

語者辨識

Speaker Recognition

陳慶瀚

機器智慧與自動化技術(MIAT)實驗室

義守大學電機系

pierre@isu.edu.tw

2005年12月20日



語者辨識課題

語者辨認系統可分為兩個課題：

- 語音特徵擷取 (Feature Extraction)
 - 分類器 (Classifier)

重要的研究方向：

- 提昇辨識性能
- 簡化計算複雜度
- 硬體架構

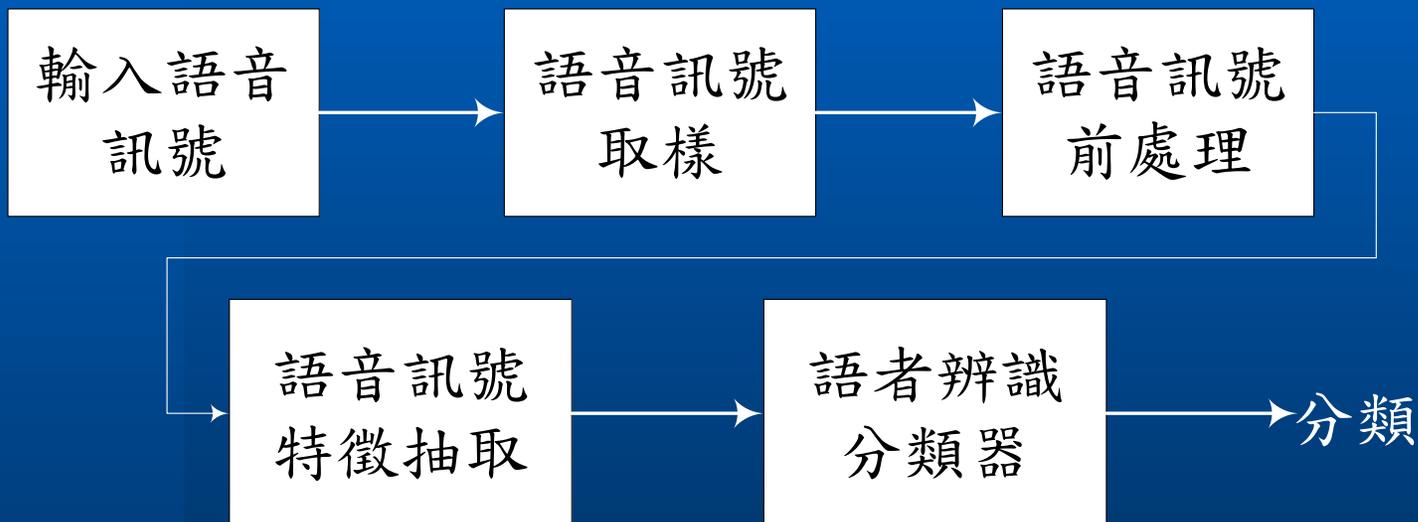


語者辨識模式

- 根據使用者的聲音訊號輸入，辨認使用者身份
- 辨識模式
 - 語者識別 (Speaker Identification)
 - 語者驗證 (Speaker Verification)
- 語音輸入模式
 - 文字相關 (Text-Dependent) 模式
 - 文字不相關 (Text-Independent) 模式

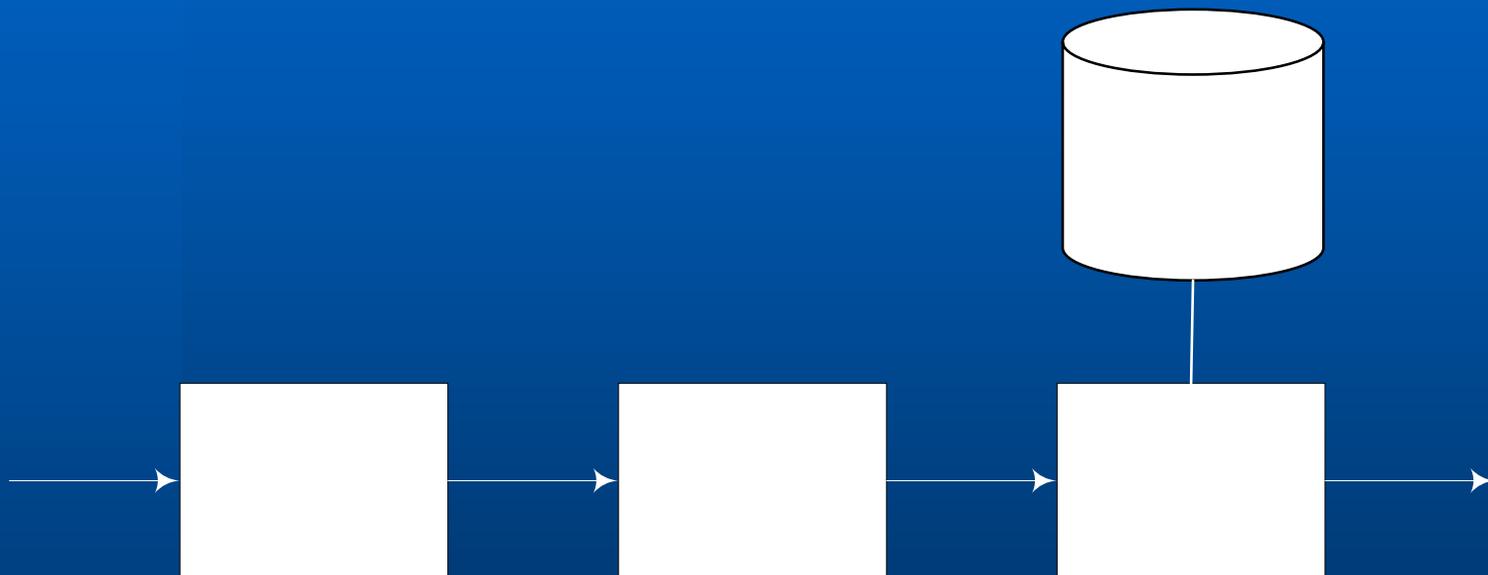


語者辨識流程



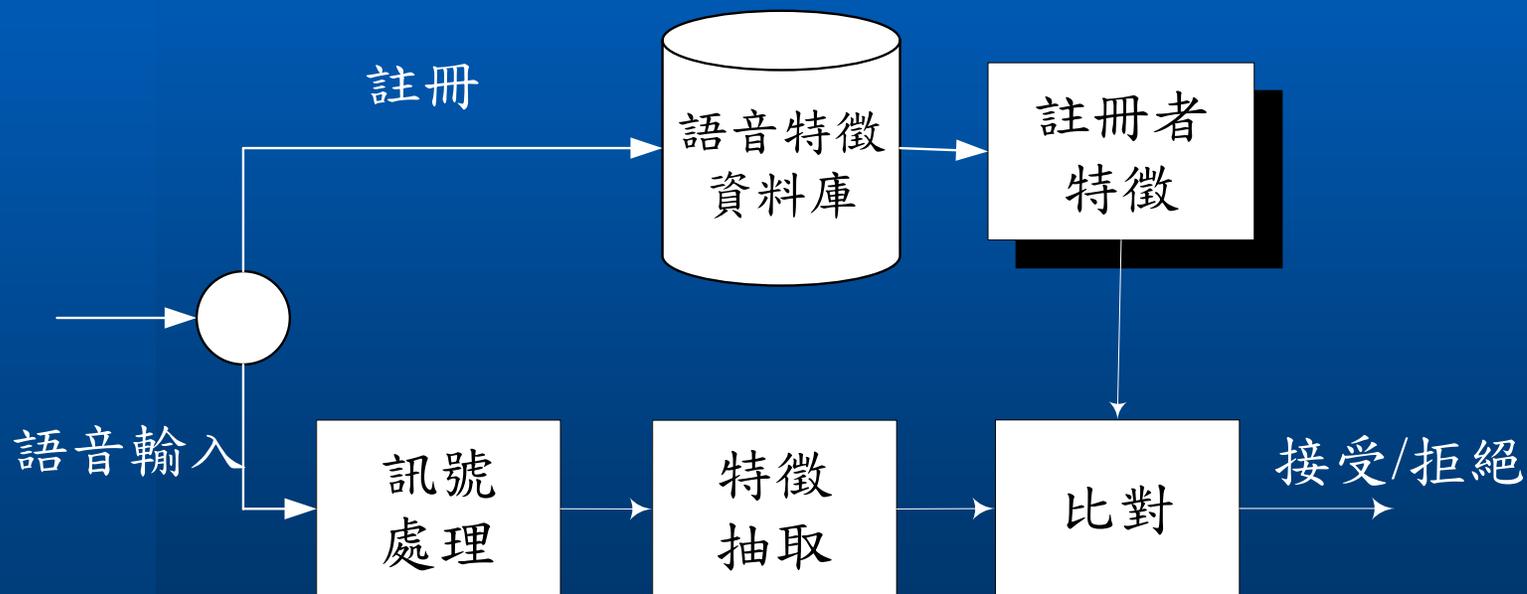


語者識別模式





語者驗證(比對)模式



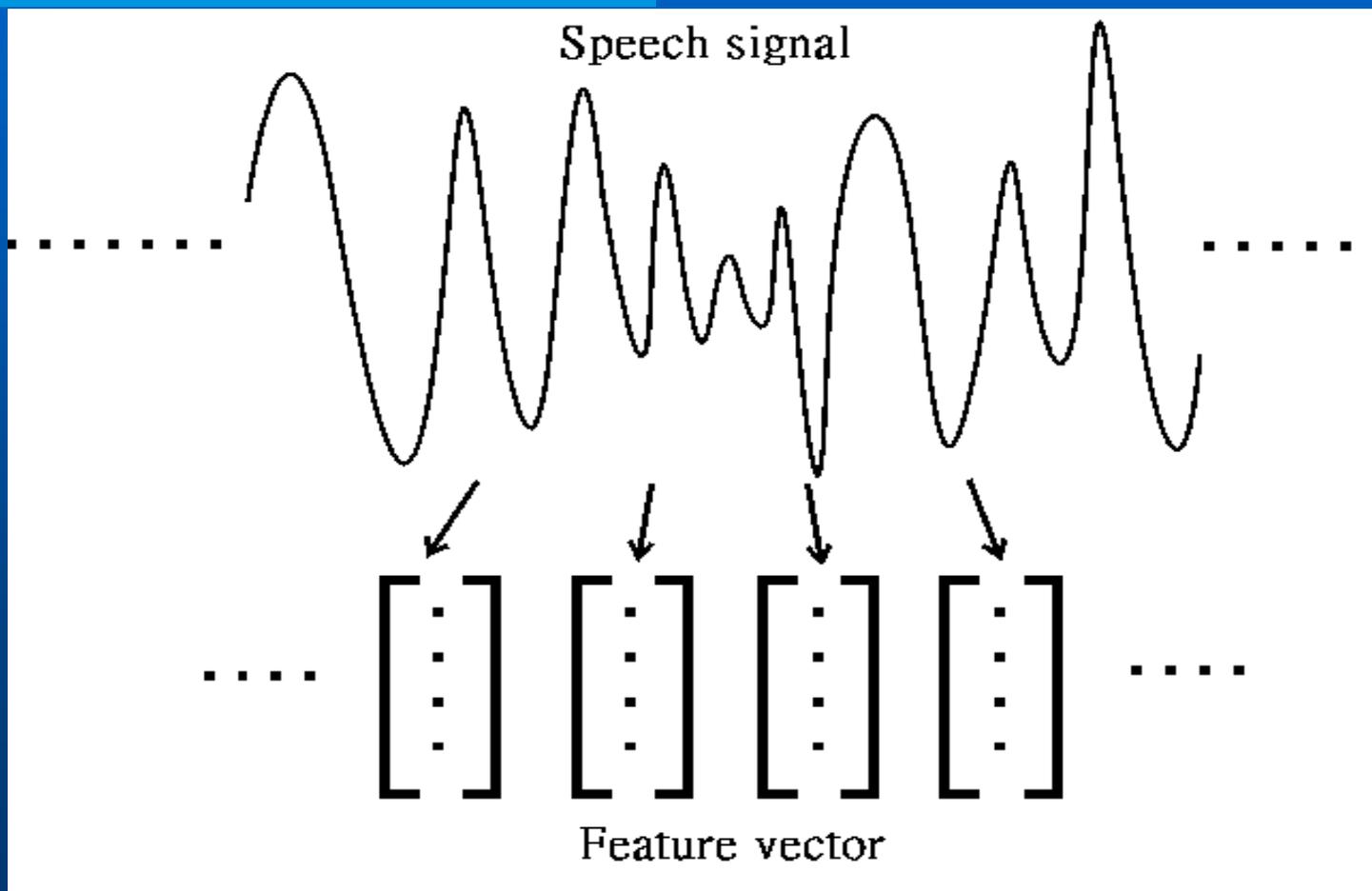


語音訊號處理及特徵抽取

- 語音訊號前處理
- 語音特徵抽取
 - 線性預測編碼導出的倒頻譜參數 (Short-Time Cepstral Coefficient - LPCC)
 - 梅爾刻度式倒頻譜參數 (Mel-Frequency Cepstral Coefficient - MFCC)



語音訊號之音框(Frame)





語音訊號前處理

- 漢明窗 (Hamming Window)

$$W(n) = \begin{cases} 0.54 - 0.46 \times \cos\left(\frac{2\pi n}{N-1}\right) & 0 \leq n \leq N-1 \\ 0 & \textit{otherwise} \end{cases}$$

- 正規化 (Normalization)
- 預強調 (Pre-Emphasis)

$$h(z) = 1 - az^{-1}$$



語音訊號特徵—LPC

1. 求相關係數：

$$R(k) = \sum_{n=-\infty}^{\infty} x(n)x(n+k) \quad k = 0, \dots, P+1$$

2. 求線性預測編碼 (LPC) 係數：

$$\alpha_i^{(P)} \quad 1 \leq i \leq P$$

Levinson-Durbin演算法

$$E^{(0)} = R(0)$$

for $i = 1 : P$

$$k_i = R(i) - \sum_{j=1}^{i-1} \alpha_j^{(i-1)} R(i-j)$$

$$\alpha_i^{(i)} = k_i / E^{(i-1)}$$

for $j = 1 : i-1$

$$\alpha_j^{(i)} = \alpha_j^{(i-1)} - k_i \alpha_{i-j}^{(i-1)}$$

end

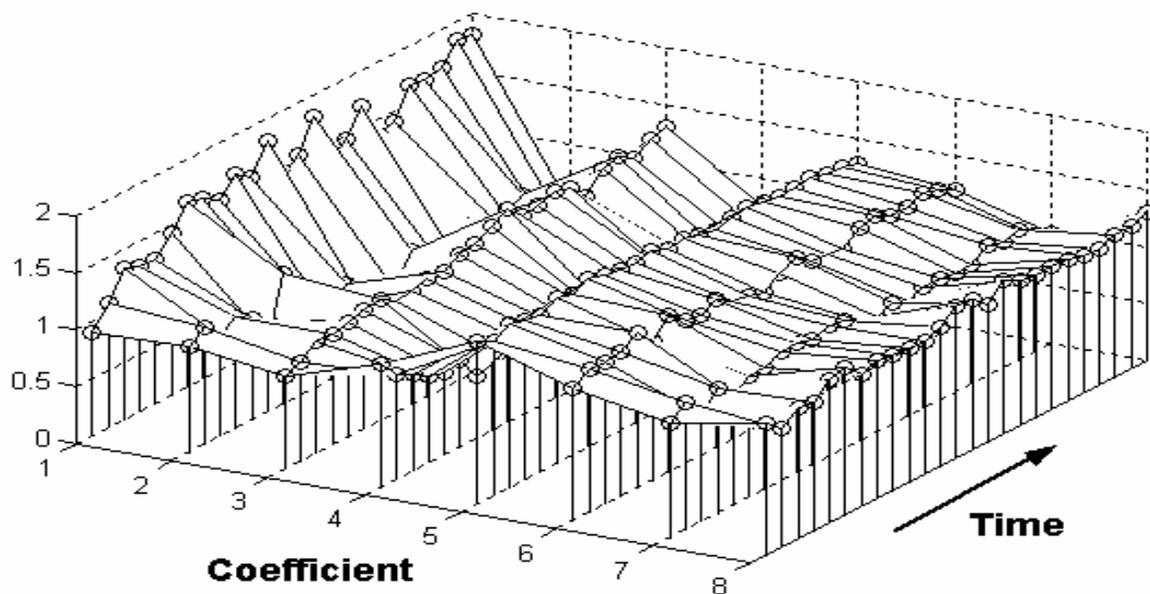
$$E^{(i)} = (1 - k_i^2) E^{(i-1)}$$

end



特徵抽取—LPCC

$$c_m = a_m + \sum_{k=1}^{m-1} \left(\frac{k}{m}\right) c_k a_{m-k} \quad 1 \leq m \leq P$$





LPCCC特徵向量

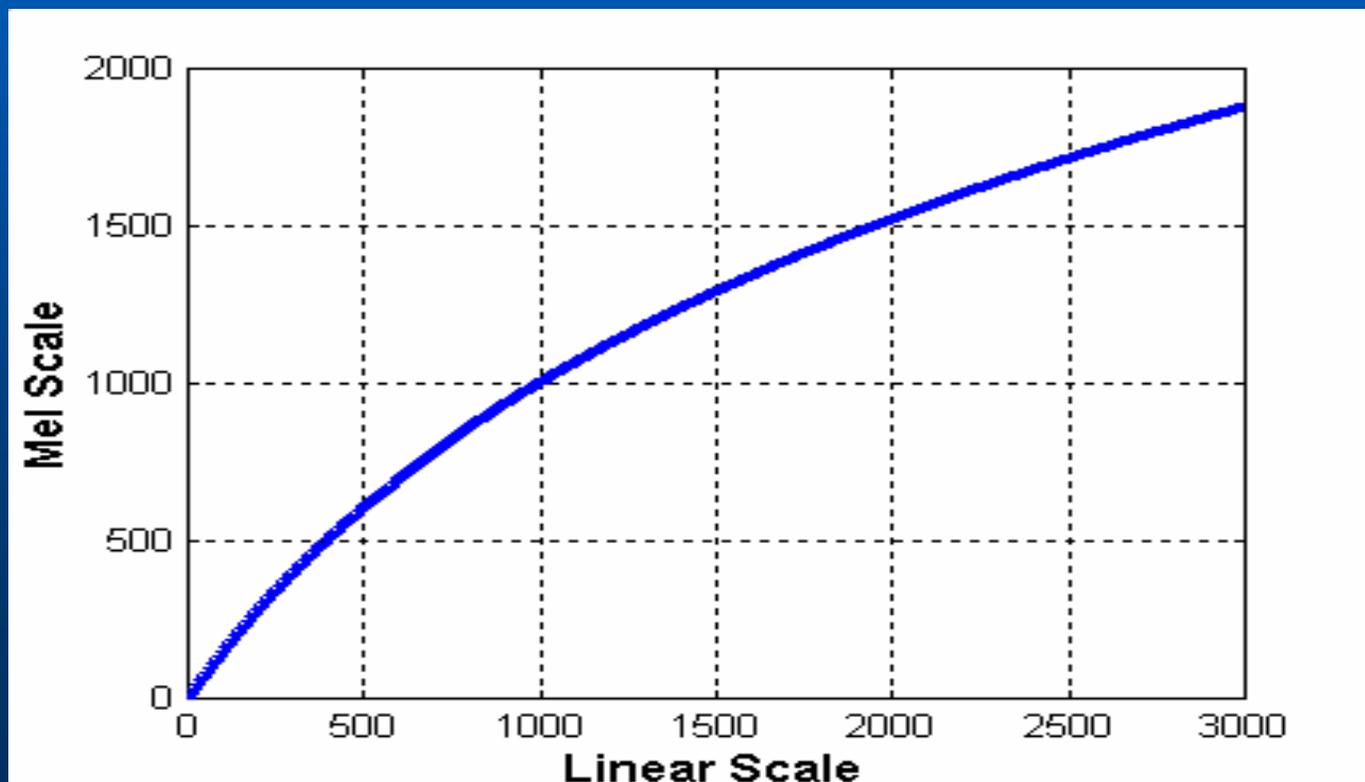
- LPCCC比LPC的強健性（Robustness）及可信賴度（Reliability）都要來的高。
- 缺點：將聲音在每段頻率上的特性視為線性，這與實際上人類的聽覺反應並不一致。



特徵抽取—MFCC

Mel Scale :

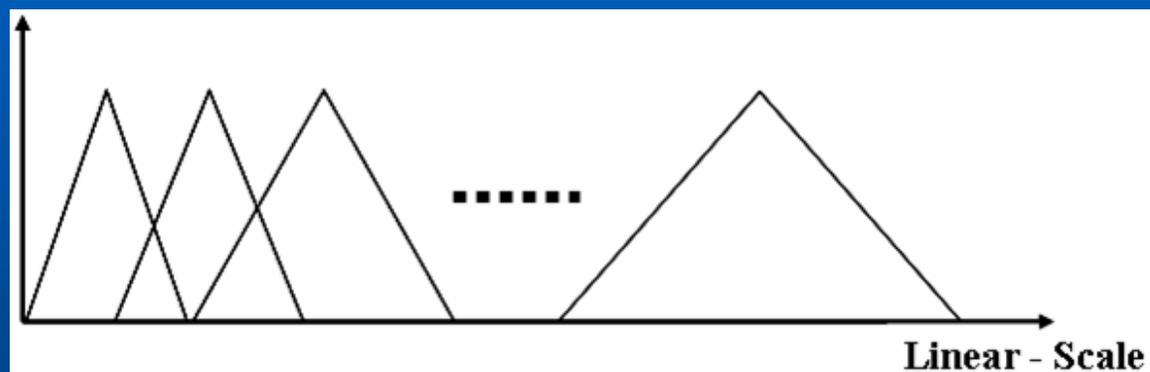
$$Mel(f) = 2595 \log_{10} \left(1 + \frac{f}{700} \right)$$



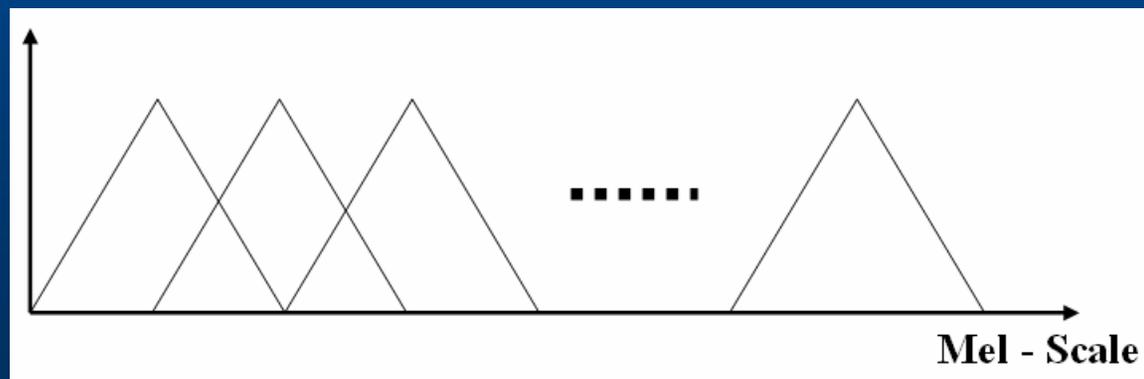


Linear-Scale和Mel-Scale

Linear - Scale



Mel - Scale





三角帶通濾波器

三角帶通濾波器的函數值：

$$U_m(n) = \begin{cases} 1 - \frac{|b_m - n|}{\Delta m} & \text{if } |b_m - n| \leq \Delta m \\ 0 & \text{if } |b_m - n| > \Delta m \end{cases}$$

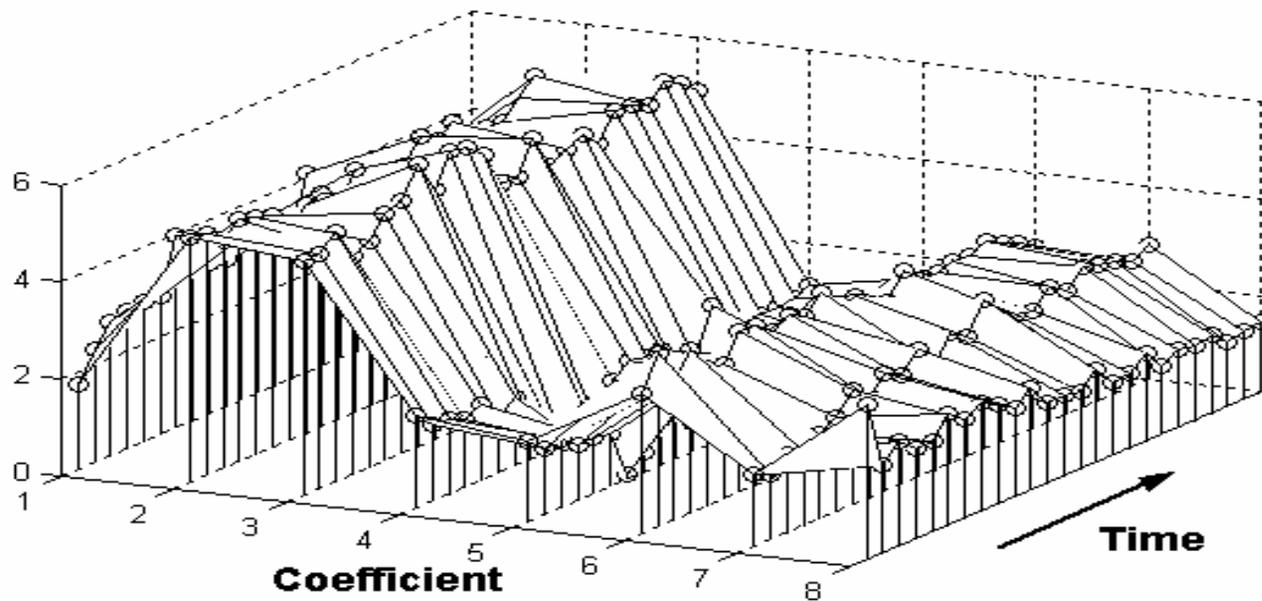
帶通濾波器之輸出：

$$Y(m) = \sum_{n=b_m - \Delta m}^{b_m + \Delta m} X(n)U_m(n - b_m)$$



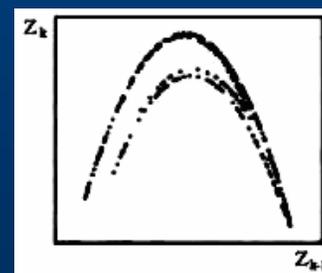
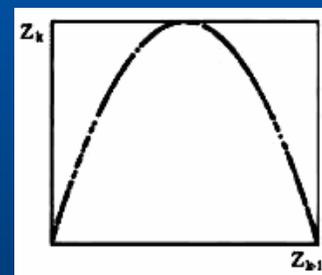
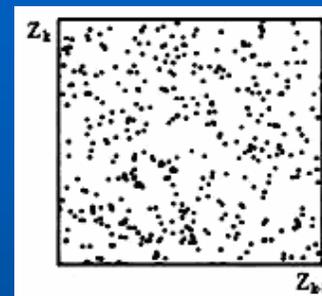
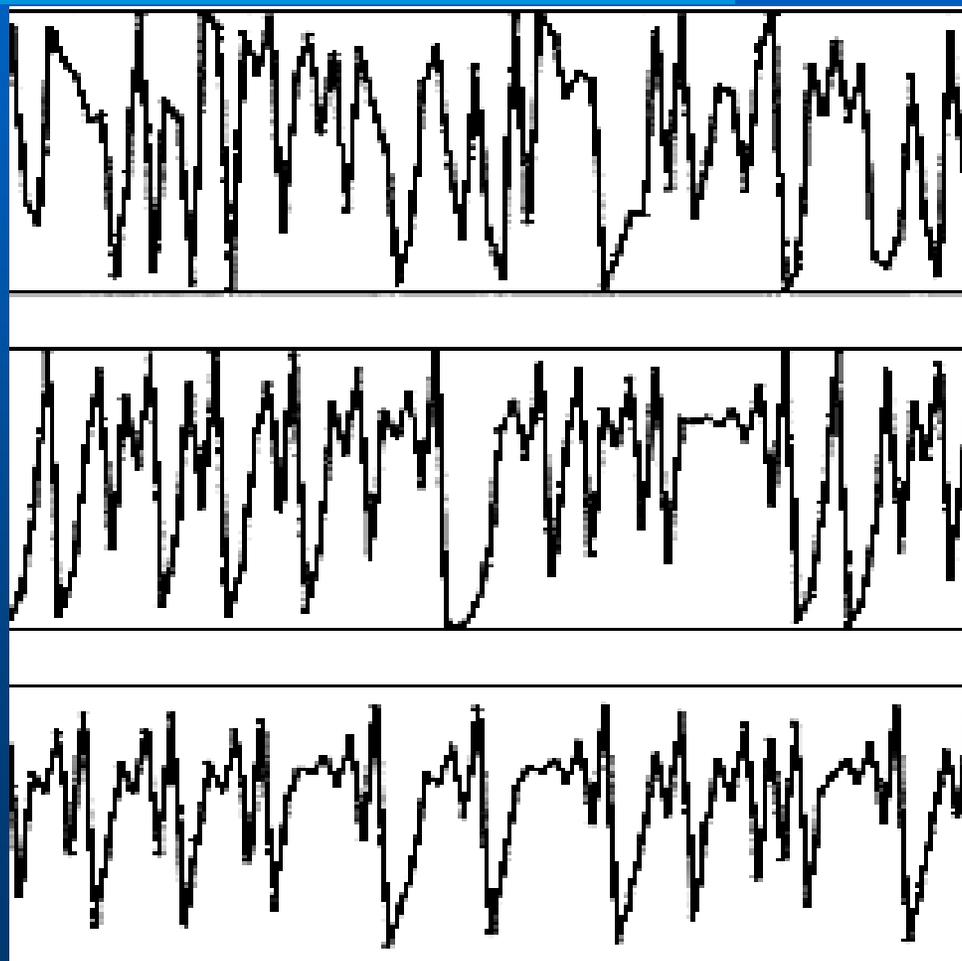
MFCC

$$C(k) = \sum_{m=1}^M \log\{|Y(m)|\} \cdot \cos\left\{k\left(m - \frac{1}{2}\right) \frac{\pi}{M}\right\}, \quad k = 1, 2, \dots, P$$





語音與混沌時間序列



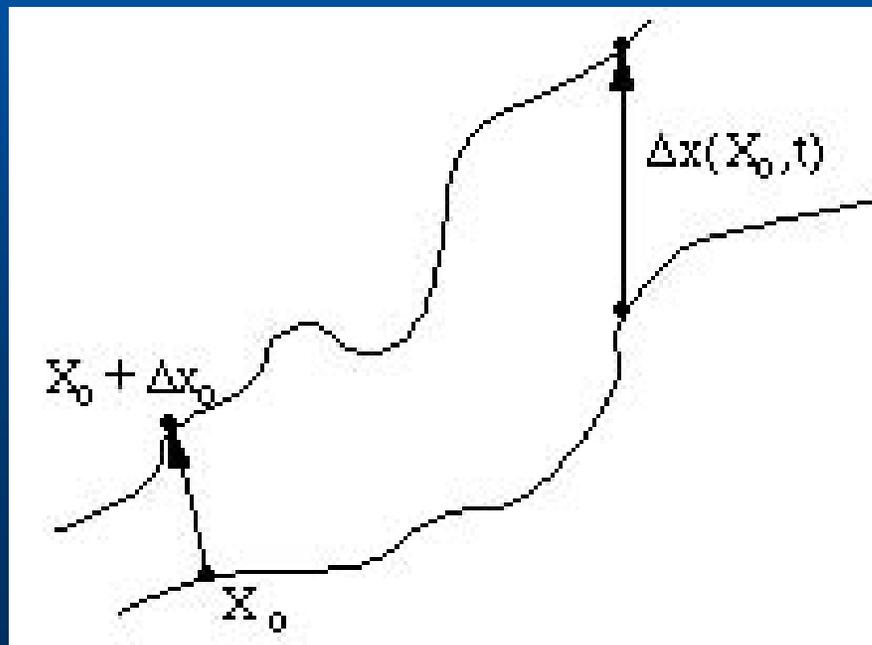


里雅普諾夫指數 (Liapunov exponents)

Benettin等人於1976年提出一個計算里雅普諾夫指數的方法：將動態系統特定時間狀態作為初始狀態，兩個軌跡相臨點間的距離為 $\Delta x_0(X_0)$ ，經過時間 t 後兩點距離變化為 $\Delta x(X_0, t)$

● 里雅普諾夫指數 λ ：

$$\lambda = \lim_{\substack{t \rightarrow \infty \\ |\Delta x_0| \rightarrow 0}} \frac{1}{t} \ln \left| \frac{\Delta x(X_0, t)}{\Delta x_0} \right|$$





里雅普諾夫指數演算法 (Liapunov exponents)

- 對一組離散的語音訊號 $x_0, x_1, x_2, \dots, x_n$ 中任一 x_i ($0 \leq i \leq n$) 搜尋一個與其值最接近的 x_j ($0 \leq j \leq n$)。

求出初始距離： $d_{ij} = |x_i - x_j|$

觀察經過時間 t 後的距離： $dt_{ij} = |x_{i+t} - x_{j+t}|$

- x_i 的里雅普諾夫指數：

$$\lambda_i = \frac{1}{t} \ln \left| \frac{dt_{ij}}{d_{ij}} \right|$$



PCA 語音特徵擷取

Mean Value :

$$m_v = \frac{1}{M} \sum_{k=1}^M X_k$$

Covariance Matrix :

$$E = (X - M)(X - M)'$$

PCA :

$$V = \left\{ \begin{bmatrix} \\ \\ \end{bmatrix}, \begin{bmatrix} \\ \\ \end{bmatrix}, \dots, \begin{bmatrix} \\ \\ \end{bmatrix} \right\}$$



碎形語音特徵擷取



Iterated Function System(IFS)

$$a_i x_0 + e_i = x_{i-1}$$

$$a_i x_N + e_i = x_i$$

$$c_i x_0 + d_i F_0 + F_i = F_{i-1}$$

$$c_i x_N + d_i F_N + f_i = F_i$$

a_i 、 c_i 、 d_i 、 e_i 和 f_i 是IFS係數



碎形語音特徵擷取



$$x_i^* = [c_i \quad d_i \quad f_i]$$

$$A^* = [x_1^* \quad x_2^* \quad \cdots \quad x_N^*]^T$$

$$m_j = \frac{1}{N} \sum_{i=1}^N x_{ij}^*$$

$$S_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{ij}^* - m_j)^2$$



小波語音特徵擷取

1. 計算每一層小波轉換後的語音訊號係數的平均能量：

$$v_i = \frac{1}{n_i} \sum_{j=1}^{n_i} w_{i,j}^2 \quad i = 1, 2, \dots, N + 1$$

2. 將每一層 i 平均能量當小波轉換的語音訊號特徵值：

$$V = \{v_1, v_2, \dots, v_i\}^t$$



期末報告專題

- Ear Biometrics
- Lips Biometrics(鄧)
- Hand Geometry Biometrics(邱)
- Palmprint Biometrics
- Gait Biometrics
- Handwriting(Signature) Biometrics(黃)
- Keystroke Biometrics(謝昇憲)
- Voice Biometrics(陳俊任)
- Multimodal Biometrics(朱家德)