

# 數位語音處理概論

## Introduction to Digital Speech Processing

## 2.0 Fundamentals of Speech Recognition

### References for 2.0

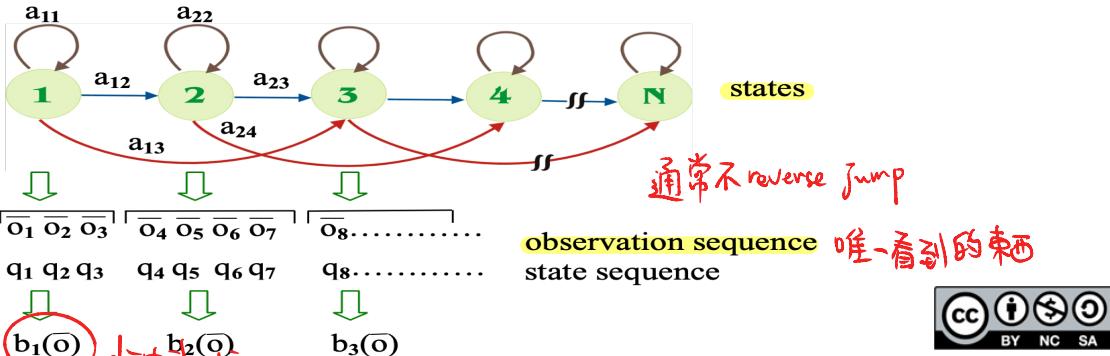
1.3, 3.3, 3.4, 4.2, 4.3, 6.4, 7.2, 7.3, of Bechetti

授課教師：國立臺灣大學 電機工程學系 李琳山 教授



【本著作除另有註明外，採取創用CC「姓名標示  
—非商業性—相同方式分享」臺灣3.0版授權釋出】

## Hidden Markov Models (HMM)



### Formulation

$$\bar{O}_t = [x_1, x_2, \dots, x_D]^T \quad \text{feature vectors for a frame at time } t$$

$$q_t \in \{1, 2, 3, \dots, N\} \quad \text{state number for feature vector } \bar{O}_t$$

$$A = [a_{ij}] , \quad a_{ij} = \text{Prob}[q_t = j | q_{t-1} = i] \quad \text{e.g. } a_{11} = \text{state } 1 \rightarrow 1$$

*參照到  $a_{ij}, a_{iz}$*

$$B = [b_j(\bar{O}), j = 1, 2, \dots, N] \quad \text{observation (emission) probability}$$

$$b_j(\bar{O}) = \sum_{k=1}^M c_{jk} b_{jk}(\bar{O}) \quad \text{Gaussian Mixture Model (GMM)}$$

*j-th state      k-th gaussian*

$b_{jk}(\bar{O})$ : multi-variate Gaussian distribution

for the k-th mixture (Gaussian) of the j-th state

M : total number of mixtures

$$\sum_{k=1}^M c_{jk} = 1 \quad \text{weight of gaussian}$$

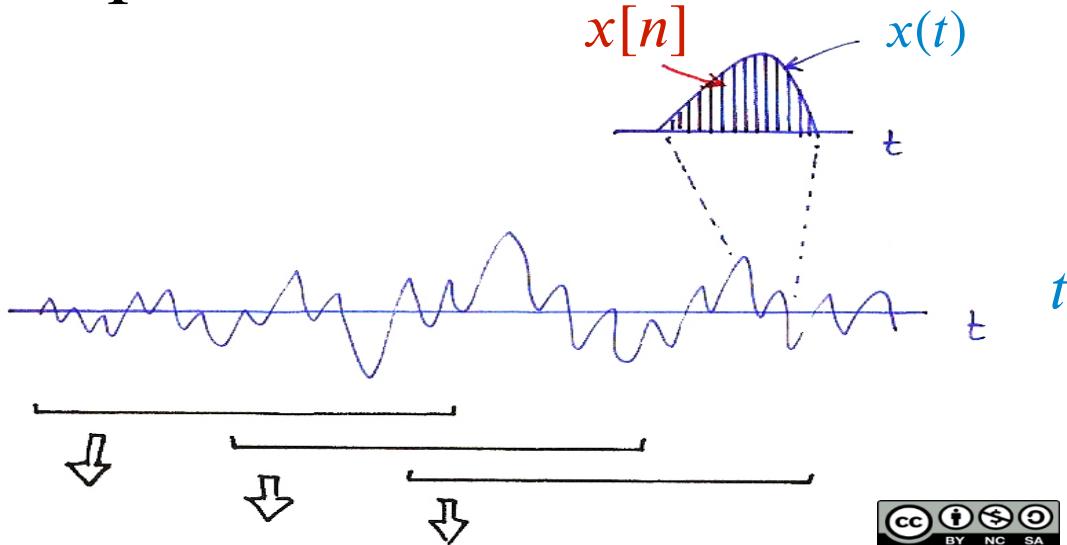
$$\pi = [\pi_1, \pi_2, \dots, \pi_N] \quad \text{initial probabilities}$$

$\pi_i = \text{Prob}[q_1 = i]$  自 i 開始之 probability

$$\text{HMM} : (A, B, \pi) = \lambda$$

e.g. 上 shang → set as initial  
 s - sh - shan - shang

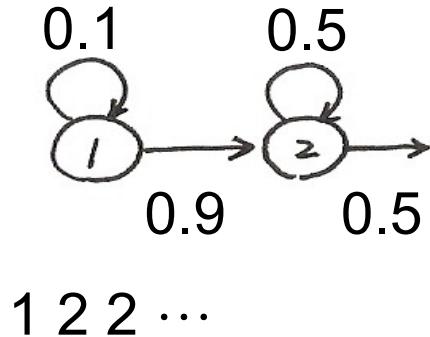
# Observation Sequences



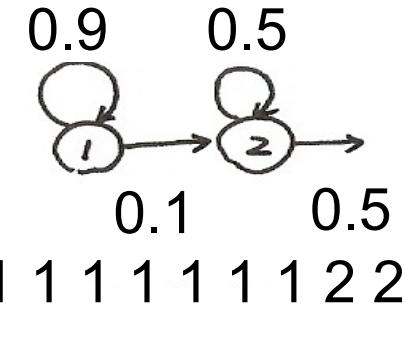
$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \stackrel{\text{observation}}{=} o_1 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = o_2 \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = o_3$$

$D \cong 39$

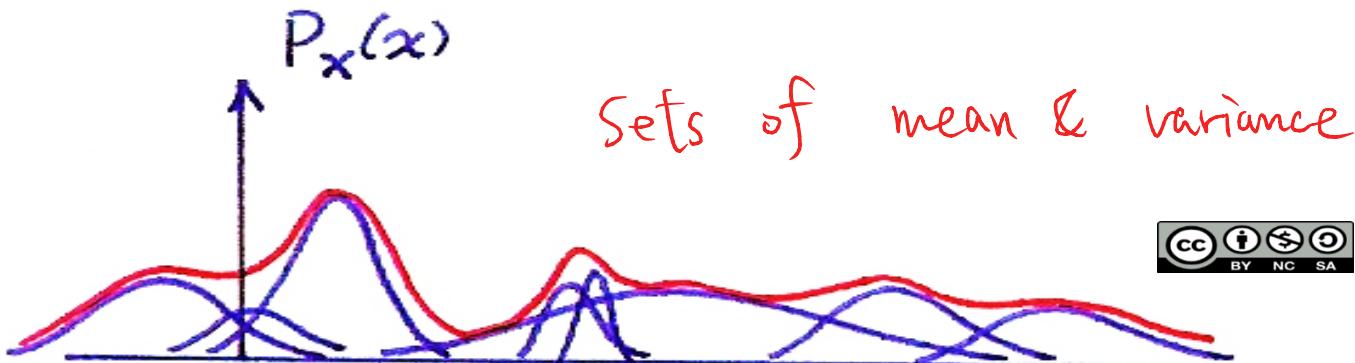
# State Transition Probabilities



GMM = Gaussian Mixture Model



1-dim Gaussian Mixtures probability model  $\Rightarrow$  use gaussian to approximate



## • Gaussian Random Variable X

$$f_X(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-(x-\mu)^2/2\sigma^2}$$

↗ value of random variable X  
 Random vector  
 ↑

## • Multivariate Gaussian Distribution for n Random Variables

$$\bar{X} = [ X_1, X_2, \dots, X_n ]^t \quad n\text{-dimension}$$

$$f_{\bar{X}}(\bar{x}) = \frac{1}{(2\pi)^{n/2} \Delta^{1/2}} e^{-\frac{1}{2}[(x-\bar{\mu})^t \Sigma^{-1}(x-\bar{\mu})]}$$

$$\bar{\mu} = [ \mu_{X_1}, \mu_{X_2}, \dots, \mu_{X_n} ]^t$$

$$\Sigma = [ \sigma_{ij} ], \text{ covariance matrix}$$

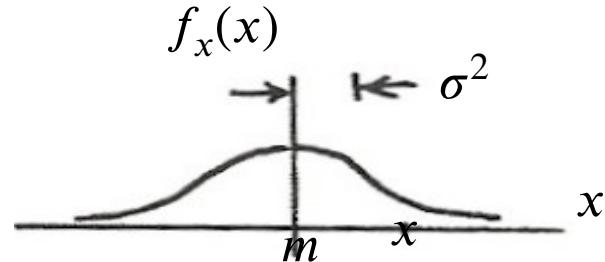
$$\Sigma = [ \sigma_{ij} ] :$$

$$\sigma_{ij} = E [ (X_i - \mu_{X_i})(X_j - \mu_{X_j}) ]$$

$\Delta$  : determinant of  $\Sigma$

# Multivariate Gaussian Distribution

$$(\bar{x} - \bar{\mu})^T \Sigma^{-1} (\bar{x} - \bar{\mu}) = \left( \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} \right)^T \Sigma^{-1} \left( \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} \right)$$



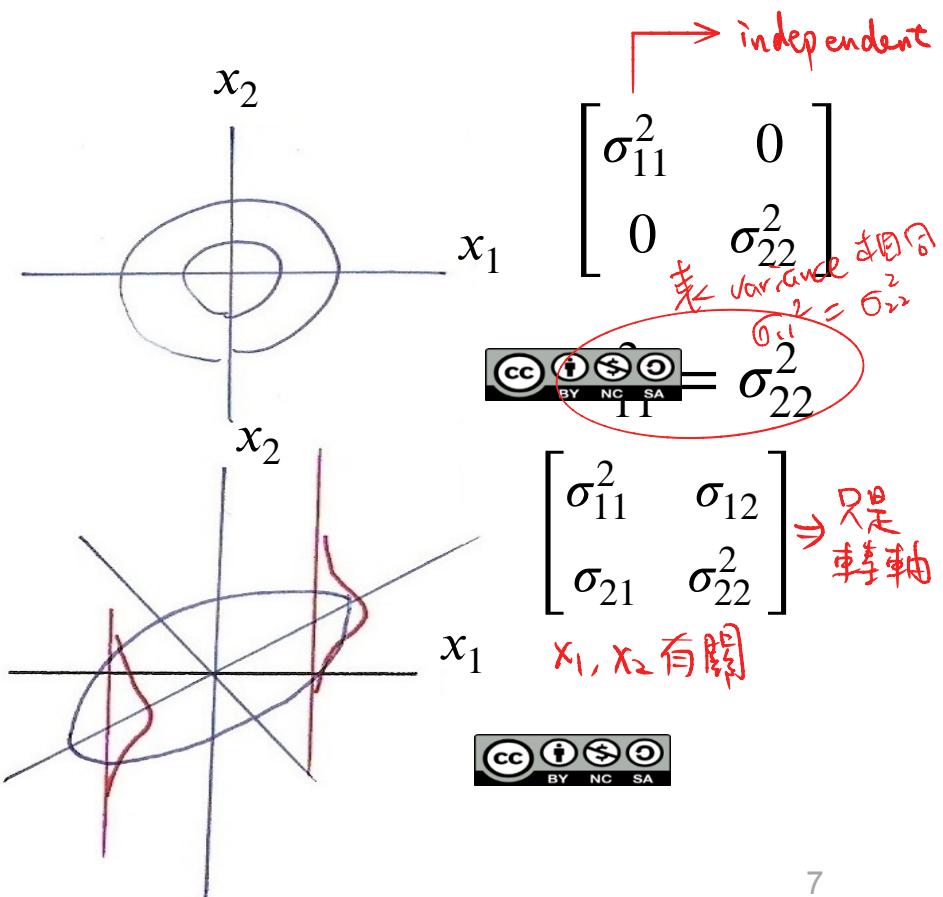
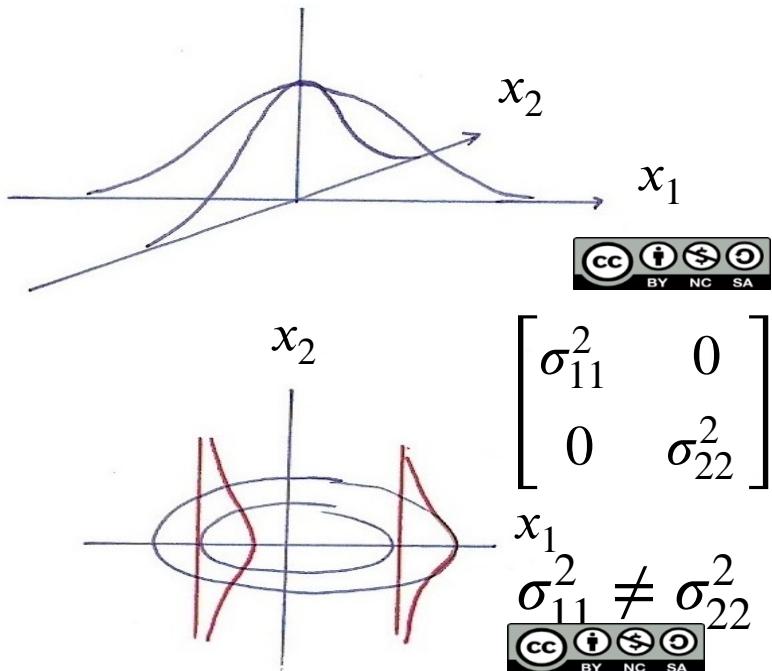
$$\begin{aligned}
 &= [x_1 - \mu_1 \ x_2 - \mu_2 \ \dots \ x_n - \mu_n] \underbrace{\Sigma^{-1}}_{\substack{\text{matrix} \\ \text{inverse}}} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ \vdots \\ x_n - \mu_n \end{bmatrix} \\
 &= (x_1 - \mu_1)^2 + (x_2 - \mu_2)^2 + \dots, \quad \text{if } \Sigma = \begin{bmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \\
 &= \frac{(x_1 - \mu_1)^2}{\sigma_{11}^2} + \frac{(x_2 - \mu_2)^2}{\sigma_{22}^2} + \dots, \quad \text{if } \Sigma = \begin{bmatrix} \sigma_{11}^2 & & 0 \\ & \sigma_{22}^2 & \\ 0 & & \ddots \end{bmatrix} \\
 &\text{normalize} \quad \Rightarrow \text{calculate the norm of probability}
 \end{aligned}$$

$$\Sigma = \begin{bmatrix} & j & \\ & \vdots & \\ \cdots & \sigma_{ij} & \cdots \\ & \vdots & \end{bmatrix} i$$

與 mean 差的遠

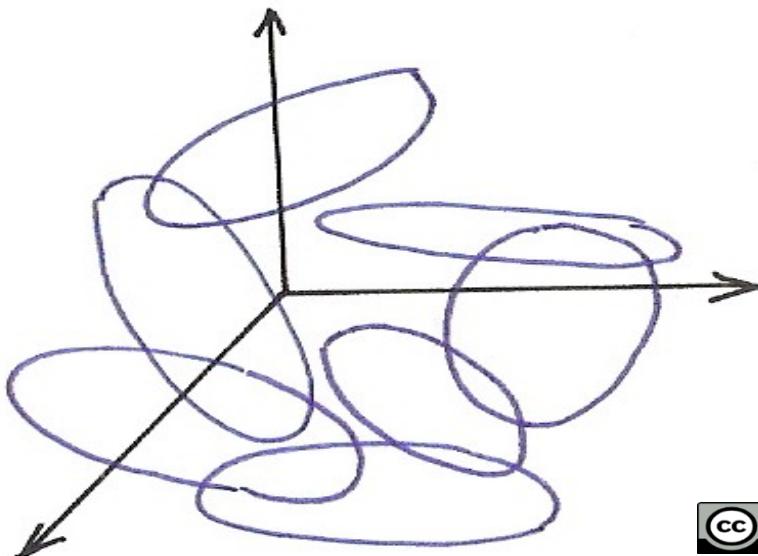
$$\sigma_{ij} = E[(x_i - \mu_{x_i})(x_j - \mu_{x_j})]$$

# 2-dim Gaussian



# N-dim Gaussian Mixtures

N-dim Gaussian Mixtures



# Hidden Markov Models (HMM) *modelize acoustic signal*

## Double Layers of Stochastic Processes

*每人音長都不同 e.g. - 長 - ~~~*

- hidden states with random transitions for time warping
- random output given state for random acoustic characteristics

## Three Basic Problems

### (1) Evaluation Problem:

Given  $\overline{O} = (\overline{o_1}, \overline{o_2}, \dots, \overline{o_t}, \dots, \overline{o_T})$  and  $\lambda = (A, B, \pi)$   
find Prob [  $\overline{O}$  |  $\lambda$  ]

### (2) Decoding Problem:

Given  $\overline{O} = (\overline{o_1}, \overline{o_2}, \dots, \overline{o_t}, \dots, \overline{o_T})$  and  $\lambda = (A, B, \pi)$   
find a best state sequence  $\overline{q} = (q_1, q_2, \dots, q_t, \dots, q_T)$

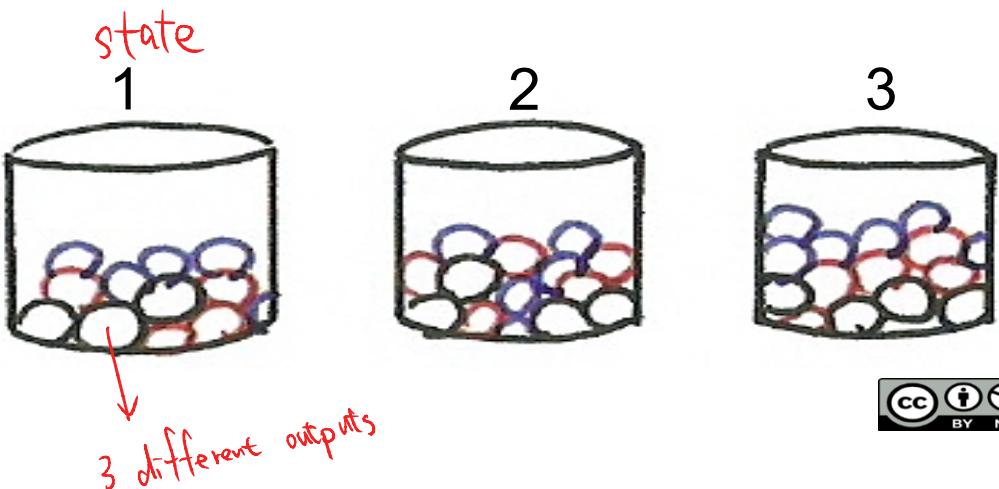
### (3) Learning Problem:

Given  $\overline{O}$ , find best values for parameters in  $\lambda$   
such that Prob [  $\overline{O}$  |  $\lambda$  ] = max

# Simplified HMM

RGBGGBBGRRR.....

⇒ hidden random output  
from hidden random state



拉高頻率  
去雜訊

## Feature Extraction (Front-end Signal Processing)

### Pre-emphasis

$$H(z) = 1 - az^{-1}, \quad 0 < a < 1$$

$$x[n] = x'[n] - ax'[n-1]$$

算對不同頻率之反應

- pre-emphasis of spectrum at higher frequencies

### Endpoint Detection (Speech/Silence Discrimination)

- short-time energy

$$E_n = \sum_{m=-\infty}^{\infty} (x[m])^2 w[m-n]$$

背景雜訊

- adaptive thresholds

⇒ 希望 machine 不要把 silence 作人聲

easiest concept: 算 Energy

### Windowing

$$Q_n = \sum_{m=-\infty}^{\infty} T\{x[m]\} w[m-n]$$

$T\{\cdot\}$  : some operator

$w[m]$  : window shape

- Rectangular window

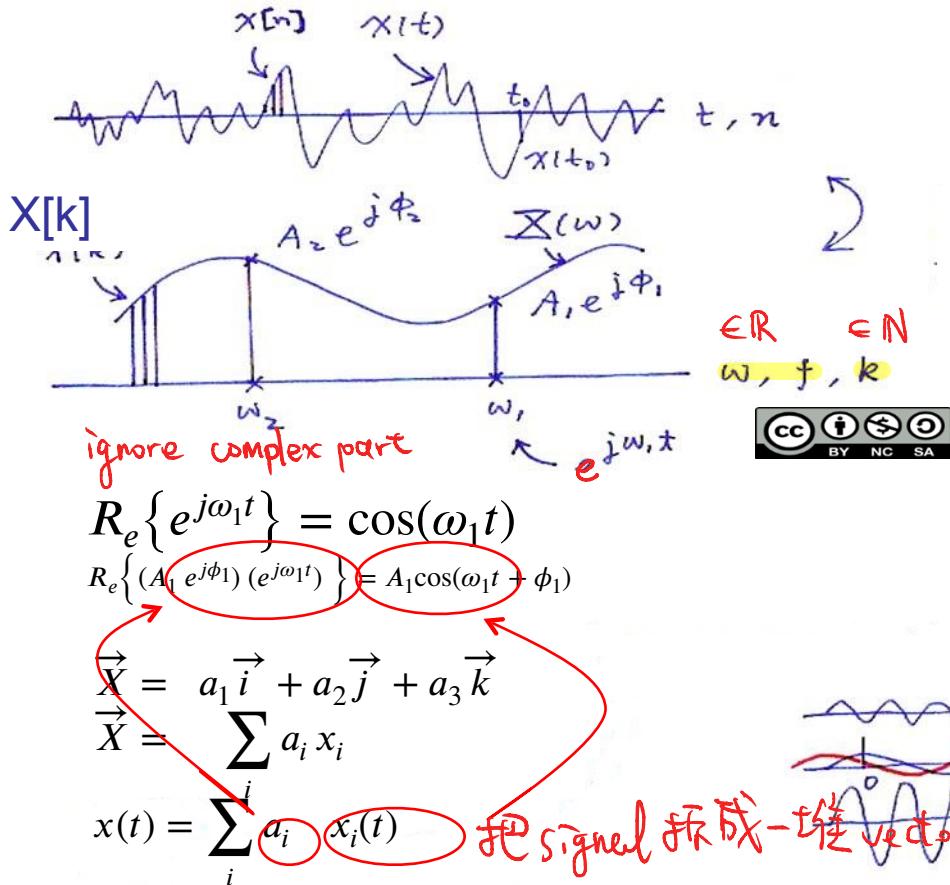
$$w[m] = \begin{cases} 1, & 0 < m \leq L-1 \\ 0, & \text{else} \end{cases}$$

Hamming window

$$w[m] = \begin{cases} 0.54 - 0.46 \cos\left[\frac{2\pi m}{L}\right], & 0 \leq m \leq L-1 \\ 0, & \text{else} \end{cases}$$

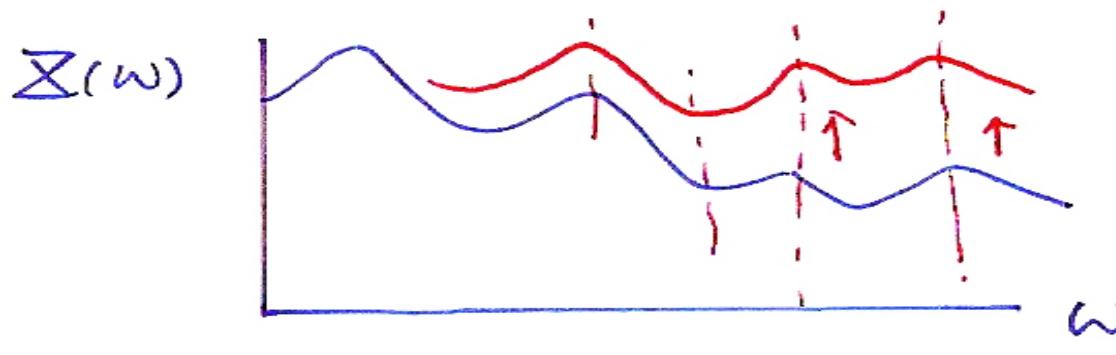
window length/shift/shape ← 可選擇

# Time and Frequency Domains



# Pre-emphasis

- 很多音在高頻較弱,  
 : 才放大



- Pre-emphasis 放大高頻

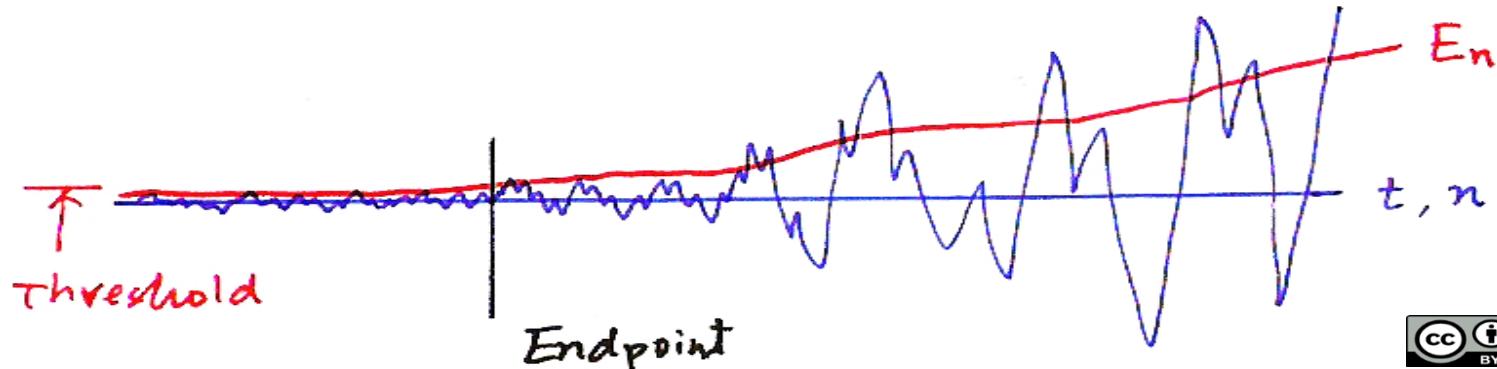
$$H(z) = 1 - az^{-1}, \quad 0 < a < 1$$

$$x[n] = x'[n] - ax'[n-1]$$

*Z-transform*

pre-emphasis of spectrum at higher frequencies

# Endpoint Detection



- **Endpoint Detection (Speech/Silence Discrimination)**

- short-time energy

$$E_n = \sum_{m=-\infty}^{\infty} (x[m])^2 w[m-n]$$

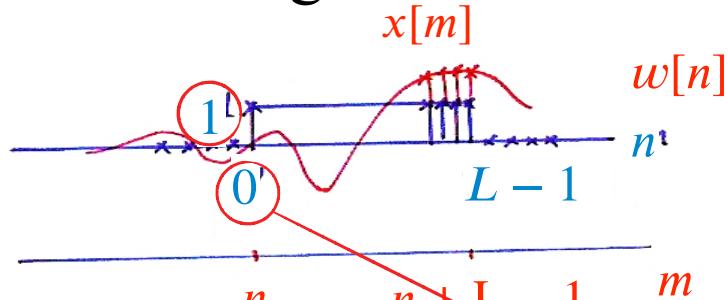
*window*

- adaptive thresholds

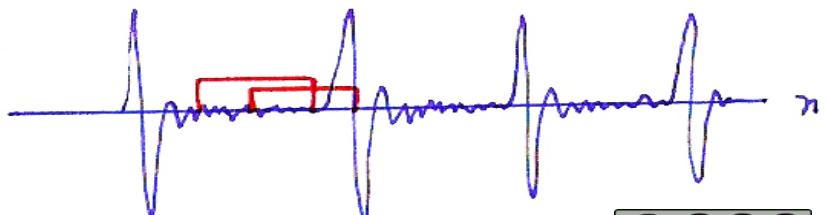
效果: hamming  $\Rightarrow$  rectangular

# Endpoint Detection

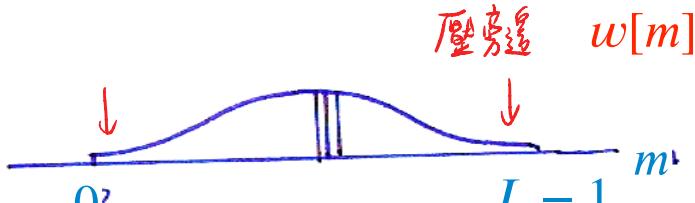
## ● Rectangular Window



$$E_n = \sum_{m=-\infty}^{\infty} (x[m])^2 w[m - n]$$



## ● Hamming Window



Hamming window

$$w[m] = \begin{cases} 0.54 - 0.46\cos\left[\frac{2\pi m}{L}\right], & 0 \leq m \leq L - 1 \\ 0, & \text{else} \end{cases}$$

$$Q_n = \sum_{m=-\infty}^{\infty} T\{x[m]\} w[m - n]$$

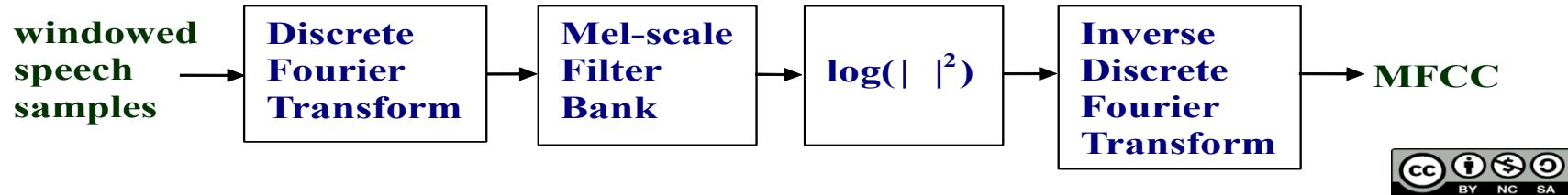
一个段一小段取，才不会变动太大

$T\{\bullet\}$  : some operator

$w\{m\}$  : window shape

# Feature Extraction (Front-end Signal Processing)

## Mel Frequency Cepstral Coefficients (MFCC)



### - Mel-scale Filter Bank

triangular shape in frequency/overlapped  
uniformly spaced below 1 kHz

logarithmic scale above 1 kHz 在高頻會失真且高頻, 2倍 freq, 會一直

維持高

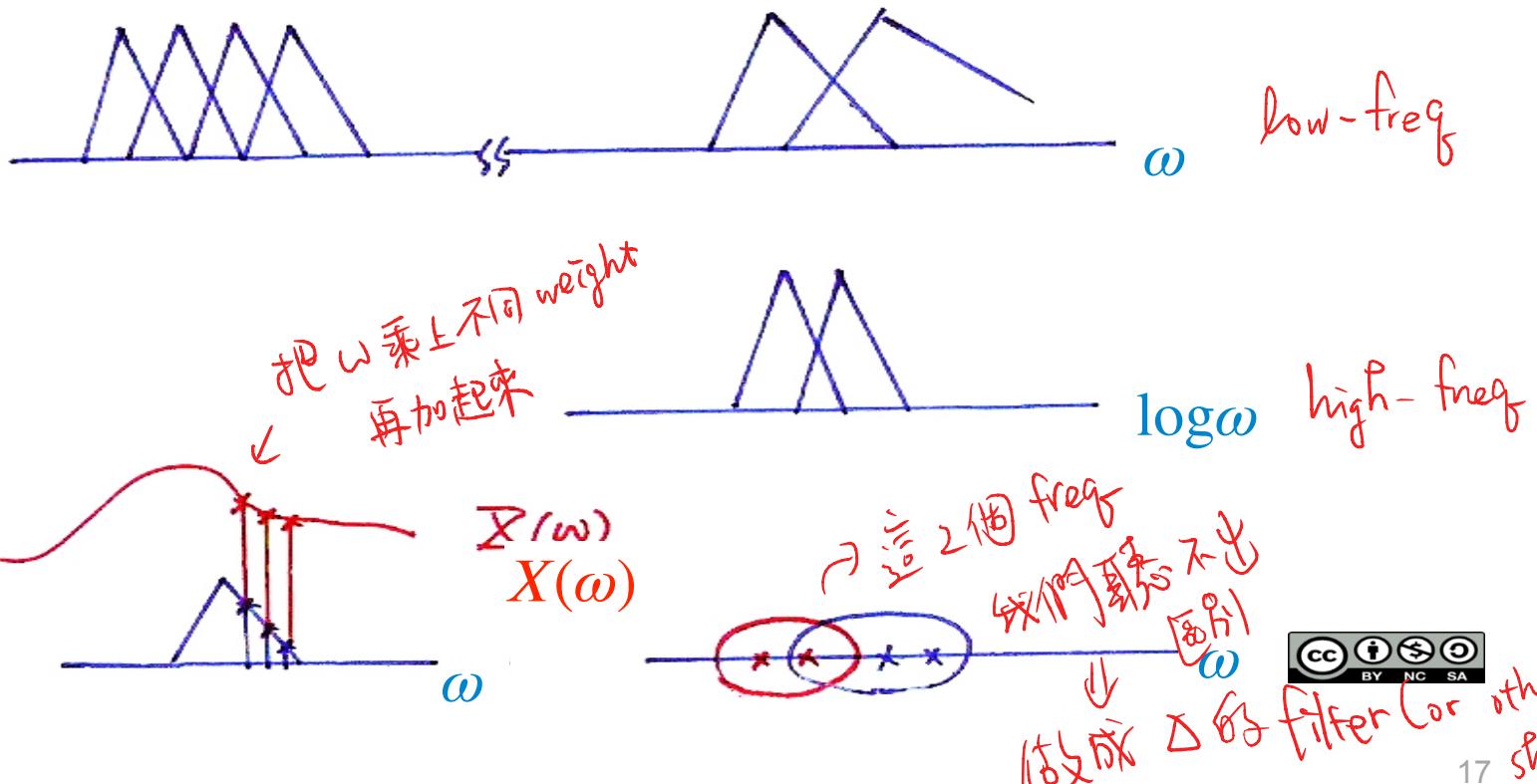
∴ 高頻再取 log

## Delta Coefficients

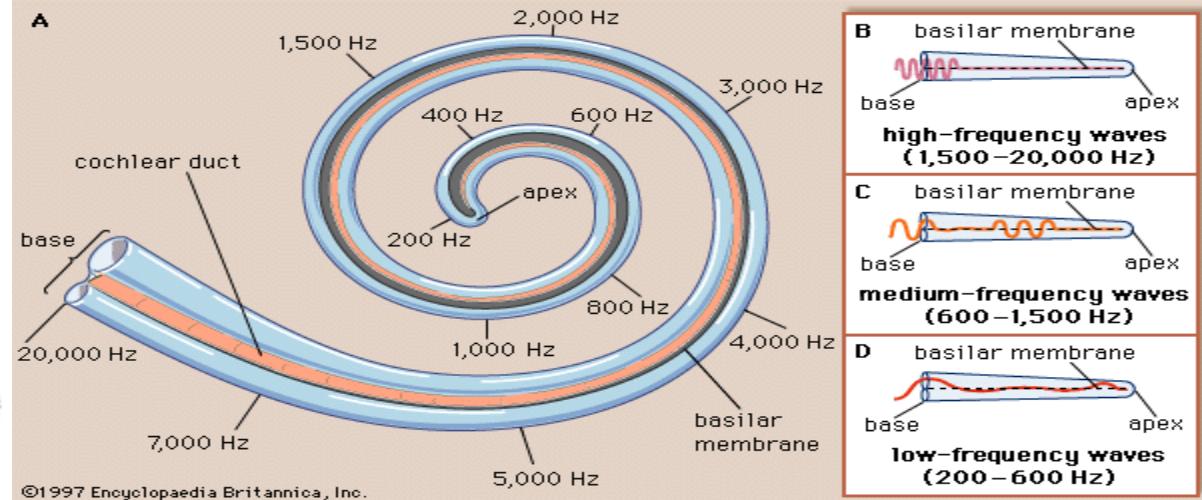
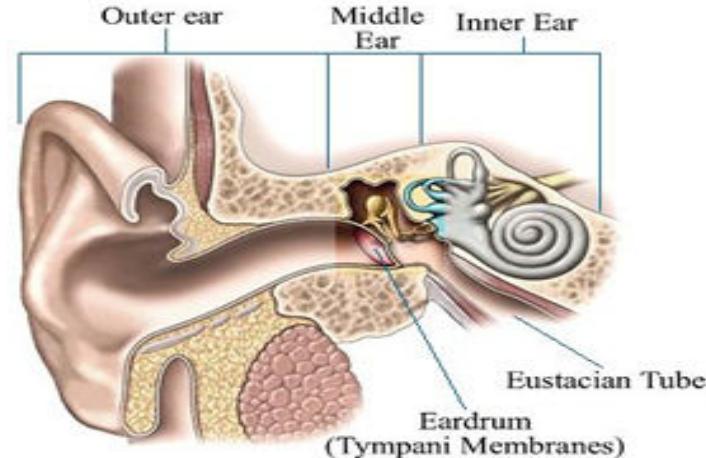
- 1st/2nd order differences

14057

# Mel-scale Filter Bank

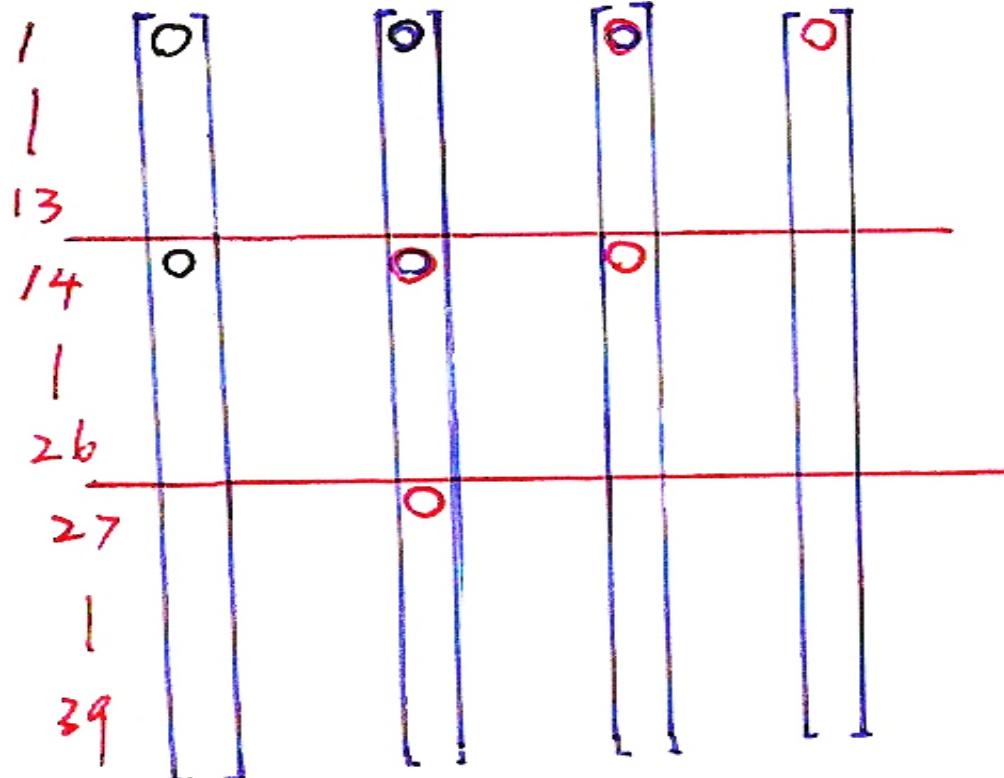


# Peripheral Processing for Human Perception (P.34 of 7.0 )



# Delta Coefficients

做 differential 找 delta



## Language Modeling: N-gram

看詞→句的構造(寫句)  
 $W = (w_1, w_2, w_3, \dots, w_i, \dots, w_R)$  a word sequence  
 詞 = character

- Evaluation of  $P(W)$

$$P(W) = P(w_1) \prod_{i=2}^R P(w_i | w_1, w_2, \dots, w_{i-1})$$

- Assumption:

$$P(w_i | w_1, w_2, \dots, w_{i-1}) = P(w_i | w_{i-N+1}, w_{i-N+2}, \dots, w_{i-1})$$

Occurrence of a word depends on previous  $N-1$  words only

N-gram language models

$$N = 2 : \text{bigram} \quad P(w_i | w_{i-1})$$

$$N = 3 : \text{tri-gram} \quad P(w_i | w_{i-2}, w_{i-1})$$

$$N = 4 : \text{four-gram} \quad P(w_i | w_{i-3}, w_{i-2}, w_{i-1})$$

⋮

$$N = 1 : \text{unigram} \quad P(w_i)$$

probabilities estimated from a training text database

example : tri-gram model

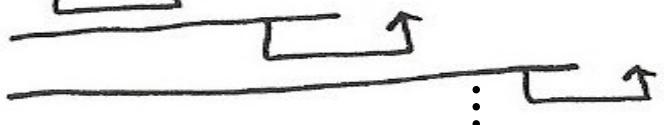
$$P(W) = P(w_1) P(w_2 | w_1) \prod_{i=3}^N P(w_i | w_{i-2}, w_{i-1})$$

# N-gram

20330

$$P(W) = P(w_1) \prod_{i=2}^R P(w_i | w_1, w_2, \dots, w_{i-1})$$

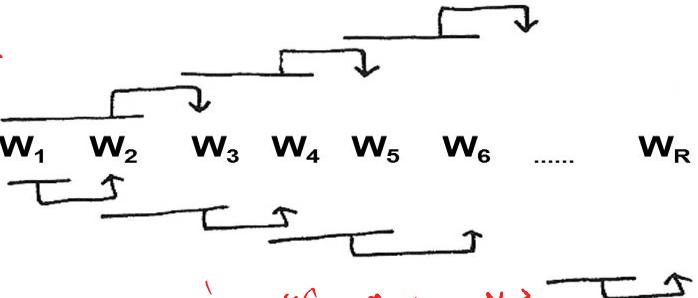
$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad \dots \quad w_R$



有 1, 2, ... 代表會到下一個  
的機率



## ◎ tri-grammer



$$P(W) = P(w_1) P(w_2 | w_1) \prod_{i=3}^N P(w_i | w_{i-2}, w_{i-1})$$

只算3連的機率  
前2個看到下一個 prob

因為太複雜，但其實不太通，只是方便有效



# Language Modeling

- Evaluation of N-gram model parameters

unigram

$$P(w^i) = \frac{N(w^i)}{\sum_{j=1}^V N(w^j)}$$

$w^i$  : a word in the vocabulary

$V$  : total number of different words in the vocabulary

$N(\cdot)$  number of counts in the training text database

bigram

$$P(w^j|w^k) = \frac{N(w^k, w^j)}{N(w^k)}$$

$\langle w^k, w^j \rangle$  : a word pair

trigram

$$P(w^j|w^k, w^m) = \frac{N(w^k, w^m, w^j)}{N(w^k, w^m)}$$

**smoothing** – estimation of probabilities of rare events by statistical approaches

... this .....	50000
... this is .....	500
... this is a ...	5

Prob [ is| this ] =

500

50000

bigram

Prob [ a| this is ] =

5

500

$$P(w^j | w^k) = \frac{N(\langle w^k, w^j \rangle)}{N(w^k)}$$

$\langle w^k, w^j \rangle$ : a word pair

trigram

$$P(w^j | w^k, w^m) = \frac{N(\langle w^k, w^m, w^j \rangle)}{N(\langle w^k, w^m \rangle)}$$

# Large Vocabulary Continuous Speech Recognition

~~$\bar{W} = (w_1, w_2, \dots, w_R)$~~  a word sequence  
 ~~$\bar{x} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_T)$~~  feature vectors for a speech utterance  
 $\bar{W}^* = \underset{\bar{W}}{\text{Arg Max}} \text{Prob}(\bar{W}|\bar{x})$  MAP principle given  $\bar{x}$ , find  $\bar{W}^*$ 's prob

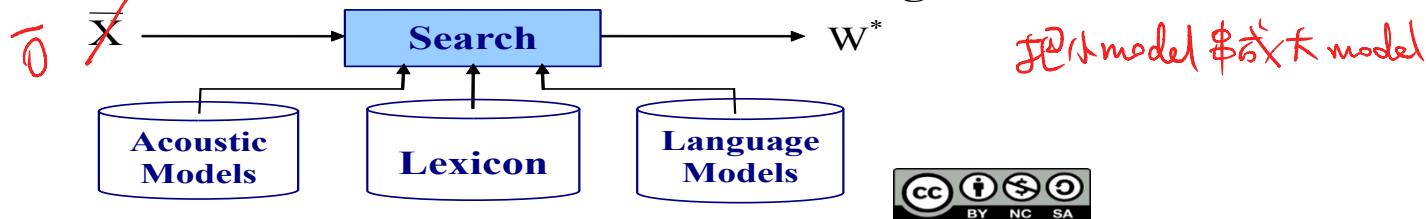
$$\text{Prob}(\bar{W}|\bar{x}) = \frac{\text{Prob}(\bar{x}|W) \cdot P(W)}{P(\bar{x})} = \max$$

A Posteriori Probability ( $= W^*$ )  
 $\Rightarrow$  Maximum A Posteriori (MAP) Principle

$$\text{Prob}(\bar{x}|W) \cdot P(W) = \max$$

$\uparrow$  by HMM       $\uparrow$  by language model

## A Search Process Based on Three Knowledge Sources



- Acoustic Models : HMMs for basic voice units (e.g. phonemes)
- Lexicon : a database of all possible words in the vocabulary, each word including its pronunciation in terms of component basic voice units
- Language Models : based on words in the lexicon

# Maximum A Posteriori Principle (MAP)

$$W : \{ w_1, w_2, w_3 \}$$

↑	↑	↑	$P(w_1)$ $P(w_2)$ $+ P(w_3)$ <hr style="width: 100px; border: 0; border-top: 1px solid black; margin-top: 5px;"/> 1.0
sunny	rainy	cloudy	

$\vec{X} = (x_1, x_2, x_3, \dots)$       weather parameters

## ● Problem

given  $\vec{X}$  today, to predict  $W$  for tomorrow

# Maximum A Posteriori Principle (MAP)

## ◎ Approach 1

Comparing  $P(w_1)$ ,  $P(w_2)$ ,  $P(w_3)$

$\vec{x}$  not used?

Approach 1

Likelihood function

Prior Probability

A Posteriori Probability

事後機率

$P(\vec{x} | w_i)$

事前機率

$$, P(w_i) \cdot P(\vec{x} | w_i),$$

compute

$w_1$

$w_2$

$w_3$

$i = 1, 2, 3$

observation

$w_1$

$w_2$

$w_3$

unknown

$w_1$

$w_2$

$w_3$

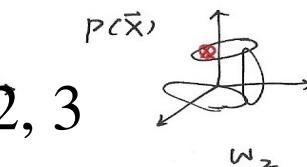
compare

$w_1$

$w_2$

$w_3$

$$P(w_i | \vec{x}) = \frac{P(\vec{x} | w_i) P(w_i)}{P(\vec{x})},$$



likelihood function  
現象の分布

prior probability

$$P(w_1, w_2, w_3 | \bar{x}) = \frac{P(\bar{x} | w_1, w_2, w_3) \cdot P(w_1, w_2, w_3)}{P(\bar{x})},$$

現在所有 parameter

$\bar{x}$  observation

$w_1, w_2, w_3$  unknown

事後確率 A posteriori probability

$$P(w_i | \bar{x}) = \frac{P(\bar{x} | w_i) P(w_i)}{P(\bar{x})},$$

$i=1 \sim 3$  doesn't matter

⇒ Compare  $P(\bar{x} | w_i) \cdot P(w_i)$

2:02:00 : CM3 Map

# Syllable-based One-pass Search

- Finding the Optimal Sentence from an Unknown Utterance Using 3 Knowledge Sources: Acoustic Models, Lexicon and Language Model
- Based on a Lattice of Syllable Candidates

