



ETIQUETADO NO FONEMATICO PARA UN SISTEMA AUTOMATICO MONOLOCUTOR DE =====

RECONOCIMIENTO DE PALABRAS AISLADAS =====

Por Antonio J. Rubio Ayuso (*) y M. Carmen Carrion Perez (**)
(*) Departamento de Electronica
(**) Departamento de Electricidad
Universidad de Granada

I. Introduccion -----

Presentamos en estas lineas una de las partes de nuestro trabajo que pueden considerarse terminadas desde el punto de vista logico. Trabajamos en otros temas dentro del mismo campo, pero eso lo presentaremos en otros numeros.

Aqui presentamos un Sistema Automatico de Reconocimiento Monolocutor con Etiquetado No Fonemático (SARMENF) (1). Funciona a base de dividir la señal sonora en segmentos de la misma duracion e intentando clasificar cada segmento entre unos prototipos obtenidos durante la Fase de Entrenamiento, mediante un criterio de distancia euclidea. La etiqueta que se asigna a cada segmento no guarda ninguna relacion con los fonemas del lenguaje. Este tipo de etiquetado puede ser ventajoso debido a la difícil localización del fonema como unidad de analisis.

II. Esquema general -----

Hay dos fases en el funcionamiento del SARMENF: fase de entrenamiento y fase de reconocimiento.

2.1 Fase de Reconocimiento -----

Durante la Fase de Reconocimiento se realizan los siguientes pasos:

1) Una vez que el locutor ha pronunciado una palabra del vocabulario permitido, el primer paso consiste en muestrear la señal, in-



roduciendo las muestras en un ordenador.

La frecuencia de muestreo es de 10 kHz. La conversión A/D se hace a través de un conversor A/D de 12 bits. Para cada palabra, se toman 15000 muestras, durante 1'5 s.

2) Con objeto de reducir el volumen de cálculo, los extremos de la señal sonora son localizados. En general, ninguna palabra ocupa completamente las 15000 muestras.

3) Una vez que los extremos han sido localizados, se establece una segmentación regular de la señal sonora. Los segmentos son de 10 ms (100 muestras) y se solapan a la mitad.

Sobre cada segmento se calcula un conjunto de parámetros. En primer lugar se hace una clasificación previa entre segmentos sonoros y no sonoros, utilizando la Energía y el Número de pasos por cero del segmento (ZCR). Se clasifican como segmentos sonoros aquellos que cumplen la inequación

$$\text{Energía} > A+B \cdot \text{ZCR}$$

donde las constantes $A = .1605E+08$ y $B = 1.075E+06$ se han determinado experimentalmente y son función de la ganancia global del sistema de audio.

Se clasifican como no sonoros aquellos segmentos que cumplen la desigualdad

$$\text{Energía} > C+D \cdot \text{ZCR}$$

donde las constantes $C = .2012E+08$ y $D = 8.047E+05$ también se han determinado experimentalmente.

Aquellos segmentos que no cumplen ninguna de las dos condiciones anteriores son clasificados como dudosos y son ignorados en el tratamiento posterior.

Si un segmento se clasifica como sonoro, se calculan sobre el otros cuatro parámetros adicionales: los tres primeros formantes y el logaritmo de la razón de las áreas de la primera y segunda sección del tubo acústico equivalente al tracto vocal. Esta determinación se hace mediante Predicción Lineal (2).

Por otra parte, si el segmento se clasifica como no sonoro, los cuatro parámetros adicionales son los valores medios del espectro de la señal dividido en cuatro zonas o bandas (0 - 1.25 kHz, 1.25 - 2.5 kHz, 2.5 - 3.75 kHz y 3.75 - 5 kHz).

4) Los parámetros correspondientes a un segmento se comparan con los de los Segmentos de Referencia obtenidos en la Fase de Entrenamiento. De esta comparación se obtiene una etiqueta para el segmento en cuestión. Así se consigue una cadena de etiquetas (una etiqueta por cada 5 ms) que representa a la palabra emitida.

5) El último paso en la Fase de Reconocimiento consiste en comparar la cadena de etiquetas con las Cadenas de Referencia, obtenidas también en la Fase de Entrenamiento (una Cadena de Referencia por cada palabra del vocabulario). Esta comparación se lleva a cabo por medio de un algoritmo de Programación Dinámica (3). Este algoritmo permite comparar cadenas de etiquetas de diferente longitud, pro-



porcionando la "Similitud" entre ellas en función de la distancia euclídea entre los diferentes símbolos de cada una de las cadenas comparadas.

Un esquema de bloques para la Fase de Reconocimiento puede verse en la parte derecha de la Figura 1.

2.1 Fase de Entrenamiento

Esta Fase está representada en la parte izquierda de la Figura 1.

El Entrenamiento consiste en la pronunciación de todas las palabras del vocabulario, cada una de ellas repetida cuatro veces.

Después del muestreo de cada palabra de entrenamiento, se calculan sus extremos y se obtienen los parámetros correspondientes a los segmentos de la señal, del mismo modo que en la Fase de Reconocimiento.

Cada segmento se considera como un punto de un hiper-espacio euclídeo de seis dimensiones. En realidad se manejan dos hiper-espacios: uno para los segmentos sonoros y otro para los no sonoros. De este modo, la voz del locutor se caracteriza por una cierta distribución de puntos a lo largo de ambos hiper-espacios.

Los pasos que siguen en la Fase de Entrenamiento son:

1) Análisis de las nubes formadas por los segmentos de las palabras pronunciadas, en ambos hiper-espacios.

Aquellas nubes claramente separadas del resto de puntos se consideran como correspondientes a segmentos de señal del mismo tipo. Su centro geométrico representará al Segmento de Referencia, con una etiqueta asociada.

La obtención de las nubes y Segmentos de Referencia se lleva a cabo mediante el algoritmo "k-medias" modificado (4).

La modificación puede describirse brevemente. La principal fuente de fallos en el algoritmo "k-medias" original está en la selección del conjunto inicial de centros. Evitamos este problema de la siguiente forma:

Se hace una normalización de todos los parámetros de modo que la media de cada parámetro sea 1. Así, el punto $(1,1,1,1,1,1)$ es el centro de la distribución completa. Entonces, se colocan los puntos iniciales de modo que estén uniformemente repartidos alrededor del punto $(1,1,1,1,1,1)$.

El número de Segmentos de Referencia (k) se ha establecido experimentalmente. Para los segmentos sonoros $k = 21$ y para los segmentos no sonoros $k = 7$.

2) De acuerdo con los Segmentos de Referencia, cada palabra del vocabulario se analiza con el fin de elaborar una cadena de etiquetas (Cadena de Referencia) para esa palabra. Estas Cadenas de Referencia son las que se usan en el algoritmo de Programación Dinámica (Fase



de Reconocimiento) para comparar una palabra emitida con todas las del vocabulario.

III. Conclusion

En este trabajo se ha presentado un sistema automatico de reconocimiento monolocator con etiquetado no fonemático. Las principales características son:

1) Es un sistema de reconocimiento monolocator para palabras aisladas.

2) Se lleva a cabo un etiquetado no fonemático. Hay 21 diferentes tipos de segmentos sonoros y 7 tipos de segmentos no sonoros.

3) El vocabulario es un conjunto de 19 palabras castellanas relacionadas con el calculo aritmético: "Cero", "Uno", "Dos", "Tres", "Cuatro", "Cinco", "Seis", "Siete", "Ocho", "Nueve", "Parentesis", "Cierra", "Mas", "Menos", "Por", "Partido", "Raiz", "Igual?" y "Punto". Este vocabulario permitira el manejo oral del ordenador en calculos aritmeticos simples.

4) Los resultados experimentales estan hasta el momento alrededor del 65% de reconocimiento correcto. Merece la pena destacar que no se ha puesto especial atencion en la seleccion de las Cadenas de Referencia.

5) El sistema ha sido implementado sobre un ordenador ECLIPSE-250 de DATA GENERAL, en lenguaje de programacion FORTRAN V.

IV. Referencias

- (1) R. Reddy
"Speech recognition by machine: a review"
Proc. IEEE, vol. 64, pp. 505-531
 - (2) J. D. Markel, A. H. Gray
"Linear Prediction of Speech"
Springer-Verlag, 1976.
 - (3) H. Sakoe, S. Chiba
"Dynamic programming algorithm optimization for spoken word recognition"
IEEE Trans. on ASSP. Vol. ASSP-26, n. 1, pp. 43-49, 1978.
 - (4) J. T. Tou, R. C. Gonzalez
"Pattern recognition principles"
Addison-Wesley, 1974.
-

Nota: Este trabajo ha sido presentado al "1983 SPAIN WORKSHOP ON SIGNAL PROCESSING AND ITS APPLICATIONS" que se celebrara en Sitges (Barcelona) en Septiembre de 1983.

ESQUEMA - RESUMEN

Martin S. Ruiperez, Catedrático de Filología Griega
Universidad Complutense, Madrid

Aparte de temas de Filología Clásica y Lingüística Comparada Indoeuropea ha trabajado en:

- Lingüística estructural y funcional (griego y español), concretamente morfosintaxis y fonología.
- Con equipo de alumnos y colaboradores, sin ordenador y sólo con calculadoras, ha establecido (1968) un vocabulario de frecuencias de griego ático de 400-350 a. C. Igualmente un índice morfológico verbal de frecuencias. Todo ello para la elaboración de un método de enseñanza del griego antiguo, que está en fase de experimentación en varios Institutos y Facultades.
- Ha dirigido una tesis doctoral de lingüística cuantitativa aplicada al problema de los arcaísmos en Homero.
- Desde 1973 se ha familiarizado con cálculo programado.
- Desde 1980 se ha familiarizado con ordenadores, concretamente con un HEWLETT-PACKARD HP250, provisto de base de datos.

Además de programas comerciales, ha realizado programas lingüísticos:

1. Construcción de cadenas a partir de segmentos conmutables. Concretamente, conversión de un número introducido en dígitos en texto español, con perfecta coherencia sintáctica.
2. Corte de palabras conforme a las reglas de la Academia y al buen uso tipográfico para programas de textos cuando al final de una línea no cabe la palabra entera. Versiones para Castellano, catalán, francés e italiano. Disponible en el mercado. Tiene un programa de calendario (que incluye la reforma del calendario gregoriano) para determinar en qué día de la semana cae cualquier fecha entre el 1 de enero del año 1 (!) y el 28 de febrero del 2100; número de días entre dos fechas y fecha a partir de otra y de un número de días de intervalo.

Todo en lenguaje BASIC y "Extended BASIC" o "Business BASIC" (en realidad BASIC enriquecido con PASCAL)

- Desde diciembre de 1982 dispone en su cátedra de un modesto equipo compuesto por:
 - Ordenador HEWLETT-PACKARD HP87 con 128 Kb de memoria RAM, más un ROM de Input/Output y otro de Assembler.
 - Periférico de una unidad dual de diskettes floppies con capacidad total de 540 Kb
 - Periférico de impresora de agujas EPSON III/100

Actualmente con este equipo está desarrollando "software" para tratamiento científico de textos: introducción, almacenamiento, modificación, impresión, búsqueda y confección de índices alfabéticos (con alfabetización distinta para el griego, en el que, por ejemplo, la z no va al final sino después de la e).

Se prevé la aplicación inmediata al análisis de textos métricos griegos codificados y de un repertorio de inscripciones.

- En 1982-83 curso impartió el curso de "Investigación lingüística y filológica" en Sevilla, 15 de abril de 1983.

