

## 2. Isolated word speech recognition & HMMs

Goal: Recognition of isolated spoken words from a vocabulary

Problems:

- variable speed
- speaker independence
- prosody etc.

How do we hear?

→ cochlea, basilar membrane, inner hair cells, ...  
auditory cortex

### 1A. Signal pre-processing

- Sample the pressure-time function  $f(t)$ , digitise.  
highest freq. in speech signal  $\sim 10$  KHz  
→ sample with 20 KHz

- Frequency analysis: apply local Fourier transforms with sliding window

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

simplest case:  $W(t) = \begin{cases} 1 & \text{if } |t| \leq b \\ 0 & \text{otherwise} \end{cases}$

choice of  $b$ : lowest freq. vs. time resolution

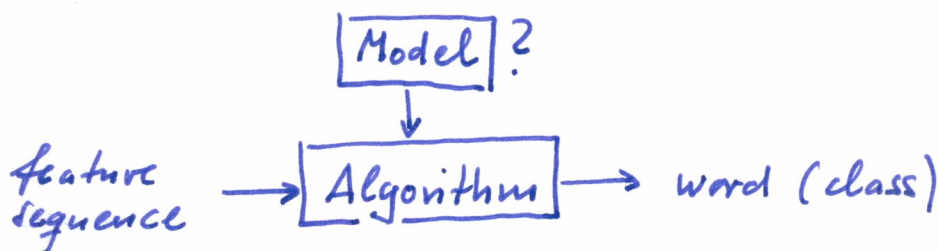
- Energy in spectra (logarithmic, dB)

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

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

- possibly also cluster spectral vectors
  - $\oplus$ : small number of feature vectors
  - $\ominus$ : dominance of stationary parts

B. Dynamic time warping & word recognition



Model: a set of prototypes ( $\cong$  feature sequences) for each word (class)

We need: distance measure for sequences of feature vectors:

prototype  $\bar{x} = (\vec{x}_1, \dots, \vec{x}_n)$ , signal  $\bar{y} = (\vec{y}_1, \dots, \vec{y}_m)$

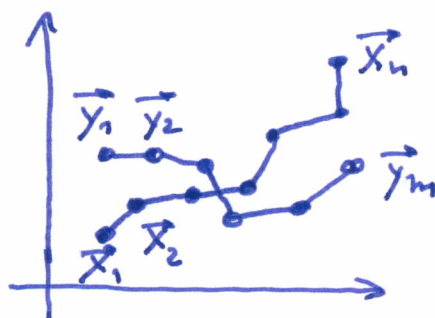
distance  $D(\bar{x}, \bar{y}) = ?$

Monotonous matching (aka time warping)

$$\mathcal{T} = ((i_1, j_1), (i_2, j_2), \dots, (i_n, j_n))$$

$\mathcal{T} \in \mathcal{T}$  if

- $(i_1, j_1) = (1, 1), (i_n, j_n) = (n, m)$
- $i_{k-1} \leq i_k \leq i_{k-1} + 1$
- similar for  $j$

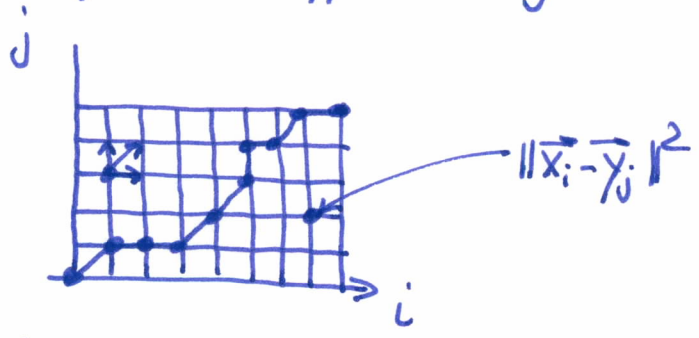


distance for a fixed matching  $\tau \in \mathcal{T}$

$$D(\bar{x}, \bar{y}; \tau) = \sum_{k=1}^{\ell_1} \|\bar{x}_{i_k} - \bar{y}_{j_k}\|^2$$

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

How to compute it efficiently?



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

Evaluate model & algorithm:

- ⊖ • many prototypes per word (class) → redundancy, slows down inference
- learning: how to choose optimal prototypes?

Better: model each word (class) by an HMM

$X = (x_1, \dots, x_n)$  - sequence of features

$S = (s_1, \dots, s_n)$  - sequence of hidden states

$$P(x, s) = p(s_1) \prod_{i=2}^n p(s_i | s_{i-1}) \prod_{i=1}^n p(x_i | s_i)$$

- ⊕ • fast inference
- feasible learning (of model parameters)