## 2. Isolated word speech recognition & HMMS

<u>Task</u>: Recognition of isolated spoken words from a given vocabulary

<u>Problems</u>:
- variable speed
- speaker independence
- prosody etc.

<u>How do we (<u>mammals</u>) hear?</u>

→ tympanic membrane → ossicles →

→ cochlea:
   basilar membrane (scala media)
   inner & outer hair cells

→ auditory cortex

## A. <u>Signal pre-processing</u>

- Sample the pressure - time function $f(t)$, digitise
   highest frequency in speech signal < 10 KHz
   → Nyquist theorem → sample with 20 KHz

- Frequency analysis: apply digital Fourier transform
   with sliding window

$$C(\omega, t) = \int_{-\infty}^{\infty} W(t-t') f(t') e^{i\omega t'} dt'$$

   Simplest window function $W(t) = \begin{cases} 1 & \text{if } |t| \leq b \\ 0 & \text{otherwise} \end{cases}$
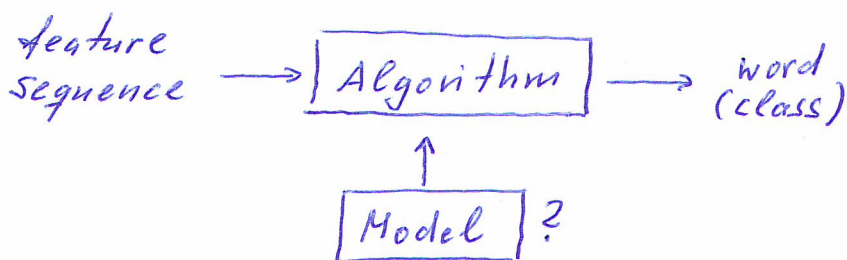
   width $b$ : lowest freq. vs. time resolution

- Energy in spectra (logarithmic, dB)

$$S(\omega, t) = 20 \cdot \log_{10} \sqrt{Re^2 C(\omega, t) + Im^2 C(\omega, t)}$$

discretise domain of $\omega$ into ~20 frequency channels
with freq. dependent width

- possibly cluster spectral vectors
    pro: small number of feature vectors
    con: dominance of stationary parts

## B. Dynamic time warping & word recognition

feature
sequence $\longrightarrow$ | Algorithm | $\longrightarrow$ word
                                              (class)
                              $\uparrow$
                          | Model | ?

Model: a set of prototypes (i.e. feature sequences)
    for each word (i.e. class)

Algorithm: nearest neighbour classifier

We need a distance measure for sequences of feature
vectors
    prototype $X = (\vec{X_1}, \vec{X_2}, ..., \vec{X_n})$, signal $Y = (\vec{Y_1}, \vec{Y_2}, ..., \vec{Y_m})$
    $\vec{X_i}, \vec{Y_i} \in \mathbb{R}^{20}$. Distance $D(x,y) = ?$

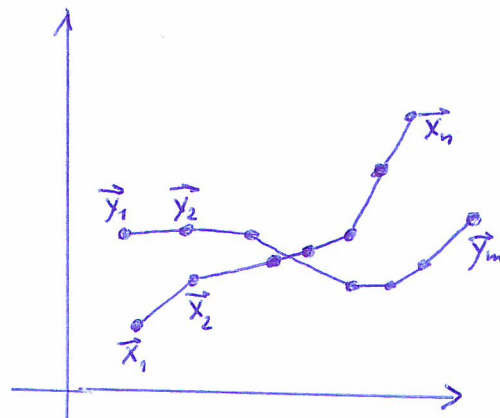Monotonous matching (aka time warping)

$\tau = ((i_1, j_1), (i_2, j_2), ....., (i_e, j_e))$

$\tau \in T$ if

(1) $(i_1, j_1) = (1, 1)$, $(i_e, j_e) = (n, m)$

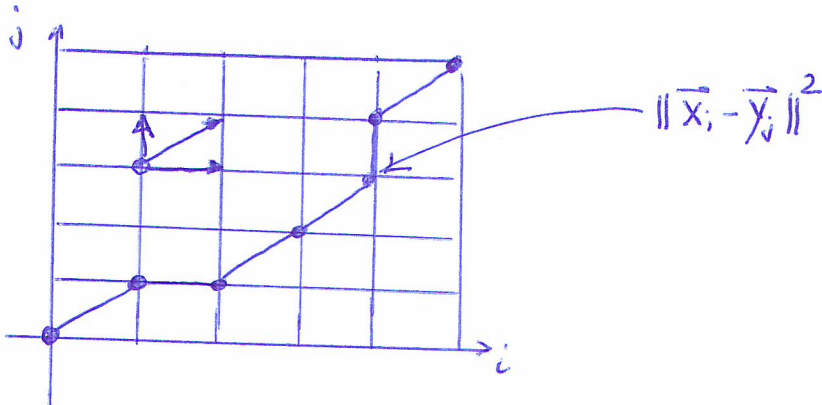(2) $i_{k-1} \leq i_k \leq i_{k-1} + 1$
    and similar for $j$

③

distance for a fixed matching $\tau \in \mathcal{T}$

$$D(x,y;\tau) = \sum_{k=1}^{|\tau|} \| \vec{X}_{i_k} - \vec{Y}_{j_k} \|^2$$

distance

$$\boxed{D(x,y) = \min_{\tau \in \mathcal{T}} D(x,y;\tau)}$$

How to compute it efficiently?



$\| \vec{X}_i - \vec{Y}_j \|^2$

i.e. shortest path, here by dynamic programming, complexity $O(nm)$

<u>Discussion</u> model & algorithm

- inference has high time complexity $O(n^2 p)$, where $p$ – total number of prototypes

- learning: how to choose optimal prototypes?

<u>Better</u>: model each word (class) by an <u>HMM</u>

$x = (\vec{X}_1, \vec{X}_2, ..., \vec{X}_n)$ – sequence of features

$S = (S_1, ..., S_n)$ – sequence of hidden states (e.g. phonemes)

$$P(x,S) = p(S_1) \prod_{i=2}^{n} p(S_i | S_{i-1}) \prod_{i=1}^{n} p(\vec{X}_i | S_i)$$

- fast inference (linear in $n$)

- feasible learning of model parameters