

Pruebas y Validación de un Sistema de Reconocimiento del Habla Basado en Sílabas con un Vocabulario Pequeño.

Dr. Sergio Suárez Guerra¹, M. en C. José Luis Oropeza Rodríguez², Ing. Mariana Del Villar³ Ing. Karen Suso⁴
j_oro2002@hotmail.com

1. Profesor Investigador del Centro de Investigación en Computación (CIC, IPN)
2. Alumno de doctorado del Centro de Investigación en Computación (CIC, IPN)
3. Ingeniera en Telemática (UPIITA, IPN)
4. Ingeniera en telemática (UPIITA, IPN)

CIC Av. Juan de Dios Batis s/n casi esq. Miguel Otón de Mendizábal, Unidad Profesional "Adolfo López Mateos" Edificio CIC. Col. Nueva Industrial Vallejo, 07738, México, D.F.
UPIITA Av. Instituto Politécnico Nacional 2580, Col. La Laguna Ticomán, 07034, México D.F.

RESUMEN.

El presente trabajo contempla los resultados obtenidos del análisis de un corpus de voz discontinuo segmentado en unidades silábicas. Los resultados contemplan la utilización de técnicas de reconocimiento tales como la codificación predictiva lineal y los modelos ocultos de Markov.

Se realizó un análisis a un corpus de 20 palabras, las cuales son: *cero, uno, dos, tres, cuatro, cinco, seis, siete, ocho, nueve, azul, blanco, café, gris, negro, rojo, rosa, verde, abrir y cerrar*, por medio de división silábica y basado en los Modelos Ocultos de Markov.

Dichas palabras generaron 31 sílabas como elementos de reconocimiento, las cuales nos permitieron crear un determinado conjunto de patrones (modelos) para ser reconocidos posteriormente.

PALABRAS CLAVE: Viterbi, HMM, LPC, sílaba.

INTRODUCCIÓN.

El reconocimiento de voz es una rama de las ciencias de la computación que se encuentra aún en proceso de crecimiento. Los distintos elementos que la conforman han llegado a generar distintas herramientas prácticas que plasman tales avances científicos. Sin embargo, la búsqueda de nuevas alternativas como sucede en otras ramas de la ciencia es de carácter relevante. Las sílabas conforman un patrón de estudio bastante útil en esta rama, el presente trabajo muestra tales hechos. En la figura 1 se plantea la solución que se le da al problema del reconocimiento de voz usando cadenas ocultas de Markov^[1]. El cual se inicia introduciendo la señal de voz por medio de un micrófono. En la siguiente etapa se segmenta en sílabas la señal obtenida, cabe destacar que

las etapas posteriores se le realizarán a cada una de las sílabas que componen la palabra introducida. El procesamiento digital de la señal (PDS) se realiza para reducir la cantidad de información que se debe de procesar, extrayendo únicamente las características acústicas sobresalientes de las sílabas.

El reconocimiento del habla se realiza con una etapa previa de entrenamiento^[2], la cual es la encargada de generar un libro código global^[3] con las muestras adquiridas, utilizando el método de Codificación Predictiva Lineal, después con este libro código se generan las Cadenas Ocultas de Markov. Una vez realizado lo anterior, se pueden reconocer las palabras contenidas en el diccionario.

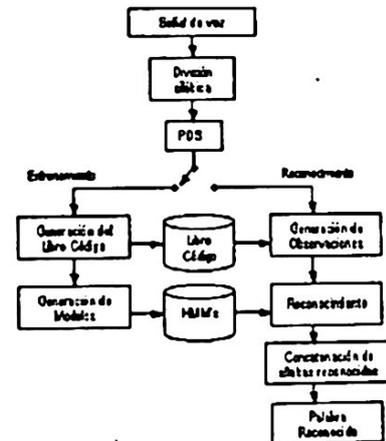


Fig 1. Diagrama a bloques de la solución propuesta.

Tal procesamiento de información es obtenido tras la solución a los tres problemas de Markov, donde el primero de ellos (obtención de la probabilidad de que una observación haya sido producida por un modelo) es de la forma^[1]:

voz original, se mostrará la caja de diálogo de la figura 5.

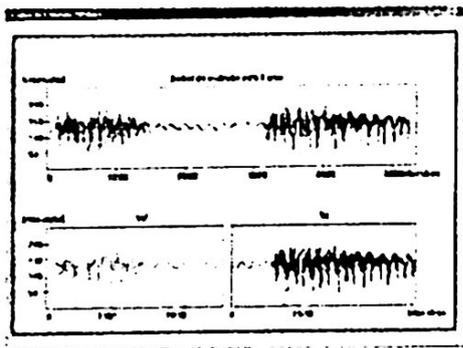


Fig. 5. Gráficas de sílabas.

III. PRUEBAS Y RESULTADOS.

En la tabla 1 se muestra el promedio de duración, en segundos, y el número de muestras que contiene cada palabra del diccionario., cabe destacar que se usaron 10 repeticiones por palabra del diccionario.

Palabra	Tempo (seg)	# Muestras
Ort	0.296	4400
Orn	0.432	6474
Ort	0.44	6572
Ort	0.445	6670
Ort	0.502	7546
Ort	0.572	8622
Ort	0.596	8938
Ort	0.626	9372
Ort	0.64	9706
Ort	0.649	9814
Ort	0.662	10024
Ort	0.669	10152
Ort	0.669	10152
Ort	0.67	10152
Ort	0.677	10294
Ort	0.677	10294
Ort	0.720	10832
Ort	0.764	11518
Ort	0.764	11518

Tabla1. Promedio de duración de las palabras del diccionario.

En la figura 6 se muestra la gráfica del promedio de duración de las palabras del diccionario. En ésta se puede ver claramente que estos parámetros (la duración y el número de muestras) se pueden utilizar para este proceso.

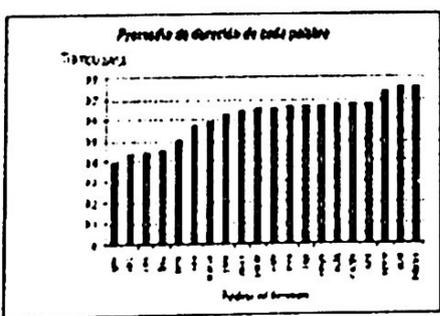


Fig. 6. Gráfica del promedio de duración de cada palabra.

Distribución de regiones

Sabemos que dado un libro código se puede partir un espacio muestra X en N regiones disjuntas y a cada una de estas se les asocia un centroide y a cada una se optimiza para convertirlo en el mejor representante de la región, el que en promedio difiere menos de cada uno de los elementos de la región^[3].

Podemos observar la distribución de vectores en figura 7, los cuales en conjunto forman 4 regiones, su centroide fue optimizado 7 veces, a medida que los vectores cambiaban de región, se obtenía un nuevo centroide, y había que recalcarlo, por ello observamos variaciones al principio de la gráfica, pero también observamos que se vuelve estable, pues el centroide ya no cambia en las últimas optimizaciones y por consiguiente, ya no hay saltos de vectores hacia otras regiones.

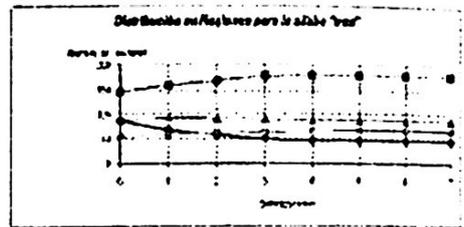


Fig. 7. Gráfica de la Distribución en regiones de la sílaba "tres" para 4 Regiones.

En la figura 8 se observa el mismo proceso, con diferencia que ahora se están modelando 16 regiones y la estabilidad se logró hasta la optimización número 12.



Fig. 8. Gráfica de la Distribución en regiones de la sílaba "tres" para 16 Regiones.

En las gráficas se puede observar que el desempeño del sistema fue el esperado porque fue distribuyendo los datos hasta lograr el centroide óptimo, sin adecuada distribución no se podrían obtener coeficientes LPC deseados para lograr reconocimiento. Cabe destacar que el proceso realizó hasta 128 regiones, pero no se muestra debido a la cantidad de datos que se está manejando, por esta razón sólo mostramos la distribución final, la cual se muestra en la figura 9.

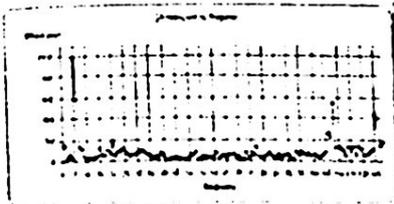


Fig. 9. Gráfica de distribución de regiones final.

En la generación de los Modelos Ocultos de Markov para sílabas se determinó el número de optimizaciones con los que trabaja mejor el sistema, se estuvo evaluando para 32 regiones, con 200 muestras con las cuales fue entrenado el sistema, 2 locutores y como se muestra en la gráfica 10 de la figura el número de errores esta por encima de 45.

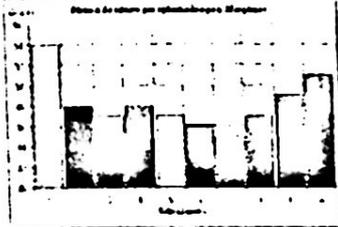


Fig. 10. Errores por optimización para 32 regiones con 200 muestras para sílabas.

En la figura 11 se muestra el desempeño del sistema para 64 regiones las mismas muestras, y como se puede observar, el número de errores ha disminuido, ahora se encuentra por encima de 14.

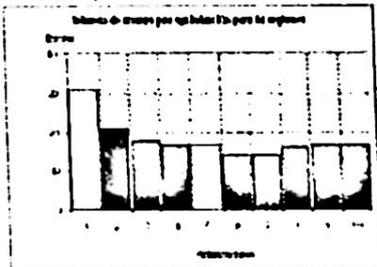


Fig. 11. Errores por optimización para 64 regiones con 200 muestras para sílabas.

Por último, la figura 12 muestra la gráfica del número de regiones elegido, el cual es de 128 pues como se puede observar, el desempeño es mucho mejor que en los casos anteriores, no continuamos incrementando el número de regiones debido al tamaño del diccionario utilizado en el sistema, otra de las razones fue que representaría un incremento innecesario en el tiempo de cómputo.

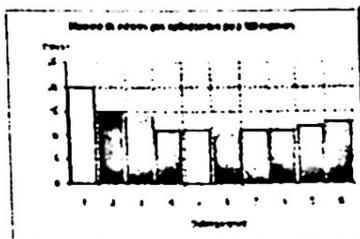


Fig. 12. Errores por optimización para 128 regiones con 200 muestras para sílabas.

Como se puede observar en las gráficas del número de errores por optimización, existe un mejor desempeño en el intervalo de optimizaciones de 5 a 7, porque después de este intervalo el número de errores comienza a incrementarse. Por ésta razón se tomó que el número de optimizaciones en la generación de los Modelos Ocultos de Markov sea de 6.

Se hizo el análisis anterior pero con más locutores para determinar de que manera influye la variabilidad de las distintas señales de voz pues ahora se realizó con 4 locutores y con 400 muestras. En la figura 13 se muestra la gráfica del número de errores por optimización para 32 regiones, como se puede observar el menor error se encuentra en la optimización 6.

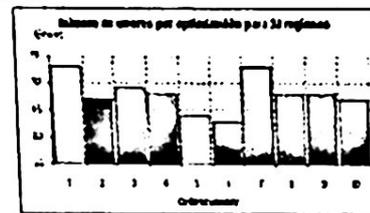


Fig. 13. Errores por optimización para 32 regiones con 400 muestras para sílabas.

En la figura 14 se muestra el desempeño del sistema para 64 regiones con las mismas muestras, y como se puede observar, el número de errores ha disminuido, ahora está por encima de 39.

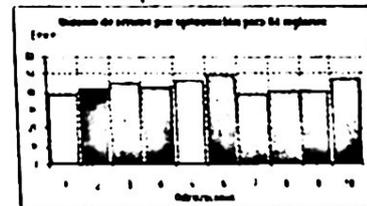


Fig. 14. Errores por optimización para 64 regiones con 400 muestras para sílabas.

En la figura 15 se muestra la gráfica con 128 regiones y como se puede observar, el desempeño es mucho mejor que en los casos anteriores, pues con las mismas muestras los errores esta por encima de 20.

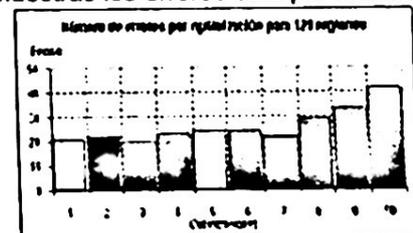


Fig. 15. Errores por optimización para 128 regiones con 400 muestras para sílabas.

A continuación se mostrarán el número de errores por optimización en la generación de los Modelos Ocultos de Markov para el caso de palabras, en la figura 16 se muestra la gráfica para 200 muestras con las

Aunque el sistema es inestable se encontró un punto en donde el índice de reconocimiento fue satisfactorio en un tiempo de cómputo razonable, se realizó una prueba para determinar el número de optimizaciones en donde se alcanza la estabilidad del sistema, la cual consiste en optimizarlo hasta que la probabilidad de una observación dado un modelo (Problema 1 de Markov, algoritmo hacia adelante) no cambie, en dicha prueba se observó que el número de optimizaciones oscila entre 700 y 800 para alcanzar su estabilidad, pero el inconveniente es el tiempo de cómputo ya que se encuentra entre 6 y 7 horas.

5. BIBLIOGRAFIA.

- [1] Fundamentals of Speech Recognition.
Lawrence Rabiner and Biing-Hwang Juang.
Prentice Hall, 1993
- [2] A Hybrid Systems with Symbolic AI and Statical Methods for Speech Recognition.
Jesus Savage Carmona
University of Washington, 1995.
- [3] Reconocimiento de palabras aisladas usando cuantización vectorial.
Ricardo Barrón Fernández.
Centro de Investigación en Computación. Octubre 1998.
- [4] Reconocimiento de voz usando estructuras silábicas.
Dr. Sergio Suárez Guerra, José Luis Oropeza
CIC 2002. noviembre 2002
- [5] Arquitecturas y métodos en sistemas de reconocimiento automático del habla de gran vocabulario.
Javier Macías Guarasa
Universidad Politécnica de Madrid 2001.
- [6] Reconocimiento de voz y fonética acústica. Pedro Gómez Vilda.
Mc Graw Hill, 2000.