiento de palabras aisladas

Para realizar este experimento se empleo una base de datos que contiene grabaciones de los diez digitos en ingles (0-9). Diez parlantes repiten cada digito un total de 26 veces cada uno, 10 para entrenamiento y 16 para prueba, generando así una base de 2600 grabaciones. Se conto por lo tanto con 1000 grabaciones de referencia con las que se reconocieron 1599 palabras de prueba (un archivo de prueba de la palabra nine: resultado estar vacío por lo que se decidió descartarlo). Se efectuó primero una detección de inicio y fin de palabra y se determinaron experimentalmente los siguientes umbrales:

$$THU = (\text{maximo de energia}) * 0.4$$

$THL = (\text{media de la energia del ruido de fondo}) + 2 * (\text{desviacion estandar de la energia del ruido de fondo})$

$TLZ = (\text{media de los cruces por cero del ruido de fondo}) + 2 * (\text{desviacion estandar de los cruces por cero del ruido de fondo})$

Despues cada grabacion fue transformada a un patron de coeficientes LTP de orden $p=7$, analizando ventanas de 128 muestras cada 128 muestras. Para cada patron de prueba se determinaron las distancias minimas con respecto a todas las referencias, asignando una distancia infinita a las referencias cuyas energias eran mayores a 4% de la longitud de la prueba o menores a 1% de la longitud de la misma. El siguiente procedimiento fue el propuesto por H. Sakoe para ser el que ha demostrado mejores tasas de reconocimiento:

Se obtuvieron las siguientes estadísticas de reconocimiento (cada renglón corresponde a 160 pruebas y las columnas corresponden al resultado del reconocimiento):

prueba 1:
ajuste lineal (r=0)

	one	two	three	four	five	six	seven	eight	nine	zero
one	111	2	1	4	0	1	1	38	1	
two	1	144	1	6	1	2	2	1	2	
three	4	9	116	0	10	3	4	11	0	
four	10	4	0	130	6	0	1	0	9	
five	10	1	1	2	100	0	5	0	41	0
six	1	8	8	2	4	127	10	6	0	0
seven	1	9	0	6	1	1	121	0	3	18
eight	4	1	2	0	0	7	0	146	0	0
nine	27	3	1	2	8	10	0	0	111	8
zero	3	2	0	2	1	0	18	0	1	136

359 errores de 1599 pruebas. Tasa de reconocimiento: 1 ASA, 99.55%

prueba 2:

Se aplicó a las veintiseis compañías un grupo de veintiseis empleados por 11

Saque un histograma de esta prueba

	one	two	three	four	five	six	seven	eight	nine	zero
one	158	0	0	1	0	0	0	0	0	0
two	0	159	0	0	0	0	0	0	0	0
three	0	0	153	0	0	0	0	0	0	0
four	0	0	0	164	0	0	0	0	0	0
five	0	0	0	0	158	0	0	0	0	0
six	0	0	0	0	0	156	0	0	0	0
seven	0	0	0	0	0	0	156	0	0	0
eight	0	0	0	0	0	0	0	153	0	0
nine	0	0	0	0	0	0	0	0	156	0
zero	0	0	0	0	0	0	0	0	0	159

24 errores de 1596 pruebas. Tasa de reconocimiento: 1 ASA, 99.50%

prueba 3:

r=14 con preenfasis después de la detección de inicio y fin

	one	two	three	four	five	six	seven	eight	nine	zero
one	159	0	0	1	0	0	0	0	0	0
two	0	160	0	0	0	0	0	0	0	0
three	0	0	157	0	0	0	0	0	0	0
four	1	0	0	159	0	0	0	0	0	0
five	1	0	0	0	157	0	0	0	2	0
six	0	0	0	0	0	158	2	0	0	0
seven	0	0	0	0	0	0	160	0	0	0
eight	0	0	0	0	0	0	0	160	0	0
nine	0	0	0	0	0	0	0	0	159	0
zero	0	0	0	0	0	0	0	0	0	160

10 errores TASA=99.33%

prueba 4:

r=14 con preenfasis antes de la detección de inicio y fin

	one	two	three	four	five	six	seven	eight	nine	zero
one	160	0	0	0	0	0	0	0	0	0
two	0	159	0	0	0	0	0	0	0	0
three	0	0	160	0	0	0	0	0	0	0
four	1	0	0	159	0	0	0	0	0	0
five	1	0	0	0	158	0	0	0	2	0
six	0	0	0	0	0	159	0	0	0	0
seven	0	0	0	0	0	0	160	0	0	0
eight	0	0	0	0	0	0	0	160	0	0
nine	0	0	0	0	0	0	0	0	159	0
zero	0	0	0	0	0	0	0	0	0	160

4 errores TASA=99.71%

Una vez demostrada la utilidad del preentasis aplicado antes de la detección de inicio y fin de palabra se realizaron pruebas con diferentes valores de r, el ancho de la región de validez, aplicando dicho preentasis, para determinar experimentalmente el valor óptimo.

prueba 5:

r=18

	one	two	three	four	five	six	seven	eight	nine	zero
one	160	0	0	0	0	0	0	0	0	0
two	0	159	0	0	0	0	0	1	0	0
three	0	0	160	0	0	0	0	0	0	0
four	1	0	0	159	0	0	0	0	0	0
five	1	0	0	0	159	0	0	0	0	0
six	0	0	0	0	0	160	0	0	0	0
seven	0	0	0	0	0	1	159	0	0	0
eight	0	0	0	0	0	0	0	160	0	0
nine	0	0	0	0	0	0	0	0	159	0
zero	0	0	0	0	0	0	0	0	0	160

4 errores: 1 ANA=99.73%

prueba 6:

r=10

	one	two	three	four	five	six	seven	eight	nine	zero
one	160	0	0	0	0	0	0	0	0	0
two	0	159	0	0	0	0	0	1	0	0
three	0	0	160	0	0	0	0	0	0	0
four	1	0	0	159	0	0	0	0	0	0
five	1	0	0	0	159	0	0	0	0	0
six	0	0	0	0	0	160	0	0	0	0
seven	0	0	0	0	0	1	159	0	0	0
eight	0	0	0	0	0	0	0	160	0	0
nine	0	0	0	0	0	0	0	0	159	0
zero	0	0	0	0	0	0	0	0	0	160

4 errores: 1 ANA=99.73%

prueba 7:

r=6

	one	two	three	four	five	six	seven	eight	nine	zero
one	160	0	0	0	0	0	0	0	0	0
two	0	159	0	0	0	0	0	1	0	0
three	0	0	159	0	0	1	0	0	0	0
four	1	0	0	159	0	0	0	0	0	0
five	1	0	0	0	159	0	0	0	0	0
six	0	0	0	0	0	160	0	0	0	0
seven	0	0	0	0	0	1	159	0	0	0
eight	0	0	0	0	0	0	0	160	0	0
nine	0	0	0	0	0	0	0	0	159	0
zero	0	0	0	0	0	0	0	0	0	10

5 errores. TASA=99.69%

prueba 8

r=2

	one	two	three	four	five	six	seven	eight	nine	zero
one	160	0	0	0	0	0	0	0	0	0
two	0	159	0	0	0	0	0	1	0	0
three	0	0	159	0	0	1	0	0	0	0
four	1	0	0	159	0	0	0	0	0	0
five	1	0	0	0	159	0	0	0	0	0
six	0	0	0	0	0	160	0	0	0	0
seven	0	0	0	0	0	0	159	0	0	0
eight	0	0	0	0	0	0	0	160	0	0
nine	0	0	0	0	0	0	0	0	159	0
zero	0	0	0	0	0	0	0	0	0	159

8 errores. TASA=99.1%

prueba 9:

r=1

	one	two	three	four	five	six	seven	eight	nine	zero
one	160	0	0	0	0	0	0	0	0	0
two	0	159	0	0	0	0	0	0	0	0
three	0	0	159	0	0	0	0	0	0	0
four	1	0	0	159	0	0	0	0	0	0
five	1	0	0	0	159	0	0	0	0	0
six	0	0	0	0	0	160	0	0	0	0
seven	0	0	0	0	0	0	159	0	0	0
eight	0	0	0	0	0	0	0	160	0	0
nine	0	0	0	0	0	0	0	0	159	0
zero	0	0	0	0	0	0	0	0	0	159

12 errores. TASA=99.2%

IV.4. Experimento III: Reconocimiento de palabras conectadas.

En este experimento se empleo una base de datos que contiene dos repeticiones de cada dígito en inglés (0-9) y la segunda versión del cero en inglés pronunciada (o/) por catorce parlantes distintos. Estas grabaciones fueron empleadas como patrones de referencia. Se seleccionaron manualmente regiones de ruido de fondo para generar los patrones de referencia para dicho ruido. Esto permite mayor flexibilidad en el caso en que la señal de prueba contenga una región de silencio entre dos dígitos. Por otra parte, la base contiene grabaciones de cadenas de desde dos hasta cinco dígitos y de siete dígitos pronunciadas por los mismos catorce parlantes, que fueron empleadas como señales de prueba. Se empleó el algoritmo descrito por H. Ney por ser el más elegante y el que menor número de cálculos requiere, preservando un tasa de reconocimiento igual o superior a las registradas con los otros algoritmos. Ensayos preliminares con resultados desastrosos demostraron la necesidad de emplear la restricción sobre la pendiente local, como se hizo en el caso de palabras aisladas. A continuación presentamos los resultados para cada longitud de cadena con número de dígitos contenidos:

Longitud de la cadena	1		2		3		4		5	
	correctos	errores	correctos	errores	correctos	errores	correctos	errores	correctos	errores
2 dígitos	136	14	2	17	1	10	1	10	1	11
3 dígitos	113	26	3	17	1	10	1	10	1	11
4 dígitos	113	26	4	15	1	10	1	10	1	11
5 dígitos	112	27	2	15	1	10	1	10	1	11
7 dígitos	111	28	2	15	1	10	1	10	1	11

Expresando estos resultados en porcentajes obtenimos:

Longitud de la cadena	1		2		3		4		5	
	correctos	errores	correctos	errores	correctos	errores	correctos	errores	correctos	errores
2 dígitos	87.5%	8.6%	1.6%	19%	100%	10%	100%	10%	100%	10%
3 dígitos	87.5%	18.6%	1.6%	19%	100%	10%	100%	10%	100%	10%
4 dígitos	87.5%	18.6%	2.6%	15%	100%	10%	100%	10%	100%	10%
5 dígitos	86.8%	19.2%	1.6%	15%	100%	10%	100%	10%	100%	10%
7 dígitos	86.4%	19.6%	1.6%	15%	100%	10%	100%	10%	100%	10%

Capítulo V

Análisis y discusión de los resultados

V.1. Experimento I

En este experimento se pudo demostrar el buen funcionamiento práctico del modelo fuente-filtro de la voz basado en los coeficientes LPC. En algunos espectrogramas se observaron ciertas estructuras en forma de bandas verticales que no existen en la señal original. Después de un análisis del patrón de periodicidad se determinó que dichas estructuras aparecen al existir errores en la detección de periodicidad. En efecto, únicamente ocurre este fenómeno en ausencia de señal de voz y ocurre cuando por algún motivo el ruido aleatorio de fondo presenta algún tipo de periodicidad en un periodo corto. Esto no es grave ya que dichas estructuras tienen una energía muy pequeña y prácticamente no son audibles.

Podríamos cuestionar aquí la validez del modelado del ruido de fondo que se efectuó. En efecto, el modelo fuente-filtro está diseñado para representar una señal de voz y no silencios o ruidos de fondo. Se decidió tratar a toda la señal (silencios y voz) con el mismo modelo ya que en la mayoría de los casos el ruido aparece representado por una señal de entrada aleatoria y el filtro diferenciado por los coeficientes LPC. Tratarse de aproximar el contenido espectral del silencio original. Entendamos esto en forma más clara: el hecho que de todas formas los silencios o ruidos de fondo representados por una energía muy pequeña, se justifica el empleo del mismo modelo.

El efecto de las bandas observadas en algunos espectrogramas, se observaron en modelos similares, fue muy interesante y sorprendente de haber ocurrido. Después de las pruebas se determinó que dichas bandas se deben a que la señal representada es aproximadamente periódica. Esto se debe a errores de precisión en el procesamiento del período. Al representar la señal con una frecuencia fuertemente aproximada (dentado) a la frecuencia exacta de la estructura, ésta estructura está periódica. Estas frecuencias totalizadas y el desarrollo de

frecuencias armónicas a la fundamental de la señal original, determinadas por los coeficientes del filtro. Como consecuencia escuchamos dos sonidos, en vez de uno solo con tono y timbre propio en el caso de la señal original. Para mayor información sobre la organización de la percepción del sonido recomendamos consultar (A. S. Bregman, 1996).

V.2. Experimento II

El primer paso del experimento consistió en realizar reconocimiento con ajuste lineal del tiempo para tener algún dato de comparación. Observamos claramente la necesidad de emplear preentasis en las señales, resultando así la información que nos permite distinguir entre sonidos sonoros. Las pruebas 2-4 respaldan esta afirmación, mientras que las pruebas 3 y 4 demuestran que el preentasis también trae un mejoramiento de la detección de inicio y fin de palabra. Explicamos esto de la siguiente manera:

Recordando que la mayoría de los errores de detección de inicio y fin ocurren en palabras que inician y/o terminan con sonidos sonoros, y que la energía de estos sonidos está concentrada por encima de los 3 kHz, observamos que al aplicarles preentasis su energía aumenta en comparación con la energía del estado de sonido, cuya energía está generalmente repartida por todo el espectro. Esto se debe a que el preentasis es un paso que un filtrado pasa alto. De esta forma, los sonidos sonoros se distinguen más fácilmente del estado de fondo cuando se le aplica preentasis a la señal, y esto obviamente mejora la detección de inicio y fin de palabra.

Las pruebas 4-7 demuestran que el ajuste lineal para la extensión de la señal es útil. En efecto, con este ajuste se observaron los mejores resultados, como es el caso también con $n=14$ y $n=18$, con un número de palabras de cuatro palabras en algunos casos y de tres palabras de palabras y algunas combinadas, a los mejores resultados para estos grupos a mayor cantidad de palabras y por lo tanto tiempos de procesamiento. La comparación de los datos de reconocimiento para palabras de cuatro palabras a los datos a que se está viendo ahora demuestra la posibilidad de aplicar el ajuste lineal en el momento de las pruebas cuando se alargan de las palabras correspondientes a los segmentos.

De estas pruebas podemos concluir que las variaciones temporales en el habla son pequeñas, de alrededor de ± 10 ventanas de análisis en este caso, lo que equivale aproximadamente a 90 milisegundos. Sin embargo, y a pesar de que estas variaciones sean pequeñas, el ajuste dinámico del tiempo es fundamental para un sistema de reconocimiento ya que mejoró en este caso la tasa de reconocimiento desde 77.55% hasta 99.75% en el caso óptimo, para el cual se registraron únicamente 4 errores. Una observación posterior demostro que estos errores se deben a errores en la detección de inicio y fin de palabra, por lo que concluimos que un aumento en el orden del modelo LPC empleado no mejoraría el funcionamiento del sistema.

Determinamos que las características óptimas para un sistema de reconocimiento de palabras aisladas basado en ajuste dinámico del tiempo y empleando una representación de la señal por patrones de coeficientes LPC son las siguientes:

- Un preprocesado debe ser aplicado a todas las señales previas cualquier procesamiento.
- El algoritmo debe ser el propuesto por H. Nakase con 12000 bins y un orden de modelo LPC de 7.

VI. Experimento III

Los resultados obtenidos muestran una alta precisión en el reconocimiento de palabras conectadas que en el reconocimiento de palabras aisladas. Es necesario sin embargo resaltar las diferencias de resultados entre ambos experimentos para concluir que los resultados no deben ser atribuidos al hecho de que las palabras con palabras aisladas se comparan con bases de reconocimiento o de referencias formadas por una única muestra de cada dígito, para cada uno de los patrones, en cambio en los patrones con palabras conectadas la base de referencias utilizada contiene representaciones de cada dígito para cada uno de los patrones. Además, en el caso de patrones conectados se

consideraron las dos versiones del cero en inglés, /zero/ y /oh/, generando así un vocabulario más grande y con mayor ambigüedad ya que /o/ es idéntica al último sonido de /zero/.

Debemos tomar en cuenta estos factores fundamentales que modificarían la precisión de cualquier sistema de reconocimiento de palabras. En efecto, al contar con menos palabras de referencia se tiene una representación mucho menos precisa de las diferentes posibilidades de pronunciación de cada palabra, aumentando así la posibilidad de confusión entre dos palabras. Este fenómeno fue muy claro ya que en su gran mayoría los errores de reconocimiento se debieron a confusión entre las palabras: five y nine y entre one y nine. Otro gran número de errores fue causado por la transformación de la palabra zero en la palabra -o-, seguida de la palabra -o-, resultando así en la inserción de un dígito. Esto se debe a la enorme similitud entre el último fonema de la palabra zero y la palabra -o-. Notamos también que la mayoría de los errores ocurrieron en los archivos de menor tamaño. Examinando con detalle estos archivos nos percatamos que la duración de algunos de sus dígitos es menor a la mitad de la duración de la mayoría de sus correspondientes referencias. Como se empleó una función similar a la de la persistencia local de la función de ajuste, esta función no pudo realizar un ajuste suficiente entre dichos dígitos y sus referencias, generando así el alto número de errores. Podríamos pensar entonces que eliminando esta restricción se podría resolver este problema, pero las pruebas preliminarmente demuestraron que esta restricción es imprescindible.

Es necesario pensar según los puntos capitales que de este experimento se deducen para que se realicen en la línea de investigaciones que se describen resultados obtenidos para palabras aisladas así como también que una base de referencias mejor distribuida permitiría obtener mejores resultados para palabras combinadas. En la línea de estas conclusiones un mayor número de referencias para cada dígito representaría las diferentes pronunciaciones y variaciones que pueden darse ante una misma referencia de acuerdo con pronunciación misma siguiendo con cada palabra aislada dentro de una misma base de palabras.

La creación de una misma base de referencia para reconocimiento de palabras aisladas y de palabras conectadas podría además permitir una comparación de los resultados obtenidos. Es importante mencionar que la creación de tal base de datos es un proceso que requiere de mucho tiempo y que podría ser, por ejemplo un trabajo, de servicio social extremadamente útil ya que además de ofrecer las ventajas mencionadas, permitiría trabajar con palabras en español. Por las características de los algoritmos empleados podemos anticipar que los resultados no deberían variar mucho, pero es importante sin embargo verificarlo prácticamente.

Conclusiones

De las experiencias obtenidas mediante el desarrollo de este trabajo podemos afirmar que el modelo fuente-filtro basado en los coeficientes LPC es una modelización razonable de la voz. Este puede ser empleado tanto para síntesis de señales como para reconocimiento. En el caso de síntesis, es deseable emplear un orden de filtro mayor a 21 para obtener señales de calidad, aunque señales comprensibles se obtienen con órdenes menores. En el caso de reconocimiento de palabras aisladas, un orden $p=7$ fue suficiente para obtener resultados excelentes.

La determinación de las bases de reconocimiento obtenidas en el caso de palabras conectadas se debe principalmente al empleo de una base de referencia poco adecuada. Es importante recordar aquí que una base de bases de referencia para reconocimiento de palabras conectadas debe incluir además de diferentes pronunciaciones que representen las diferentes entonaciones para cada letra, pronunciaciones que representen las otras vocalidades de pronunciación observadas en la pronunciación de palabras de palabras. Por estos motivos podemos anticipar un aumento substancial en las bases de reconocimiento de palabras conectadas, mediante el uso de una base de bases más adecuada.

Una alternativa a los algoritmos de búsqueda de palabras conectadas son los algoritmos basados para el estudio de la síntesis, que al igual que en el caso de palabras aisladas, determinan así una medida de similitud entre palabras, pero al igual que requieren un cómputo de gran complejidad, en particular cuando el número de palabras es grande. Sin embargo, como se puede llegar a determinaciones razonables de las más apropiadas que producen la obtención de distancias pequeñas para señales desconocidas. En el caso de palabras aisladas el algoritmo más eficiente es el presentado por Li, cuando una palabra es desconocida se calcula un M que es el caso de palabras aisladas, cuando se desconoce una palabra.

La creación de una base de datos completa en español es necesaria y urgente para poder continuar con la investigación en esta área en la Facultad. Esta base podría ser creada a través de un programa de servicio social ya que se trata de una tarea que requiere de mucho tiempo.

Bibliografía

Libros:

Thomas W. Parsons. Voice and Speech Processing. McGraw-Hill, 1987

L. R. Rabiner and R. W. Schaefer. Digital Processing of Speech. Prentice-Hall, 1978

F. J. Owens. Signal Processing of Speech. McGraw-Hill, 1993

J. R. Deller, J. G. Proakis and J. H. Hansen. Discrete-Time Processing of Speech Signals. MacMillan, 1993

Alex Waibel and Kai-Fu Lee. Reading in Speech Recognition. Morgan Kaufmann Publishers, 1990

Albert S. Bregman. Auditory Scene Analysis. MIT Press, 1990

Artículos:

D. R. Reddy. Speech Recognition in Machines. A. Karim. IEEE Proceedings, April 1976

R. W. Schaefer and L. R. Rabiner. Digital Representation of Speech Signals. The Institute of Electrical and Electronics Engineers, 1975

L. R. Rabiner and S. E. Levinson. Discrete-Time and Continuous-Time Representations of Speech and Selected Applications. IEEE, 1971

F. Itakura. Minimum Prediction Residual Principle Applied to Speech Recognition. IEEE, 1974

H. Sakoe and S. Chiba. Dynamic Programming Algorithm Optimization for Spoken Word Recognition. IEEE, 1978

H. Sakoe. Two-Level DP-Matching--A Dynamic Programming-Based Pattern Matching Algorithm for Connected Word Recognition. IEEE, 1979

H. Ney. The Use of a One-Stage Dynamic Programming Algorithm for Connected Word Recognition. IEEE, 1984