



INSTITUTO TECNOLÓGICO
DE
BUENOS AIRES

INGENIERÍA ELECTRÓNICA
(22.47) - PROCESAMIENTO DE VOZ

TRABAJO PRÁCTICO 1

VOCODER DE APLICACIÓN MUSICAL EN TIEMPO REAL.

Grupo:

Díaz, Ian Cruz

Rodriguez Turco, Martin

Sebastián

Legajos:

57515

56629

Contenido

1. Introducción	2
2. Enfoque LPC	3
2.1. Modelo de producción de voz	3
2.2. Implementación	3
2.3. Pulsos glotales	5

1. Introducción

Para el siguiente proyecto se busco implementar un **Vocoder** de aplicación musical. Un vocoder es un sistema que analiza el sonido de una persona hablando, y mediante diferentes algoritmos realiza cambios en los pulsos glotales sin alterar el filtro articulatorio, y de esta manera logra sintetizar un sonido artificial de la señal de voz de la persona que este hablando. Junto con esto, y eligiendo frecuencias fundamentales de los pulsos glotales acordemente, se pueden generar melodías y acordes para utilizar en ambientes musicales.

2. Enfoque LPC

Para poder realizar este sistema se decidió utilizar el enfoque *Linear Prediction Coefficients* (LPC), que se detallará brevemente a continuación.

2.1. Modelo de producción de voz

Si recordamos del modelo de producción de voz, podemos simplificar el modelo en los bloques de la Figura 2.1.

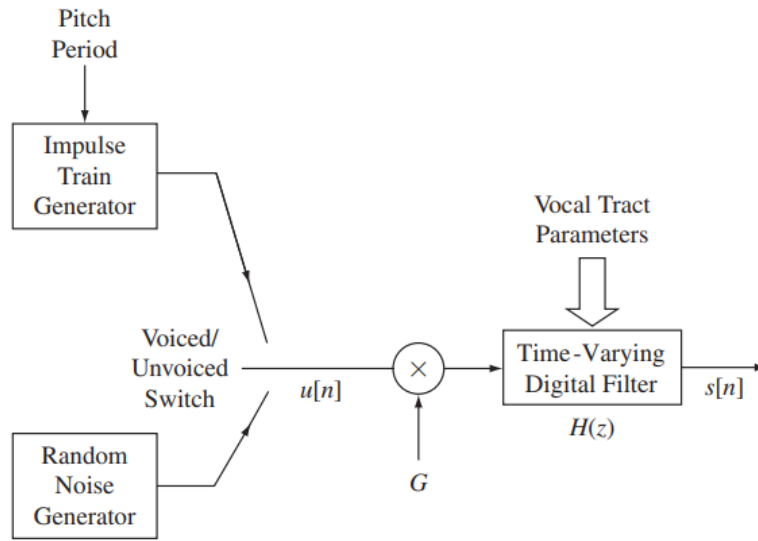


FIGURA 2.1: Diagrama en bloques del modelo de producción de voz

Realizando algunos cálculos [2] [3], se llega a la conclusión que el filtro error de predicción estima la inversa del filtro articulatorio, por lo tanto, simplemente estimando el filtro error de predicción podemos utilizar el filtro articulatorio obtenido de una porción de audio para sintetizar unos pulsos glotales artificiales.

2.2. Implementación

Para la implementación de este sistema, se realizó el método de la autocorrelación para estimar los coeficientes del filtro de error de producción y una ventana de *hann* para realizar la sintetización¹. Esto se puede ver en el código a continuación²

¹Recordemos que para sintetizar una señal correctamente debemos cumplir con el requisito de que la suma de las ventanas en el tiempo debe ser constante, y por lo tanto, con las ventanas de hann, el overlap debe ser del 50 %.

²Este código fue basado en [1], donde ese código fue implementado para poder sintetizar la voz con pulsos artificiales, pero intentando modificar lo menos posible la voz de salida respecto de la entrada,

```

1 def vocode(signal, fs, f_custom, block_len, overlap, order, prev_block
    , p_coverage=0.01, unvoiced2zeros=True, glotal_type="triang"):
2
3     out = zeros(len(signal))
4     out[:len(prev_block)] = prev_block
5     prev_delta = 0
6     T_samples = int((fs / f_custom))
7
8     for block, idx in block_process(signal, fs, block_len, overlap):
9         gain_correction = (1 - overlap) * 2 # *2 due to hann window
10        block *= hann(len(block)) * gain_correction
11
12        rxx = correlate(block, block, mode='full')
13        rxx = rxx[len(rxx) // 2:]
14        period_samples, is_voiced = fundamental_period_estimate(rxx, fs)
15        a = -solve_toeplitz(rxx[:order], rxx[1:order + 1])
16        a = np.concatenate([1], a)
17        error_power = rms(lfilter(a, (1,), block))
18        if is_voiced:
19            vocoded, new_delta = pitch_maker(len(block), T_samples,
20            prev_delta, overlap=overlap)
21            prev_delta = new_delta
22            impulse_response = glotales[glotal_type](len(block), p_coverage=
23            p_coverage)
24            vocoded = np.convolve(vocoded, impulse_response, mode="same")
25        else:
26            if unvoiced2zeros:
27                vocoded = np.zeros(len(block))
28            else:
29                vocoded = randn(len(block)) / 2
30
31        vocoded = lfilter((error_power,), a, vocoded)
32        vocoded *= hann(len(block))
33        out[idx:idx + len(block)] += deemphasis(vocoded)
34    return out

```

LISTING 1: Vocoder

En este código se puede ver que, lo primero que se realiza es dividir la señal de entrada en bloques de *block_len* milisegundos, con un overlap de *overlap* % ³, luego, se trabajará por bloques. Como se puede ver, en cada bloque, lo primero que se hace es obtener la autocorrelación $R_{xx}(\tau)$ para $\tau \geq 0$, y se llama a la función *fundamental_period_estimate*, esta función definida por el creador del código original se utilizaba

por ejemplo para utilizarla en codificación de voz.

³En este caso, el overlap siempre debería ser del 50 % ya que solo se utiliza ventanas de *hann*

para obtener el período de los pulsos glotales así como para saber si la señal es sonora o sorda. Si ahondamos mas profundo en la función y con lo que nos interesa (saber si la señal es sonora o no), podemos ver que el método que se realiza para determinar si la señal es o no sonora, simplemente se busca cual es el máximo de la autocorrelación, y si esta supera un cierto threshold, se determina que es sonora, esto se puede ver en el Listing 2. Una vez hecho esto, se determinan los coeficientes del filtro error de producción con *solve_toeplitz*, y se estima la potencia de error para hacer el control automático de ganancia (AGC) al final de la síntesis. Finalmente, cuando la señal de entrada es determinada como sonora, se realizan los pulsos glotales deseados y se los filtran por el filtro $H(z) = \frac{G}{A(z)}$, siendo $A(z)$ el filtro error de predicción y $G^2 = e(n)$.

```

1 def fundamental_period_estimate(rxx, fs):
2     """
3     Calculates the fundamental frequency of an auto correlation array.
4     rxx    the auto correlation array.
5     fs     the sample rate in hertz.
6     """
7     f_low, f_high = 50, 250
8     f_low_idx = round(fs / f_low)
9     f_high_idx = round(fs / f_high)
10    period_idx = argmax(rxx[f_high_idx:f_low_idx]) + f_high_idx
11    is_voiced = max(rxx) > 0.20
12    return (period_idx, is_voiced)

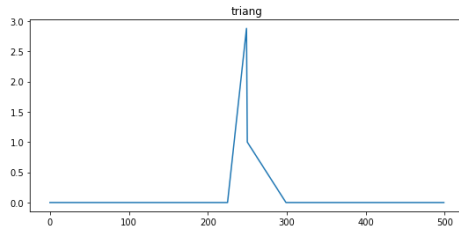
```

LISTING 2: Fundamental Period Estimate

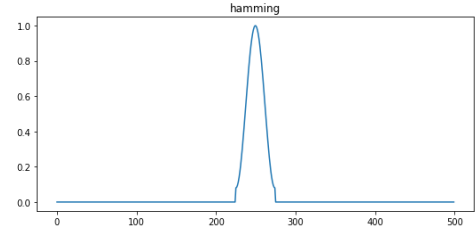
2.3. Pulsos glotales

Para la implementación de pulsos glotales sintéticos, se utilizaron 4 formas de pulsos distintas, ellas son:

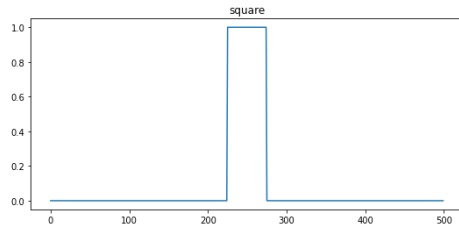
- **Triangular:** Cuenta con una amplia personalización donde se es capaz de ajustar la pendiente de crecimiento, la de decrecimiento, la altura, el ancho, el final de la primera rampa, y el comienzo de la segunda rampa.
- **Hamming:** Es una ventana de hamming centrada, con ancho ajustable y amplitud ajustable.
- **Cuadrada:** Es un pulso cuadrado con personalización del ancho y la altura del mismo
- **Exponencial:** Un pulso centrado en 0 donde la parte negativa es una exponencial creciente y la parte positiva una exponencial decreciente.



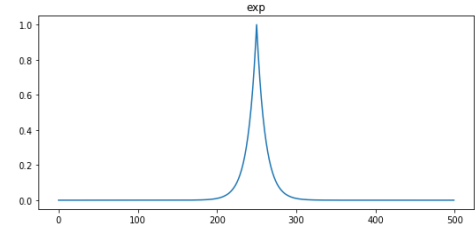
(A) Pulso glotal Triangular



(B) Pulso glotal Hamming



(C) Pulso glotal Cuadrada



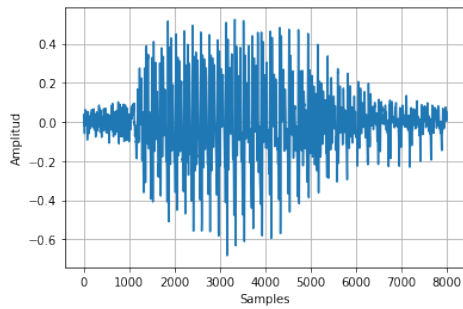
(D) Pulso glotal Exponencial

FIGURA 2.2: Pulsos glotales

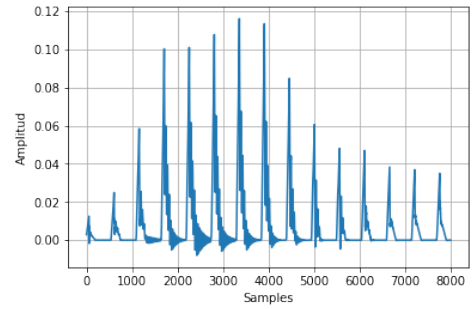
Ejemplos de estos pulsos se pueden ver en la figura 2.2. Por otro lado, si sintetizamos una señal de voz de ejemplo con estos pulsos glotales, con los siguientes parámetros:

- Frecuencia = 500 (Hz)
- Overlap = 50 %
- Bloque = 32 (ms)

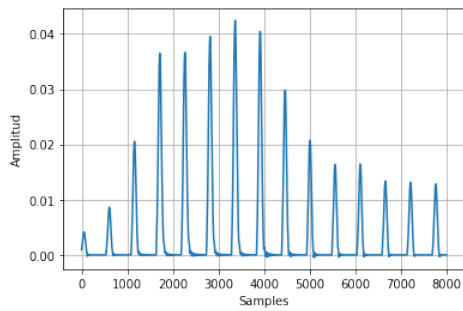
Obtenemos los resultados que se pueden observar en la Figura 2.3.



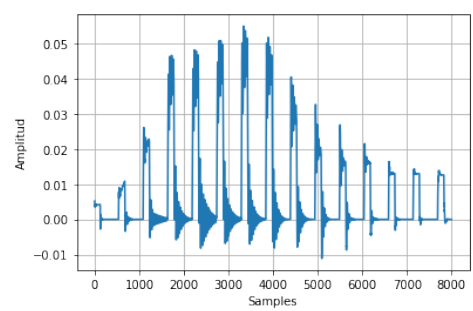
(A) Señal original



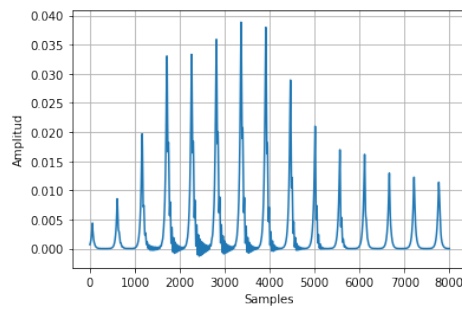
(B) Síntesis con triangular



(C) Síntesis con Hamming



(D) Síntesis con Cuadrada



(E) Síntesis con Exponencial

FIGURA 2.3: Síntesis de señales con diferentes pulsos glotales

Referencias

- [1] Bastian Bechtold. *Pocoder*. URL: <https://github.com/bastibe/pocoder>.
- [2] Thomas F. Quatieri. *Discrete Time Speech Signal Proccesing: Principles and practice*. 1.^a ed. Prentice-Hall, Inc, 2002. ISBN: 0-13-242942-X.
- [3] Lawrence R. Rabiner y Ronald W. Schafer. *Introduction to Digital Speech Processing*. 1.^a ed. now Publishers Inc., 2007. ISBN: 978-1-60198-070-0.