installation Tensorflow

CPU GPU TPU

librispeech  SER - Sentence error rate

WER - word error rate

WSJ - wall street journal

corpus - collection of trained
dataset

librivox - open source platform
where get speech dataset

open SLR  -

| Daniel Povey |     100, 360, 500

LM  language Model

Acoustic Model (it is trained)
Model WSJ

related to sound

| Subset | hours | per/sp min | Male | femal | total |
|---|---|---|---|---|---|
| train.clean.100 | 100.6 | 25 | 126 | 125 | 251 |

# Automatic Speech Rego
## Recognition

DTW. GMM, HMM DNN

Processing Audio Signal
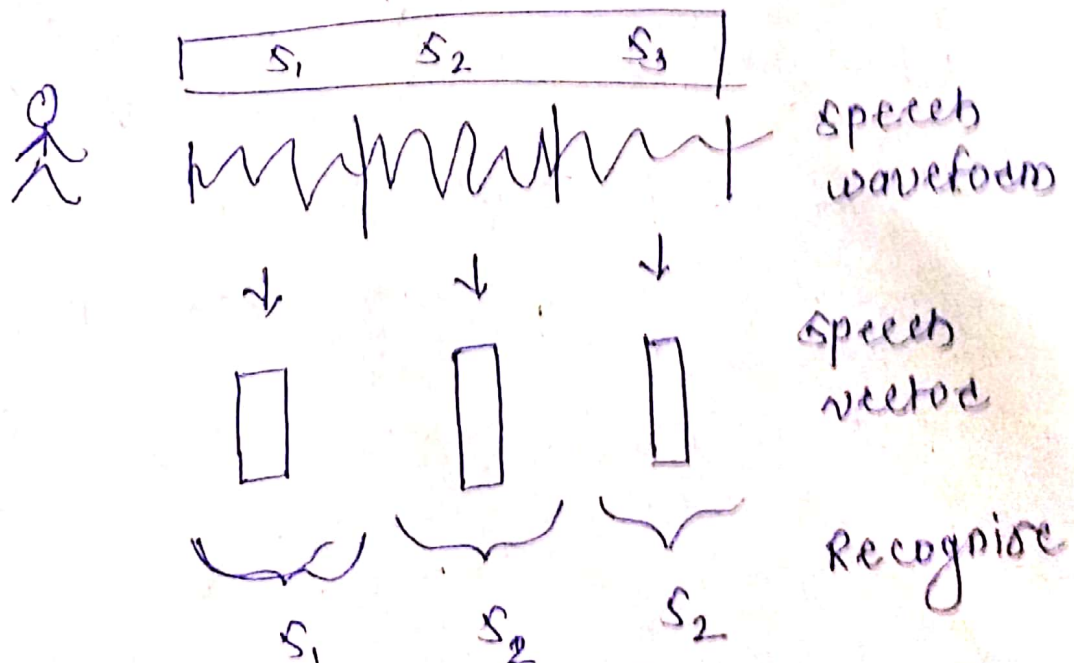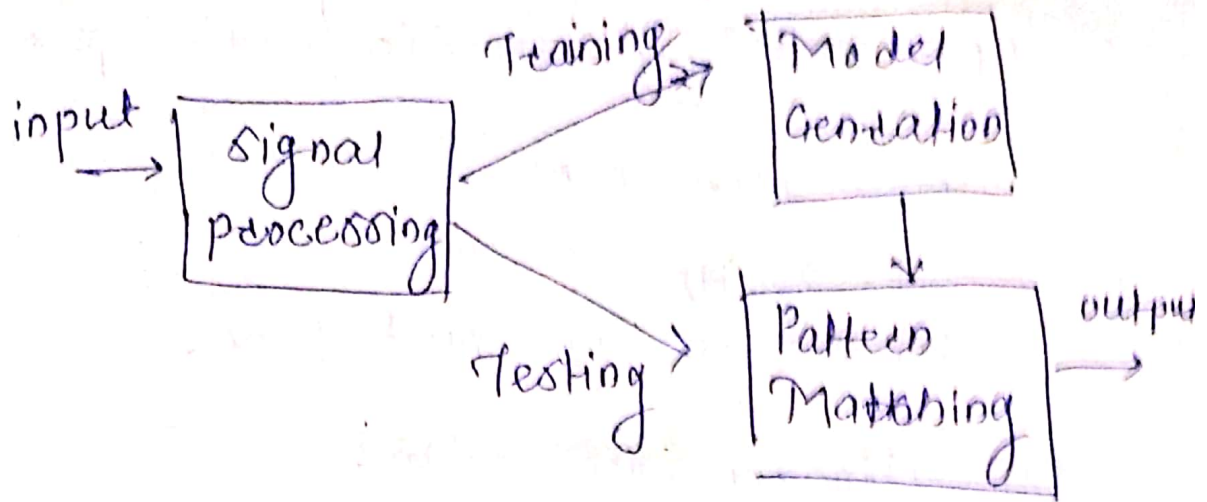
Extraction

Techniques

Statistical Model - GMM

Recognition word - DTW

Recognition sentence - HMM
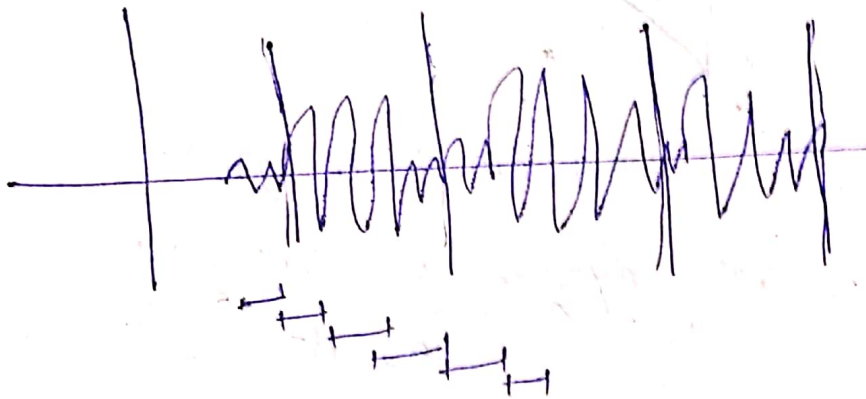
DNN - HMM

Language Model for SR

Symbol

| S₁ | S₂ | S₃ |
|---|---|---|



→ speech waveform

↓          ↓          ↓

speech vector

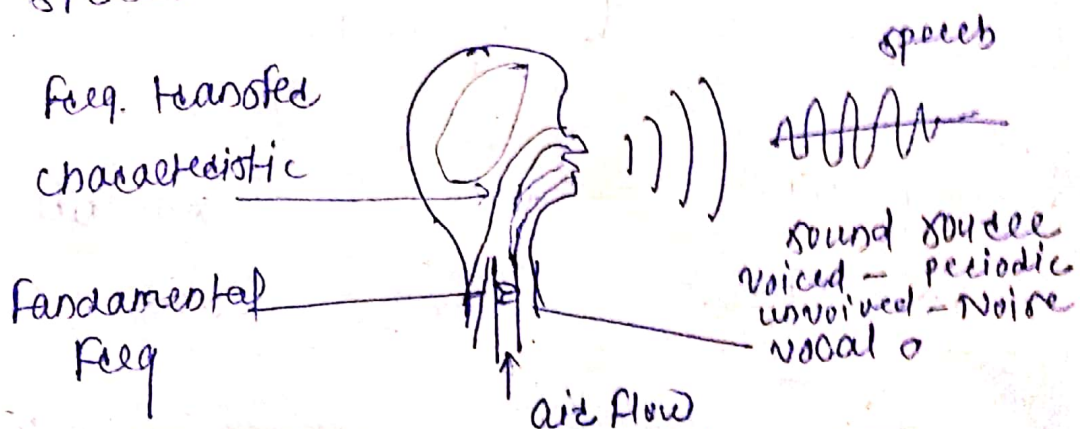S₁          S₂          S₂          Recognise

Training - learning

Testing - Recognise

\* Short time processing of SR

perform frequency analysis of short
  segment

  frame size
  frame shift
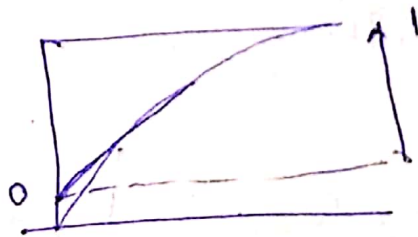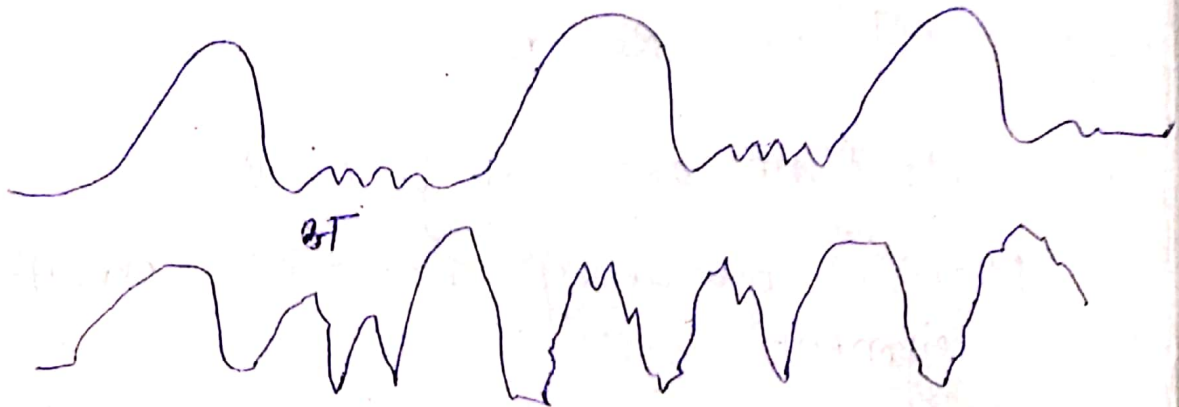


Speech Production Mechanism

freq. transfer
characteristic

fundamental
freq

air flow

speech

sound source
voiced - periodic
unvoiced - Noise
vocal o

input sound is periodic / seq. of pulse

pitch frequency

300 Hz
Female is high pitch freq.

Production of voiced sound



3T

uniform tube model

$$v = c/\lambda = 34000/4*17$$
$$= 500 Hz$$

source → | Filter | → output

glottal          vocal          speech
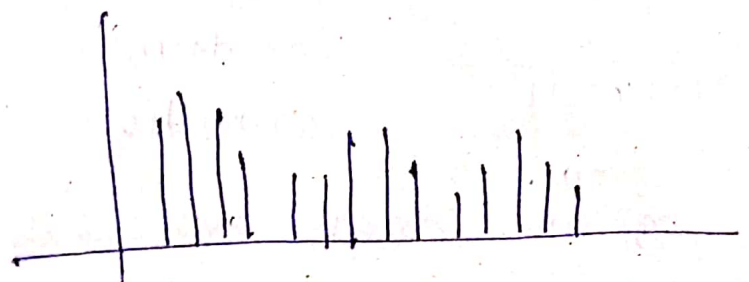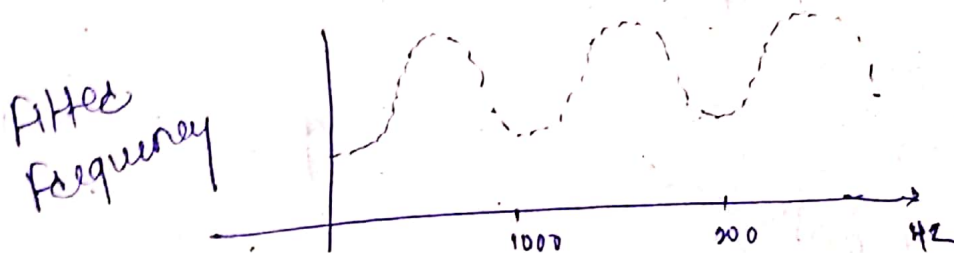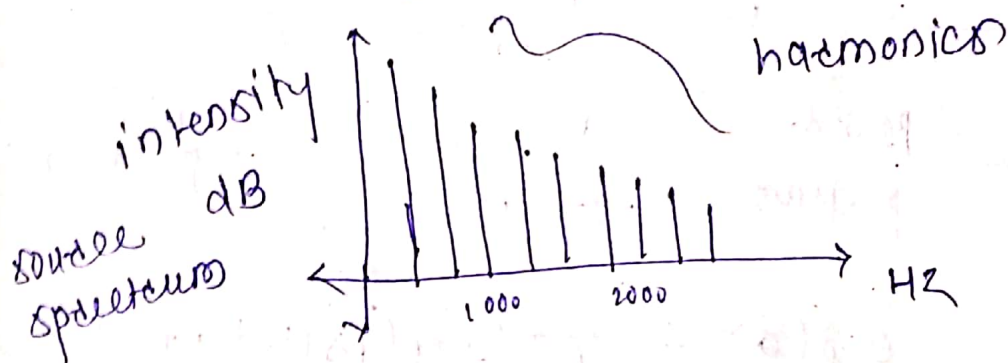vibration        tract

$$S(n) = e(n) * h(n)$$

exitation signal — impulse
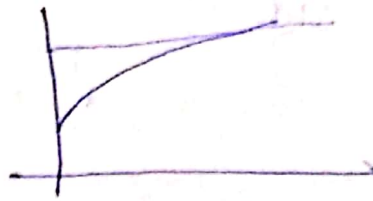
convolution

$$S(k) = E(k)H(k)$$

$$\log(|S(k)| **2) = \log(|E(k)**2) +$$
$$\log(|H(k)|**2)$$

exitation signal is periodic



intensity
dB

source spectrum

harmonics

1000    2000    Hz

filter frequency

1000    200    Hz

output energy spectrum

Glottal air flow


vocal tract


time waveform

$$\rightarrow \boxed{FFT} \rightarrow \boxed{\log} \rightarrow \boxed{IFFT} \rightarrow$$

| wave form | power spectrum | log spectrum | cepstrum (~~cepte~~) |

$$cep(q) = IFFT(\log(|S(k)| ** 2))$$

$$q = 0, 1, -- N-1$$

formants
___
peaks in spectrum

कॉकलिया

3 cm
30 mm

ear drum
vibration

20 Hz - 20,000 Hz
audio freq.

Basilar membrane

Back scale/mel scale
Half of part in linear
& the half is logarithmic

## MFCC



# triangles = # mel filters = length of
mel spectrum

$$B(m) = \sum_{k=lo(m)}^{hi(m)} |X(k)|^2$$

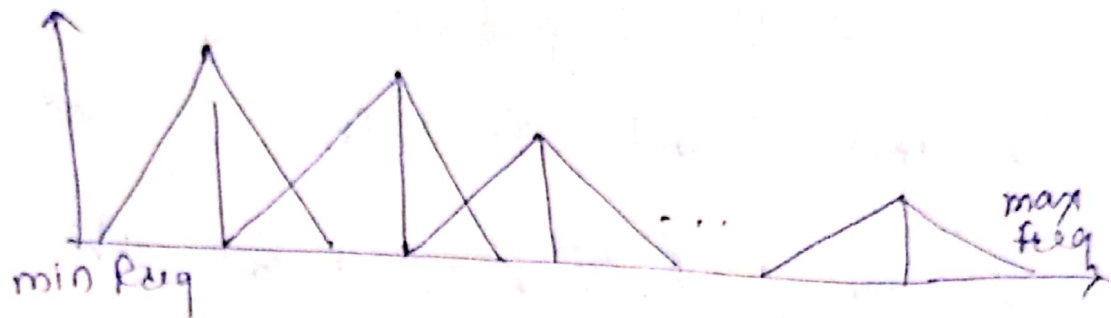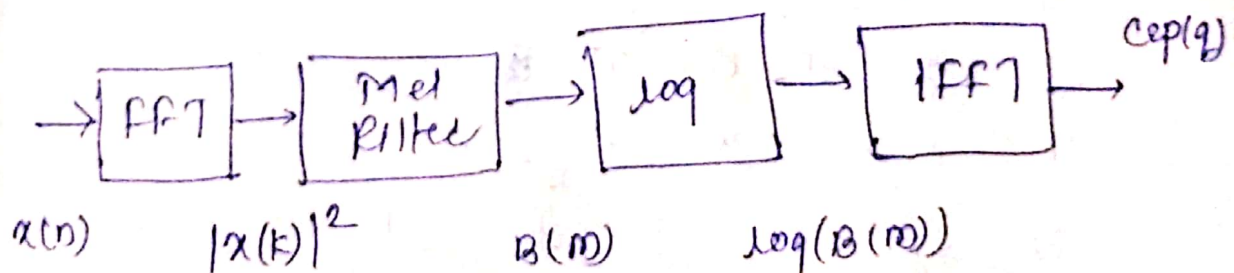$$cep(q) = IFFT\{ log(|B(m)|^2) \}$$

$$q = 0, 1, \dots N$$

| mel frequency cepstral coefficients |



extraction of features (MFCCs) that will be
used for representation as well as recognition
of speech sound.

# Acoustic phonetics

## phones and phonemes

human language → / ← smallest meaningful contrastive unit

allophones = p & ph

Aspirate Sound प & फ

| अ | आ | इ | ई | उ | ऊ | ए | ऐ | ओ | औ |
|---|---|---|---|---|---|---|---|---|---|
| a | A | i | I | u | U | e | E | o | O |

| क | ख | ग | घ | ङ |
|---|---|---|---|---|
| k | kh | g | gh | ng |

| च | छ | ज | झ | ञ |
|---|---|---|---|---|
| c | ch | j | jh | nj |

| ट | ठ | ड | ढ | ण |
|---|---|---|---|---|
| T | Th | D | Dh | N |

| त | थ | द | ध | न |
|---|---|---|---|---|
| t | th | d | dh | n |

| प | फ | ब | भ | म |
|---|---|---|---|---|
| p | ph | b | bh | m |

| य | र | ल | व | श | ष |
|---|---|---|---|---|---|
| y | r | l | w | sh | |

| स | ष | ह | ळ | श्र | ज्ञ |
|---|---|---|---|---|---|
| s | s | h | | | |

delta coefficient

$$y(x) = mx + c$$

$$\Delta cep(n,l) = \frac{\sum\limits_{l=-L}^{L} l \, cep(n,l)}{\sum\limits_{l=-L}^{L} (l)^2}$$

<u>sequence of feature vector</u>

Digitisation of analog speech signal

Blocking signal into frames

39 13   FFT → mel filter → log → IFFT → MFCC

13+13  slope and velocity curvature

sequence of feature vectors

$$: x_1, x_2 \cdots x_T \quad \boxed{39\text{-dim}}$$

$$: 0_1, 0_2 \cdots 0_T \quad \boxed{T = 150}$$

<u>linear perceptron</u>



activation fn

$$y > \theta \Rightarrow 1$$
$$y < 0 \Rightarrow 0$$

## Loss / cost Function

$$f(\omega) = 0.5 * (t(n) - y(n))^2$$

where,

$$y(n) = Sum(\omega_i * x_i)$$

$$\omega_{new} = \omega_{old} - \frac{df}{d\omega}$$

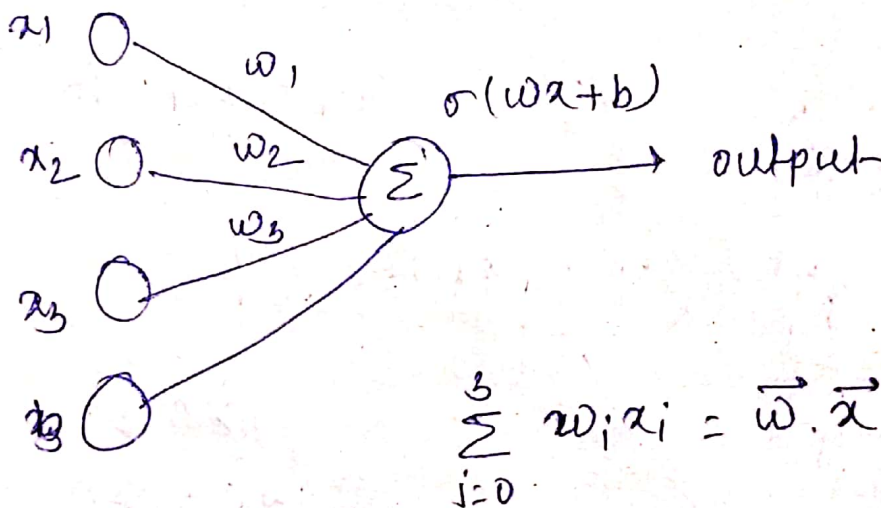A/c to gradient decent algorithm,
weight update rule,

$$\omega(n+1) = \omega(n) - \eta(d(n) - y(n)) * x(n)$$

$$y[-\infty, \infty] \rightarrow [0, 1]$$

logistic Function

$$f(x) = \frac{1}{1 + e^x}$$

## Neural Network basis : single unit



$$\sigma(\omega x + b)$$

$$\sum_{j=0}^{3} \omega_i x_i = \vec{\omega} . \vec{x}$$

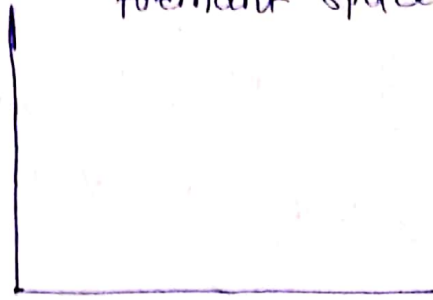$$\sigma(\omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + b) = \sigma(w'x + b)$$

$$w \in \mathbb{R}^{1 \times 3}$$

" logistic regression as a neuron "
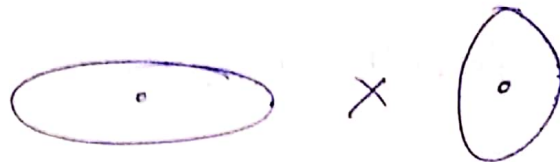
★ vowel ~~tage~~ recognition

formant space of vowels



classification criteria ( deterministic view)

Euclidean distance

$$x \in C_k \quad \text{if} \quad (x - \mu_k)^2 \leq (x - \mu_j)^2 \quad \forall j$$



weighted euclidean distance

$$d^k = \sqrt{\left(\frac{x - \mu^k}{\sigma^k}\right)^2}$$
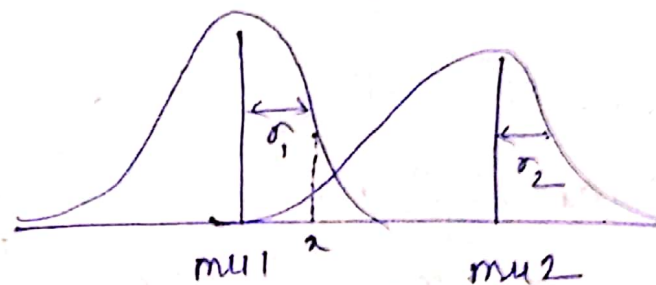
extention to multiple features

$$d^k = \sqrt{\sum_i \left(\frac{x_i - \mu_i^k}{\sigma_i^k}\right)^2}$$

DTW : Matching sequence

lexical order

give to valid word

sequence of word

Two class problem _____ (probalistic view)

Normal distribution : $N(\mu, \sigma)$

Gaussian distribution

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ -\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2 \right\}$$



mu1          mu2

$N(\mu, \sigma_1)$          $N(\mu_2, \sigma_2)$

maximum likelihood classification
criteria :

$x \in C_k$  if  $P(x/N(\mu_k : \sigma_k)) \geq P(x/N(\mu_j : \sigma_j))$

$$\frac{1}{(2\pi)^{n/2}|\Sigma|} \exp\left(-\frac{1}{2}\left\{(\vec{x}-\vec{\mu})^T \Sigma^{-1} (\vec{x}-\vec{\mu})\right\}\right)$$

x one speech frame
test data

isolated word recognition

EX. Name dialing

→ End-point detection errors
→ speaking rate variations
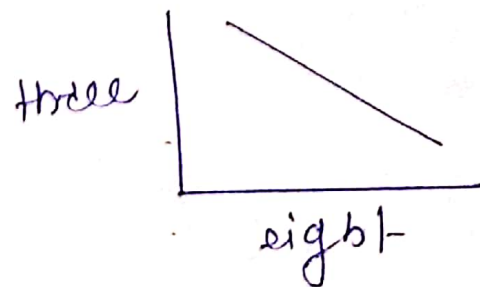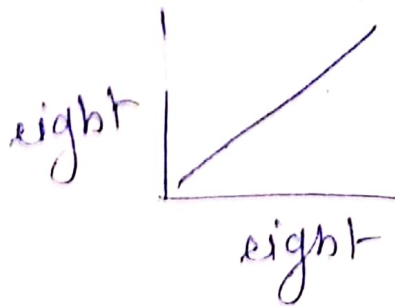→ within word variation

* linear warping → Note
matching of feature vectors
in test & reference

Greatest similarity ( lesser distance)

'eight' versus 'eight' : A path diagonal exist

'eight' versus 'three' : A path diagonal doesnot
exist.



* Dynamic programming

$$D(n,m) = d(n,m) + \min \begin{cases} D(n-1, m) \\ D(n-1, m-1) \\ D(n, m-1) \end{cases}$$