

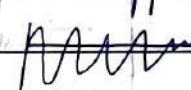
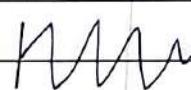
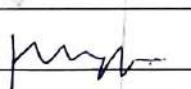
classmate
Date 23/07/05
Page

$f(x_1, x_2, \dots, x_n)$
 → dependent variable or function
 by one or more independent variables
 → A physical quantity that varies with time

- i) Basic Understanding of Signals - time implicit, amplitude
- ii) Find meaning out of these signals - how to do it

	which protocols
	which technique
	if they work
	how much
	Where they fail?

- N_o attendance - N_o 75% }
 iii) N_o midsem - 5% }
 N_o end sem - 15% } \Rightarrow continuous evaluation (N_o makeup quiz)
 • (2 quizzes/week) quiz
- iv) 40% weightage - Assignment : everything 'uttered'
 \Rightarrow (C-prereq) \Rightarrow Needs to be implemented

- v) Probability concepts applied like, consider sentence
 if hello :  we can have all three variations


 but, how can we identify if all 3 corresponds to 'Hello'
 \Rightarrow for that we look at how probability produces that
 probability with probability produces that
 most dominant signal

- vi) Speech Recognition System: Final Deployed project //

- vii) Books : ~~Rabiner & Juang~~ \rightarrow (Read this)

* Fundamental of Speech Recognition

(i) Rabiner & Schafer

Digital Processing of Speech Signals

(iii) Vector Quantization & Signal Compression - Gersho & Gray

(iv) Khalid Sagood - Data Compression

viii) Skills to develop for assignment

QR code on

Notice ←

Board

in CSE

dept

Visual Studio 2010

- Cooledit 2000 - for editing the speech waveform
- :

Q) Drones are detectable. Can we build a system to detect? → Signature of different drones (frequency)
 → And do something (....)

12/08/25

Fact 1

~~Review~~ i) What are the things that go into
 SP, P → Speech processing (CLIPS)

Areas:
 A) Signal Processing : It involves
 (SP_1) extracting relevant parameters
 from the speech signal

→ It should be efficient & robust
 SP_2 algorithms ↓ no ↓ failure
 real time

B) Physics : (Acoustics) Relationship
 between the physical speech signal
 & physiological mechanism - vocal
 tract, hearing mechanism

C) Pattern Recognition : These are algorithms to cluster data and
 create templates or models and matching procedures.

D) Communication and Information Theory : Procedures for estimating
 parameters of statistical models, detecting the presence or
 absence of speech patterns, Coding & decoding

E) Linguistic Knowledge : Relationship between sounds (phonology), words in a language (which is a syntax), meaning (semantic), and sense derived from meaning (pragmatics)

CD stores music digitally & not mp3.

~~What is mp3 format?~~

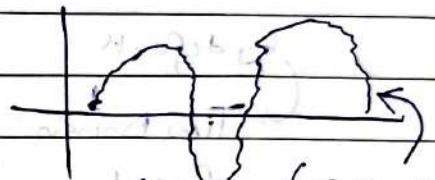
Pulse code modulated (PCM)

data = cd stores

→ no compression

→ grammar req for it

F) Physiology : Understanding higher order mechanisms within the human central nervous system that account for speech production & perception in humans (embed this knowledge in ANNs)



(actual groove)

G) Computer Science : Study of efficient

algorithms for implementing in software & hardware, the various methods of speech processing systems

⇒ Compression : Huffman Coding

types

lossy

lossless coding

coding (mp3)

H) Psychology : Understanding the factors that enable human to use speech processing systems

Speech requires language

Different Approaches to Speech Processing

- Acoustic Phonetic Approach
- Pattern Recognition Approach
- AI based Approach

DSP

Step 1: Speech Production (Speaking)

message formulation

Physical form → Mathematical form

Speak - selection of a language = language code

neuro-muscular action

• Speech

→ Homework ←

classmate

Date _____
Page _____

give us signal: { Regularity
ture of an alphabet in spacing } \leftarrow

Draw the wave for each
alphabets as you speak

\Rightarrow (pitch period)

\rightarrow Record

\rightarrow from wave (select some)

\rightarrow Save selection \rightarrow in ASCII
format

Step 2: Speech Perception (Listener)

- Capture of (longitudinal wave) air waves through ear

\rightarrow starts off with the acoustic wave, hits the basilar membrane,

Modular

steps

In step 2

\rightarrow does a neural transduction

\rightarrow thus language code is decoded

\rightarrow finally message comprehension

Assignment

by default

\rightarrow Time Domain

Board

Approach

another is: frequency

how many

\rightarrow oscillations per

sec

for this: view + spectral view

(i) record

(ii) select some part - the

signature (search

it up) = pitch period

(iii) save in ASCII format

(iv) run it again

(v) write characteristics for
Spectrogram & the wave

After opening, the
cool edit

• F10 btn - press : Red

CD1 or ~~bar~~ bar

- measured in decibels

- to measure noise
(on record btn press)

• Sample rate - samples/s

\hookrightarrow as digital $\&$ (...)

• Mono channel - single
channel

• each sample 16bit - reso-

lution

How to deal with the
noise (background noise)

human

what do we the negative represent?

• diaphragm movement

• Record vowels - 10 times

(τ , 10 different files)

\Rightarrow A E I O U

\Rightarrow {अ आ इ ई ऊ उ श्व ष्ट ष्ट्ट } ??

Sampling theorem: Maximum frequency that can be
achieved at sampling rate n is

$\frac{n}{2}$

\Rightarrow (Homework)

• Record vowels - 10 times

(τ , 10 different files)

\Rightarrow A E I O U

\Rightarrow {अ आ इ ई ऊ उ श्व ष्ट ष्ट्ट } ??

Go to analysis \rightarrow ()

\hookrightarrow instantaneous analysis

\hookrightarrow 2D graph for 3D view

↑ peaks are the resonating
frequency

Assignment 1: finding out difference b/w CLASSMATE

You & no

→ build a speech recognition system which differentiates you & no

Date _____

Page _____

→ Algorithm will be provided

→ how does yes & no look like chunk?
↓ it doesn't have

1 Logic:

has ssh.. noise
(missing sound)

→ Clip off the noise

$$\bullet \text{ Energy} = \frac{1}{N} \sum_{i=1}^N x_i^2$$

↑ Sample height

(we can remove as well)

(like N = fix to 100)

$$\bullet \text{Zero Crossing Rate} = \frac{\text{no. of zero-crossings}}{\text{total no. of samples}} \times 100$$

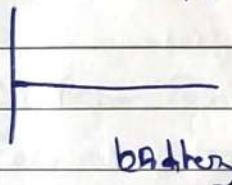
(ZCR) or $\frac{\text{no. of zero-crossings}}{\text{total no. of samples}} \times 100$

$m + n$

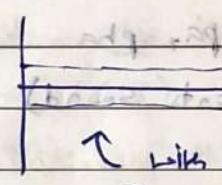
Steps to take before recording

20/08/25

a)



battery

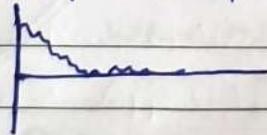


↑ with charger

Can be corrected by DC shift correction

(ii) get a good mic (find it out) before recording

(iii) a big pulse before the actual recording



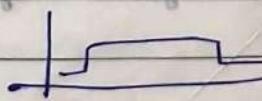
ignore
Starting
data

↑ no sound; why does it come?

⇒ When tap on Start button, the electromagnetic pulse results in that case

→ cut it out & then throw it (like 1 sec etc)

(iv) finding out safe limit: not too louder too low (close to the mouth & away from hand)



right near to chin

(iv) Normalisation : (i) find highest or lowest no. in the recorded of date

(ii) Scale it to ± 5000

(This is known as equal loudness)

\rightarrow (here we only control volume not signal distortion)

This can overflow on

(-----)

Different Types of Sounds & Their Characteristics

i) Vowels

\rightarrow how is a vowel pronounced — Quiz

\rightarrow Steady sound but [i] [e] [a] [o] [u] etc. = phone : inbuilt microphone | condenser microphone
like tongue movement etc, power supply
alignment of vocal chords etc

\rightarrow regular pitch periods

ii) Consonants — pa, pha

\rightarrow plosives (consonants sound)

\rightarrow closure somewhere in the vocal tract

\rightarrow closure somewhere in the vocal tract : noisy sound

Voiceless

iv) Fricatives

v) Nasals — m, n, ñ (ni)

\rightarrow both mouth & nose used

\rightarrow Antiresonance : wave coming from mouth & nose
if have same freq \rightarrow if opposite in sign

\Rightarrow they cancel each other

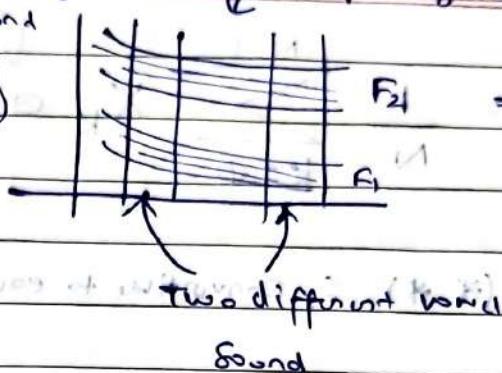
\rightarrow can be used for identification of ppl

two vowel sounds that are pronounced together to make one sound

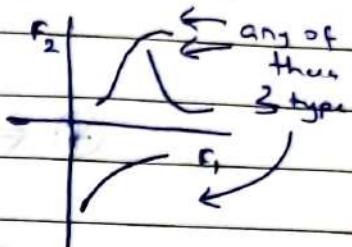
vii) Diphthong \rightarrow

- Sliding formant (??)

frequency



\Rightarrow Plot F_1 & F_2 graph



(viii) Semivowel \rightarrow aka lateral (ऽ ः)

How to parametrically represent waveform speech sounds

i) Energy

ii) ZCR - extremely low in vowels, high in fricatives & plosives.

iii) Short term Spectral Envelope - instead of point eval, do in intervals

• Time Domain Based Approach } either one of them only used in

features of speech } Frequency Domain " " deployment

• based on Fourier transform

Discrete Fourier Transform (DFT) \rightarrow give the frequency content of a signal

For a discrete signal $x(n)$, if the signal is periodic with period N
 \hookrightarrow signals having finite values at fixed point of time
 then $x(n)$ can be written as :

$$\tilde{x}(n) = x(n+N) \quad -\infty < n < \infty$$

signal repeats after N time

(vowels have this property) - i.e. periodic
 where N is periodicity

Then $\tilde{x}(n)$ can be represented by a discrete sum of sinusoids.

The Fourier series representation of a periodic sequence is given by,

$$\tilde{X}(k) = \sum_{n=0}^{N-1} \tilde{x}(n) e^{-j \frac{2\pi}{N} nk} \quad \dots \star$$

{ Cubic operation }

(for a particular k we try to sum up)

With this we can regenerate $\tilde{x}(n)$

$$\tilde{x}(n) = \frac{1}{N} \sum_{k=0}^{N-1} \tilde{X}_k e^{\frac{j2\pi}{N} \cdot kn} \quad (*)$$

$(*) \Leftrightarrow (**)$: Convolution to each other

- Fourier Transform (FT) - (not a periodic signal)

The Fourier Transform representation of a signal $x(n)$ is given by :

$$X(e^{jw}) = \sum_{n=-\infty}^{\infty} x(n) e^{-jwn} \quad \text{--- (1)}$$

(not periodic) $\rightarrow n = -\infty$ where all $x(n)$ $\rightarrow 0$
 turned to 0

$$\text{Observe compute } x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{jw}) \cdot e^{jwn} dw$$

with a factor e^{jwn} in front where w = frequency

(T2A) \rightarrow constant value \rightarrow FT

NOTE: if noise there in environment - FT not good

-> friction (like h), case of friction and slope \rightarrow

\rightarrow no constant at $n=0$ ($x(0) \neq 0$)

A sufficient condition for the existence of FT is

$$\sum_{m=-\infty}^{\infty} |x(m)| < \infty$$

$$\sum_{m=-\infty}^{\infty} |x(m)| < \infty$$

$m = -\infty$

1) write for a condition of FT through $x(n) \neq 0$ $\forall n$

so if $x(n) \neq 0$ \rightarrow to all the right side must $\neq 0$ $\forall n$

$$x(n) \neq 0 \rightarrow x(n+1) \neq 0 \rightarrow x(n+2) \neq 0 \rightarrow \dots = (1) \neq 0$$

can do same for the left side remaining $\neq 0$

so $x(n) \neq 0$

✓ gives the overall picture of the system

Z-Transform

• mathematical tool to convert discrete time signals
to

its Z-domain representation

$$\text{Discrete Time} \cdot X(z) = \sum_{n=-\infty}^{\infty} x[n] z^{-n}$$

Time Domain Signal Sampled Periodically

Domain



Two sided Z transform

$$X(z) = \frac{z}{z-1} \quad (\text{Z domain})$$

eqn

(for an integrator system which outputs 1)

- Discrete Time Domain $\iff X(z) = \frac{z}{z-1}$

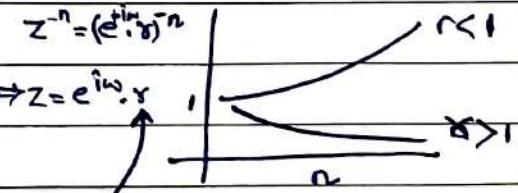
\Rightarrow Z-domain representation

Part 1

- How to solve Z-transform eqns

NOTE:

DFT $\star r^{-n}$ = make it
 $\Rightarrow \sum_{n=-\infty}^{\infty} x[n] e^{-j\omega n}$ the Z-transform
 $\star r^{-n}$ (a positive real)
 $\Rightarrow r < 1$



- This r has convergence pt!

The Z-transform

It is given by

$$X(z) = \sum_{n=-\infty}^{\infty} x(n) \cdot z^{-n}$$

The complete signal

$$x(n) = \frac{1}{2\pi j} \int_C X(z) z^{n-1} dz$$

$C: (z-a)^2 + (y-\omega)^2 = R^2$ $\omega = \alpha j$

If we replace

$\omega = j 2\pi f$ we get Fourier Transform

depending upon ω we can get different transform (or signal)

Fundamental frequency - first frequency which is maximum from the lower frequency side

→ It represents the vibration of vocal cord

Time Domain Analysis of Speech

Analysis of speech wrt time only (no frequency)

Largely based on Linear Predictive Coding (LPC)

Why only LPC?

⇒ Cellphones running this algo (very old; it actually works)

i) It's a good model for periodic parts of speech

ii) Good source - Vocal tract, Separation

↳ lung system / pressure + functioning of vocal chord

remaining part of vocal tract &

articulators

iii) It's a mathematically stable model - 0 chance of system crashing

iv) Pronunciation method (it is)

Driving LPC

let $s(n)$ be the speech signal at time point n , this can be approximated as a linear combination of past ' p ' samples such that,

$$s(n) \approx a_1 s(n-1) + a_2 s(n-2)$$

$$\dots + a_{p-1} s(n-p) \dots \quad (1)$$

\oplus

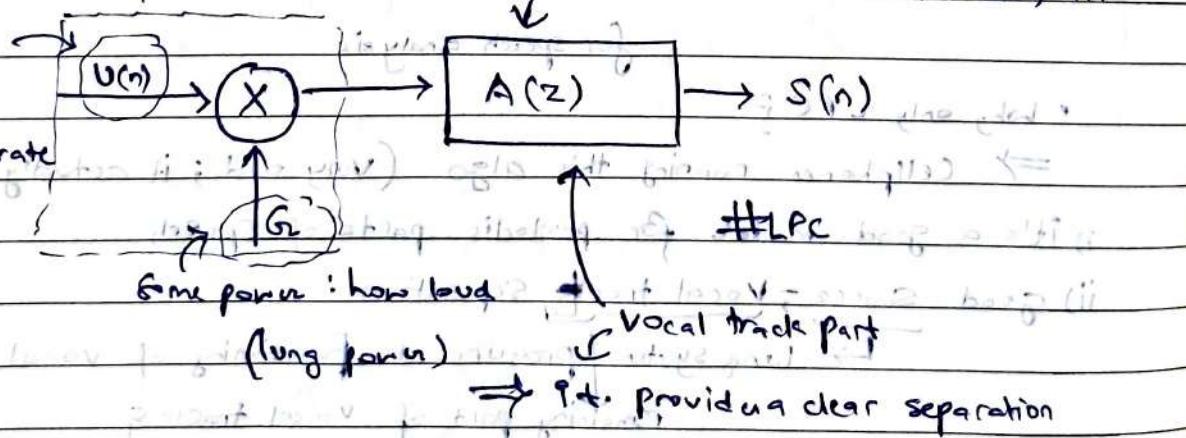
for predicting this, we need past p samples. Where a_i , ($i = 1, 2, \dots, p$) are constants. Now convert eqn (1) to an equality by including a term called excitation term given by, $-G \cdot v(n)$ giving resultant

$$s(n) = \sum_{i=1}^p a_i s(n-i) + G \cdot v(n) \quad (2)$$

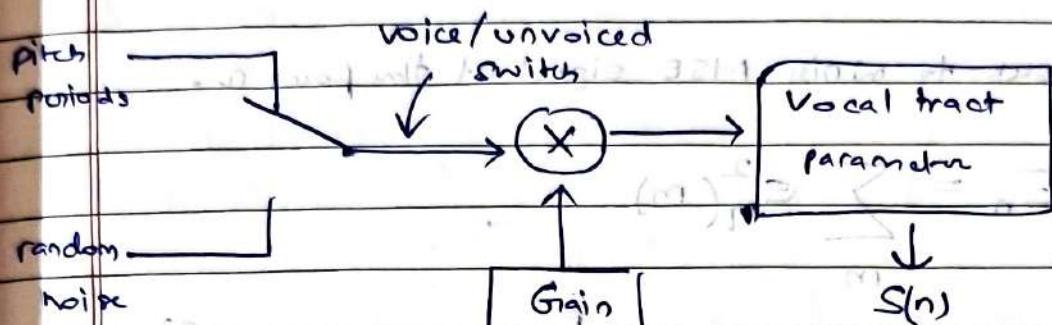
(determined for that analysis frame)

Where $v(n)$ is the normalized excitation. G , G is the gain term with high enough value.

decision if vocal chord should vibrate or not



Vocoders - read little bit about it (A para)



Let us consider,

$$\tilde{S}(n) = \sum_{i=1}^p a_i s(n-i) \dots \quad (3)$$

(Assuming no source
not absolute) $G_i, U(n) = 0$

speech signals

Now the prediction error will be given by,

$$e(n) = S(n) - \tilde{S}(n)$$

$$= S(n) - \sum_{i=1}^p a_i s(n-i)$$

how much signals should we analyze?

for this

- We need to analyze the error for a short segment of the speech

Signal

$s(n)$ is fixed

a_i is variable to make $e(n)$ low

$s(n-i)$ is fixed

- now let us take $S_n(m) = S(n+m)$

$$e_n(m) = e(n+m)$$

m tells the range

of analysis

- Now we seek to minimize MSE signal at time point n ,

i.e.

$$E_n = \sum_m e_n^2(m)$$

$$= \sum_m \left\{ S_n(m) - \sum_{k=1}^P a_k S_n(m-k) \right\}^2$$

$\rightarrow P$ quantities unknown (a_k)

$\rightarrow P$ equations { since we want to minimize } \Rightarrow differentiate

Differentiating partially wrt each a_k & setting it equal to zero

$\frac{\partial E_n}{\partial a_i} = 0$, will give arbitrary constraint

$$\Rightarrow \frac{\partial (E_n)}{\partial a_i} = \frac{\partial}{\partial a_i} \left[\sum_m \left\{ S_n(m) - \sum_{k=1}^P a_k S_n(m-k) \right\}^2 \right] = 0$$

$$= \frac{\partial}{\partial a_i} \sum_m \left[S_n^2(m) + \left\{ \sum_{k=1}^P a_k S_n(m-k) \right\}^2 - 2 S_n(m) \cdot \left\{ \sum_{k=1}^P a_k S_n(m-k) \right\} \right] = 0$$

$$= 0 + \frac{\partial}{\partial a_i} \sum_m \left(\sum_{k=1}^P a_k^2 S_n^2(m-k) \right)$$

$$+ 2 \cdot \sum_{k=1, k \neq i}^P a_i a_k S_n(m-i) \cdot (S_n(m-k))$$

$$= \sum_m S_n(m) \cdot \frac{\partial}{\partial a_i} \left\{ \sum_{k=1}^P a_k S_n(m-k) \right\}$$

thus 2 can be clubbed

classmate

Date _____

Page _____

$$= \sum_{i=1}^m x a_i (S_n^2(m-i))$$

$$+ x \sum_{m=1}^p \sum_{k=1, k \neq i}^p a_k S_n(m-i) (S_n(m-k))$$

$$= \sum_{m=1}^p S_n(m) (S_n(m-i))$$

$$\Rightarrow \sum_{m=1}^p \sum_{k=1}^p a_k S_n(m-i) \cdot S_n(m-k)$$

$$= \sum_{m=1}^p S_n(m) S_n(m-i)$$

$$= \sum_{m=1}^p S_n(m-0) S_n(m-i)$$

$$\Rightarrow \sum_{k=1}^p a_k \phi_n(i, k) = \phi_n(0, i) \quad \dots \quad (1)$$

where $\phi_n(i, k) = \sum_{m=1}^p S_n(m-i) \cdot S_n(m-k)$

$$\Rightarrow p = 2 \text{ to } 12$$

$$i = 1 \text{ to } p$$

P eqn's - how many unknowns? - P of them

$$(x-i+y) \cdot 2 (y) \cdot 2 \sum_{i=1}^p (a_i)$$

$$(i-j+1-y) \cdot 2 (i) \cdot 2$$

$$(i-j+1-y) \cdot 2 (i) \cdot 2$$

~~No²~~ Autocorrelation Method

→ To find a_i 's

i) In this method, we take, $0 \leq m \leq N-1$ ($\Rightarrow N$ speech values)

⇒ determinant of cofactor matrix > 0 (normally $N=320$)

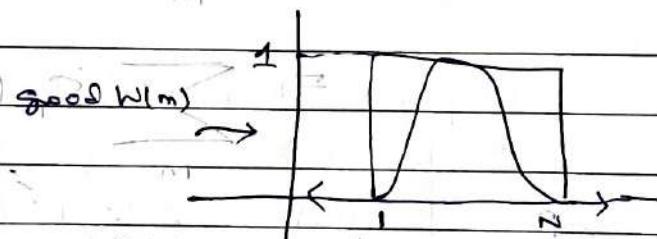
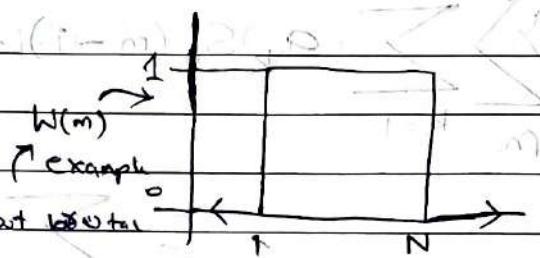
ii) Signal ceases to exist beyond the interval

i.e. available in the frame only

aka short term
Speech analysis

$$S_n(m) = w(m) S(m+n)$$

↑
Controls signal when $m < 0$ or $m \geq N$
aka windowing function



① We can now write,

$$\Phi_n(i, k) = \sum_{m=0}^{N-1} S_n(m-i) S_n(m-k)$$

does the signal ever reach to $N-1$? (No) \Rightarrow infinite sum!

⇒ as it ceases to exist

\Rightarrow therefore add $p : \{N-1+p\}$ limit of m

$$y+1 \text{ put } m-i = y$$

so $y = -i \rightarrow$ windowed signal \rightarrow $y = 0 \rightarrow$ if

$$\Phi_n(i, k) = \sum_{y=0}^p S_n(y) S_n(y+i-k)$$

$y=0$ but not significant so (0)

as $y = -i$ to $(N-1+p-i)$ but not significant

$$y = \underbrace{-i}_{0} \text{ to } \underbrace{N-1+p-i}_{N-1-(i-k)}$$

as beyond $(i-k)$ something not

$N-1-(i-k) > i-1$ significant

$$\phi_n(i, k) = \sum_{y=0}^{\infty} s_n(y) \underbrace{s_n(y+i-k)}_{\rightarrow (\text{Auto correlation fn of } s_n(y) \text{ at a lag of } i-k)} \quad (6)$$

Correlation: Degree of agreement b/w two sinus

$$\begin{matrix} x: \\ y: \end{matrix} \left. \begin{matrix} \\ \end{matrix} \right\} 2 \text{ sinus}$$

$$R(x, y) = \frac{1}{N} \sum_{i=1}^N x_i y_i$$

random variables (actual phenomenon values)

Ques 1: If y replaced with x , {high value for the correlation?}

$\rightarrow R(x, x)$ is energy {possible value for -ve "}

\rightarrow it is the maximum called energy

overlap

Case 2: $x_1: 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0$

$$\downarrow \downarrow \downarrow \downarrow \downarrow \downarrow \downarrow$$

$$x_2: (0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0) \ 0$$

$$R(x_1, x_2) = \frac{1}{N} \sum_{i=1}^N x_1 i x_2 i+1$$

Auto correlation of

lag 1

with lag = 0 \Rightarrow energy

\Rightarrow lag \rightarrow if increased

Auto correlation

\nwarrow pitch pulses

$$- (0.002 \text{ sec})^{-1} \leftarrow -f = \text{pitch unit frequency} \text{ decreases}$$

Now we can write (6) is the short term autocorrelation function with lag ' $i-k$ ' i.e., $\phi_n(i, k)$

$$\phi_n(i, k) = R_n(i-k) \xrightarrow{\text{put modulus}} \\ = R_n(|i-k|) \xrightarrow{\text{operator}}$$

It is an even function

Finally eqn (4) becomes,

$$\sum_{k=1}^p a_k R_n(|i-k|) = R_n(i) \quad i=1, \dots, p$$

In matrix notation, we get

$$\begin{bmatrix} R_n(0) & R_n(1) & R_n(2) & \dots & R_n(p-1) \\ R_n(1) & R_n(0) & \dots & \dots & R_n(p-2) \\ R_n(2) & R_n(1) & R_n(0) & \dots & R_n(p-3) \\ \vdots & & & & \vdots \\ R_n(p-1) & R_n(p-2) & \dots & \dots & R_n(0) \end{bmatrix}_{P \times P} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \vdots \\ \alpha_p \end{bmatrix}_{P \times 1} = \begin{bmatrix} R_n(1) \\ R_n(2) \\ \vdots \\ \vdots \\ R_n(p) \end{bmatrix}_{P \times 1}$$

Cofactor matrix

(P1) Dimensionality: $P \times P$ (P by P , not P cross P)

(P2) Symmetric matrix: So can be solved efficiently

(P3) Diagonal values all same

Non-diagonal values - non-principle all same

Choose any diagonal - all elements same

(P4) How many unique elements $P - p$ \Rightarrow space footprint reduced

Total no. of elements $= P^2$

Time complexity $= P^3 \rightarrow P^2$ (we need) - we use now!!

What is that algorithm which reduces from p^3 to p^2 ? 1-11

- This particular cofactor matrix has got a special name called, Toeplitz matrix.

A very efficient soln to the set of normal equation is provided by Levinson Durbin's Algorithm (Can check derivation)

(Algorithm)

Initialise $E^0 = R(0)$ ← energy $\leftarrow p=3$

$$k_i = \{R(i) - \sum_{j=1}^{i-1} x_j^{(i-1)} R(i-j)\} / E^{(i-1)}$$

$x_i^{(i)} = k_i$ $1 \leq i \leq p$

2nd matrix

$$x_j^{(i)} = x_j^{(i-1)} - k_i x_{i-j}^{(i-1)} \quad 1 \leq j \leq i-1$$

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

$\Rightarrow R(i) = \text{energy} = E^0$ # Hardcode it

$$k_1 = R(1)$$

$$R(0)$$

$$x^{(p)}$$

• T.C = p^2

Induction, it is guaranteed true

Homework

Do for $p=2, 3$

$$\{(-1), (1), (-1)\} = 1$$

$$S = S_{\text{avg}}$$

Method to
find LPC
coefficients

M-1 The AutoCorrelation Method

- Levinson Durbin Algorithm

- ensure det is not singular

- if $\det = 0$, error not divided by 0

linear PLS — bottleneck with inverse calculation

• When will the big matrix of left is 0?

↳ if $R(0) = 0 \Rightarrow \text{Energy} = 0$? no signal at all

$R(0)$ very small ? ignore it

$\approx 10^{-4}, 10^{-3}$ — (i) $f_0 = 0$

Normally we will take $p = 12$ $\rightarrow R(0) \rightarrow R(12)$

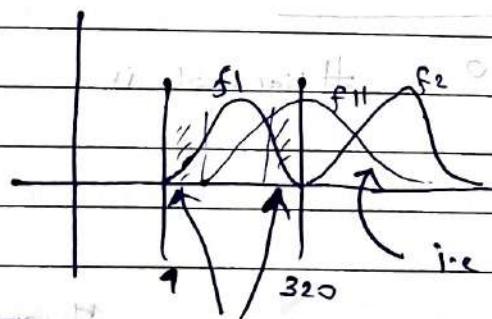
$N = 320$ (frame size) (20ms)

Why?

Sampling rate = 16000 samples/sec

16 bits/sample

Skip = 80 samples \rightarrow slide by 5ms



i.e. why overlapping we should do

just to avoid loss

\Rightarrow gives reliable answer

(take 25%)

Exercise

Consider $p=2$,

$$R = \{ r(0), r(1), r(2) \}$$

$$a_1, a_2 = ?$$

To compute a_1 & a_2 solve this pb using cross multiplication

Method 8, the Levinson Durbin Method.

\Rightarrow eqn for cross multiplication:

$$\begin{bmatrix} r(0) & r(1) \\ r(1) & r(0) \end{bmatrix} \begin{bmatrix} a[1] \\ a[2] \end{bmatrix} = \begin{bmatrix} r(1) \\ r(2) \end{bmatrix}$$

M-2 Covariance Method

\rightarrow (signal is untouched)

$$E_n(i, k) = \sum_{m=0}^{N-1} e_n^2(m) \quad (1, 1), \phi$$

$$\Phi_n(i, k) = \sum_{m=0}^{N-1} s_n(m-i) s_n(m-k) \quad 1 \leq i \leq p \quad 1 \leq k \leq p$$

now we get,

$$\Phi_n(i, k) = \sum_{m=-i}^{N-1-i} s_n(m) s_n(m+i-k) \quad \text{or}$$

p extra seg

$$\Phi_n(i, k) = \sum_{m=-k}^{N-1-k} s_n(m) s_n(m+k-i) \quad \text{or}$$

We need a few samples which is prof to the interval under consideration.

now the eqn becomes,

$$P (x_{t+1}, x_t) \phi = (x_t, \phi)$$

$$-P \leq m \leq N-1$$

$$\sum_{k=1}^P a_k \Phi_n(i, k) = \Phi_n(i, 0) \quad \text{Can } i \leq k \leq P$$

$$(x_t) = \dots \quad i = 1 \dots P$$

This gives

$$\begin{bmatrix} \phi_n(1,1) & \phi_n(1,2) & \cdots & \phi_n(1,p) \\ \phi_n(2,1) & \phi_n(2,2) & \cdots & \phi_n(2,p) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_n(p,1) & \phi_n(p,2) & \cdots & \phi_n(p,p) \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \vdots \\ \alpha_p \end{bmatrix}$$

This is the

new cofactor

matrix

NOTE: $\phi_n(1,1), \phi_n(2,2)$

$\cdots \phi_n(p,p)$

won't be same now

seen

out!

as signal
range changes

If we do
 (itm, jtm)

$\phi_n(1,2) \neq \phi_n(2,1)$

will it be same?

→ yes (identically)

$(1-i-m) \geq (m)_n \Rightarrow R = (1,i)_n \phi$

• Symmetric matrix — but not toeplitz

here $\phi_n(i,k) = \phi_n(k,i) \quad \forall i, k = 1, \dots, p$

here cofactor matrix is symmetric

but,

$\phi_n(i,j) \neq \phi_n(itk, j+tk)$

\rightarrow not toeplitz

T.C = (n^3)

$$2. \begin{bmatrix} 5 \\ 5 \\ 6 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 6 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 3 & 4 \\ 5 & 6 \end{bmatrix}$$

$$= \begin{bmatrix} 3 & 0 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

classmate

Date _____
Page _____

M-3 Cholesky's Decomposition Method - applying to only symmetric positive definite matrices

$$\phi \alpha = \psi$$

$P \times P$ $P \times 1$ $P \times 1$

pivotal part of the algo

factorize ϕ as : $\{\phi = V D V^T\}$ where

how to do it ?

$V \rightarrow$ lower Δ matrix

→ search it up

$D \rightarrow$ diagonal " "

We get :

$$(V D V^T) \alpha = \psi$$

$P \times P$ $P \times P$

Let

$$D V^T \alpha = y \leftarrow \text{has the variables}$$

Then,

$$V y = \psi$$

$$\left[\begin{array}{cccccc} - & 0 & 0 & 0 & 0 & 0 \\ - & - & 0 & 0 & 0 & 0 \\ - & - & - & 0 & 0 & 0 \\ - & - & - & - & 0 & 0 \end{array} \right] \rightarrow (\text{since first row only 1 val} \Rightarrow a_1 \text{ will be found})$$

\therefore remaining elements one at a time)

in worst case (n^3)

$$L = \begin{bmatrix} 1 & 0 & 0 \\ L_{21} & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix} \quad d = \begin{bmatrix} d_1 & 0 & 0 \\ 0 & d_2 & 0 \\ 0 & 0 & d_3 \end{bmatrix}$$

$L_{ij} = a_{ik} - \sum_{j=1}^{k-1} l_{ij} l_{kj} d_j$

$L_{ik} = \frac{1}{d_k} (a_{ik} - \sum_{j=1}^{k-1} l_{ij} l_{kj} d_j)$

$i > k$ $k = 1, 2, \dots, n$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 4 & 1 & 0 \\ 7 & 8 & 1 \end{bmatrix}$$

$\therefore \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 4 & 1 & 0 \\ 7 & 8 & 1 \end{bmatrix}$

$$a_1 x + b_1 y = c_1$$

$$a_2 x + b_2 y = c_2$$

Date _____
Page _____

$$\begin{array}{l} \cancel{b_1} \cancel{x} + \cancel{b_1} y = c_1 \\ \cancel{b_2} \cancel{x} + \cancel{b_2} y = c_2 \end{array}$$

$$\frac{b_1}{-b_2} \cancel{x} + y = \frac{c_1}{-c_2} \quad \begin{array}{l} a_1 \\ a_2 \end{array}$$

$$-b_1 c_2 + b_2 c_1 = c_1 a_2 - c_2 a_1 \quad \begin{array}{l} a_1 b_2 - a_2 b_1 \\ a_1 b_2 - a_2 b_1 \end{array}$$

$$r(0) = n^2$$

$$r(1) = n(n-1) \rightarrow V$$

$$r(2) = n(n-2)$$

$$E_0 = R(0)$$

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

$$\alpha_i = k_i$$

$$\alpha_j = \alpha_j - k_i \alpha_{i-j}$$

$$\alpha_j = \alpha_j - k_i \alpha_{i-j} \quad i=1, 2$$

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

$$E_0 = R(0)$$

$$k_1 = R(1)$$

$$= 1 - \frac{\alpha_1^2}{R_0}$$

$$\alpha_1^2 = k_1 = \frac{R(1)}{R_0}$$

$$k_2 = R(2) - \frac{\alpha_1^2}{1} R(1)$$

$$= R(2) - \frac{R(1)^2}{R_0}$$

$$= R_0 - \frac{R(1)^2}{R_0}$$

$$E_1 = (1 - k_1^2) E_0$$

$$= \left(1 - \frac{R(1)^2}{R_0^2}\right) E_0$$

$$= R_0 \left(1 - \frac{R(1)^2}{R_0^2}\right)$$

$$\alpha_2^2 = k_2$$

$$\alpha_1^2 = a_1 - k_2 \alpha_1'$$

$$= (1 - k_2) \alpha_1'$$

$$= R(1) \left(1 - R(2) + \frac{R(1)}{R_0} \left(1 - \frac{R(1)^2}{R_0^2}\right)\right)$$

$$\alpha_1 = R(1) \left(1 - R(2) + \frac{R(1)^2}{R(0)(1 - \frac{R(1)^2}{R_0^2})} \right)$$

classmate

Date _____
Page _____

$$\alpha_2 = R(2) - \frac{R(1)^2}{R(0)(1 - \frac{R(1)^2}{R_0^2})}$$

~~$$= \frac{R(0) R(2) R}{R(0)} - \frac{R(0) R(2) R(1)^2}{R(0)(1 - \frac{R(1)^2}{R_0^2})}$$~~

$$R(0) \alpha_1 + R(1) \alpha_2 = R(1)$$

$$R(1) \alpha_1 + R(0) \alpha_2 = R(2)$$

$$\frac{\alpha_1}{R(1) R(2) - R(1) R(0)} = \frac{\alpha_2}{R(1)^2 - R(0) R_0} = \frac{1}{R(0)^2 - R(1)^2}$$

$$\alpha_1 = \frac{R(1)(R(2) - R(0))}{R(0)^2 - R(1)^2} \quad \alpha_2 = \frac{R_1^2 - R_0 R_2}{R_0^2 - R_1^2}$$

$$R(0) = \frac{1}{N} \sum x_i^2 \quad R(1) = \frac{1}{N} \sum x_i x_{i+1}$$

$$R_0 = \frac{1}{N} \sum x_i x_{i+2}$$

$$1 - \frac{R_1^2 - R_0 R_2}{R_0^2 - R_1^2} = \frac{R_2 - R_0}{R_0^2 - R_1^2}$$

$$R_0^2 - R_1^2 - R_1^2 + R_0 R_2 = R_2 - R_0$$

$$R_0(R_0 + R_2) - 2R_1^2 = R_2 - R_0$$

$$R_0(R_0 + 1) - 2R_1^2 + R_2(R_0 - 1) = 0$$

$$\alpha_1 = R_1 \left(1 - R_2 + \frac{R_1^2 R_0}{R_0^2 - R_1^2} \right)$$

$$R_0^2 - R_1^2 - R_2 R_0^2 - R_2 R_1^2 + R_1^2 R_0$$

CM

$$\alpha_1 = \frac{-R_1(R_2 - R_0)}{R_0^2 - R_1^2}$$

$$\alpha_2 = \frac{R_1^2 + R_0 R_2}{R_0^2 - R_1^2}$$

classmate

Date _____
Page _____

~~$$LD \quad i=1 \quad E_0 = R_0$$~~

X

$$K_1 = \frac{R_1}{E_0} = \frac{R_1}{R_0}$$

$$\alpha'_1 = K_1 = \frac{R_1}{R_0}$$

$$E_1 = (1 - K_1^2) E_0 = \left(1 - \frac{R_1^2}{R_0^2}\right) R_0 \\ = \frac{R_0^2 - R_1^2}{R_0}$$

~~i=2 j=1~~

$$K_2 = \frac{R_2}{R_0} - \alpha'_1 R_1 \\ = \frac{R_2 R_0 - R_1^2}{R_0^2 - R_1^2}$$

$$\alpha''_2 = R_2 \checkmark$$

$$\alpha''_1 = \alpha'_1 - K_2 \alpha'_1$$

$$= \alpha'_1 (1 - K_2)$$

$$= \frac{R_1}{R_0} \left(1 - \frac{R_2 R_0 - R_1^2}{R_0^2 - R_1^2}\right)$$

$$= \frac{R_1}{R_0} \left(\frac{R_0^2 - R_0 R_2}{R_0^2 - R_1^2}\right)$$

$$\begin{bmatrix} 9 & 2 & 5 \\ 2 & 4 & 5 \\ 3 & 5 & 6 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2/7 & 1 & 0 \\ 3/7 & 35/24 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2/7 & 0 \\ 0 & 0 & 6 \end{bmatrix}$$

classmate
Date _____
Page _____

$$d_{kk} = a_{kk} - \sum_{j=1}^{k-1} L_{ij}^2 d_j$$

$$L_{ik} = \frac{1}{a_k} \left(a_{ik} - \sum_{j=1}^{k-1} L_{ij} L_{kj} d_j \right)$$

$i > k$

$$\frac{4}{7} + 2 = 4$$

$$x = 4 - \frac{4}{7}$$

Krg symmetric matrix

$$\begin{bmatrix} 7 & 2 & 3 \\ 2 & 4 & 5 \\ 3 & 5 & 6 \end{bmatrix}$$

$$y \frac{24}{7} = 5 - \frac{6}{7} \frac{9}{7} + \frac{5 \times 35}{24} + 6$$

$$y = \frac{35}{29}$$

$$\Rightarrow \begin{bmatrix} 1 & 0 & 0 \\ L_{21} & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix} \begin{bmatrix} d_{11} & 0 & 0 \\ 0 & d_{22} & 0 \\ 0 & 0 & d_{33} \end{bmatrix} \Rightarrow 4 - \frac{(4)}{7}$$

$$\Rightarrow \frac{28 - 16}{7} = \frac{12}{7}$$

$$\Rightarrow L_{21} = \frac{1}{a_{11}} \left(a_{21} - \sum_{j=1}^{1-1} \right) = \frac{2}{7} \quad d_{11} = a_{11} = \frac{7}{7}$$

$$L_{31} = \frac{1}{a_{11}} \left(a_{31} - \sum_{j=1}^{1-1} \right) \quad d_{22} = a_{22} - \sum_{j=1}^{2-1} L_{21}^2 d_j$$

$$= \frac{3}{7} \quad \Rightarrow a_{22} - \left(\frac{a_{21}}{a_{11}} \right)^2 a_{11}$$

$$L_{32} = \frac{7}{24} \left(3 - \sum_{j=1}^{3-1} L_{31} L_{21} d_{22} \right) = 4 - \frac{4}{7} = \frac{12}{7}$$

$$= \frac{7}{24} \left(3 - \frac{3 \times 2 \times 29}{7 \times 7} \right)$$

$$= \frac{21}{24} - \frac{3 \times 2 \times 2 \times 29}{7 \times 7 \times 7} = \frac{21}{24} - \frac{6}{49} = \frac{12}{7} \quad \text{X}$$

- Analysis for LPC

Durbin

Choleski

- i) Storage

Samples

N

N+p

Matrix

 $\propto p$ $\propto \frac{p^2}{2}$

proportional to

symmetric
matrix

Windowing

N

pre-store (extra

O (no clamping
req.)

Storage)

- ii) Multiplications

Windowing

N

O

Correlation

 $\propto N \cdot p$ $\propto N \cdot p$

- iii) Matrix Solution

 $\propto p^2$ $\propto p^3$

- Another benefit of Durbin

1) No divide by 0

 ~~$F = \frac{a}{b} = \frac{1}{1}$ (is it a transformation of signal?)~~

- What is clamping window?

→ we shall use a window called

hamming window.

(Weights)

$$w(m) = 0.54 - 0.46 \cos\left(\frac{2\pi m}{N-1}\right)$$

(Plot the

graph for this)

, $0 \leq m \leq N-1$

NOTE: precompute

, otherwise $w=1$ here N will be fixed

Cost a lot of computation!!

What happens when we apply hamming window?

→ the signal gets compressed

on both sides

Why necessary?

f

p

M-1 • FT

M-2 • LPC $a_i \ i=1 \dots p$

• What is the val of $a_i \ ?$ (can we generate $S \ ?$) : Regenerate speech

classmate

Date _____

Page _____

→ Take the signal of $N = 320$

Speech Representation

→ Apply the Hamming Window

→ find Auto Correlation

→ Apply Levinson - Durbin

→ Take first (p) as it is: generate the $(p+1)$ sample

↑ i.e. 13th sample

① { → keep the 13th (synthetic sample) in array

→ Skip the first sample, take next 12 → generate the 14th sample

→ Continue ① till $N = 320$ reached : (synthetic frame will

- listen to it ← be generated!)

- data too small: so nothing

so plot both the original ←

{ synthetic samples!!

→ check where varying



NOTE: use original data to generate signals &

nothing else !! ← if use synthetic → will be bad !!

M-3 0: Cepstral Coefficient

$C(m)$

↳ LPC Coefficient converted to this

$$C(0) = \log \sigma^2$$

where σ is the gain term of
the LPC mode

Gain = $R(0)$

$$C(m) = C_m + \sum_{k=1}^{m-1} \binom{k}{m} C_k a_{m-k} \quad 1 \leq m \leq p$$

$$= \sum_{k=m-p}^{m-1} \binom{k}{m} C_k a_{m-k}, \quad m > p$$

Why beneficial?

$C_i, i = 1, 2, \dots, N$

→ source is intact (which is previously in speech regeneration source not considered)

→ all a_i included

(we don't use $\{ \}$ ← we are restricting to 13 values
 $\Omega = 3/p$ as well) per frame //

The cepstral Coefficient have the following properties:

- (i) The lower order C-c are sensitive to the spectral slope
- (ii) The higher order C-c are less sensitive to noise, or noise like sounds

In order to diminish the sensitivity, a tapered window is used.

$$w(m) = \left\{ 1 + \frac{Q}{2} \sin\left(\frac{\pi m}{Q}\right) \right\}, \quad 1 \leq m \leq Q$$

(after calculating)

$w(m)$ apply it

We do, $C(m) w(m)$

$m = 1 \dots Q$

on $C(m)$ giving $P=12$ samples

not Q)

No distortion of

here we are only amplifying it

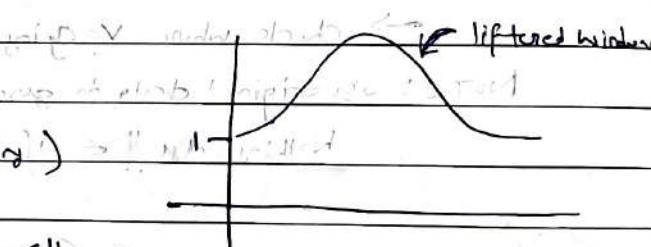
Floating point

Signal

→ called raised sine window

(12 bits)

This process is sometimes called lifting (opposite of filtering)



• here we allow some

part (e.g. low freq or high freq)

• but in lifting → (allows all as well as

amplify it) - No data loss

$$m > n \geq 1$$

$$\sum_{k=0}^{n-1} x(k) \hat{x}(k)$$

→ Mid-sum ←

• no more delay

→ T/F_{n-1}

$$q \leq m$$

$$\sum_{k=0}^{n-1} x(k) \hat{x}(k)$$

NTF: Matching

FIR

MOS

standard cell

SQ

LQ

interpolation is 1D →

- classmate
Date 03/03/25
Page _____
- feature set
Static & properties dynamic
- Instead of 320 samples \rightarrow $C_0 C_1 C_2 \dots C_p$ $\Delta C_0 \Delta C_1 \dots \Delta C_p$
 $P = 12$
 - # Orthogonal polynomial? How to do it \times
What is the use \checkmark
 - The derivation or Delta cepstrum is approximately given by,

$$\frac{d C_m(t)}{dt} = \Delta C_m(t) \simeq \gamma \sum_{k=-K}^K k \cdot C_m(t+k)$$

This is computed using orthogonal polynomial fit,
where γ is an appropriate normalisation factor $\in (2K+1)$ is the no. of frames (This is computed for over which the computation is performed, current frame & from back & front)

Usually K is treated as 3 (per analysis $i=7$)
 $\gamma = 0.375$ (empirically found)

The complete feature set also includes $13 \cdot \Delta^2 C_i$, or acceleration coefficients.

$$[C_0 \dots C_p] [\Delta C_0 \dots \Delta C_p] [\Delta^2 C_0 \dots \Delta^2 C_p] = 39 \text{ dimensional}$$

NOTE: feature vector for speech (size 39) \rightarrow 39 = industry set

Write a para on ΔC_p & $\Delta^2 C_p$

Q What can we do with cepstral coefficient?

Assignment 3
Extension \rightarrow voice recognition

How to differentiate b/w 2 speech

- Distance Measures

(Too easy) (i) Euclidean Distance

The first distance is the euclidean d

$$D_{CEP} = \sqrt{\sum_{i=1}^q \{C_i(t) - C_i(r)\}^2}$$

(Kepstral) P = 12 t → test data
r → reference data

Take 1 frame

- compute cepstral coefficients
- store it

Do the same for test data

→ Repeat for remaining frames

if distance is low \Rightarrow Same speech!

In this distance computation, equal weightage is given to each component.

It is easy to implement & fairly straight forward.

(Too good) (ii) Mahalanobis Distance

$$D_M = (\tilde{C}_t - \tilde{C}_r)^T \quad \begin{matrix} \leftarrow \text{transpose} \\ 1 \times q \end{matrix}$$

vector

• PC Mahalanobis : Write a

short note on him

(Mid sum)

$$V^{-1} (\tilde{C}_t - \tilde{C}_r) \quad \begin{matrix} \leftarrow \text{inverse} \\ q \times q \end{matrix} \quad \begin{matrix} \leftarrow \text{vector} \\ q \times 1 \end{matrix}$$

Where,

V = Covariance matrix of test feature vector (C_t 's)

Note: Inverse take (P_3^3)

Variance = how far or close

↑ ↑

proportion = $\frac{\text{high}}{\text{small}}$

i) it is extremely inaccurate (computing itself)

ii) difficulty in computing the inverse - (cubic) & $\underbrace{\text{underflow prob}}_{\text{wrong}}$

iii) $\sum_{i=1}^N \frac{1}{w(i)} = \begin{cases} \infty & \text{if diagonal} \\ \text{else} \rightarrow 0 \end{cases}$

(iii) Tokhura's distance

- Take ~~variance~~, $w(i) = 0$ (except few out)

- He suggested using only the diagonal elements of V ,

$$D_{T_{CEP}} = \sum_{i=1}^{N+1} w(i) [C_i(t) - C_i(r)]^2$$

where

$w(i)$ is the inverse of the variance of the
(computed) i th element
(for test data)

$$w(i) = \frac{1}{\sigma_i^2}, \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (u_i - \bar{u})^2$$

$$\bar{u} = \frac{1}{N} \sum_{i=1}^N u_i$$

Distinguishing speech based on a_t , using prev method not good ...

⇒ Distance Measure based on LPCs (a_t 's)

(i) Itakura Minimum Prediction Residual Distance

$$D_{I_{LPC}} = \frac{a_t' R_t a_t}{\|a_t\|^2} - 1$$

error

(here most optimal) $a_t' \rightarrow$ LPCs of ref = ref

⇒ if $R_t = R_p \Rightarrow$ error is $a_t' \rightarrow$ LPCs of test

(NOTE: $a_t' R_t a_t$ is also error $R_t \rightarrow$ correlation matrix
but of reference data)

(ii) Itakura Saito Distance

$$D_{IS} = \left(\frac{a_t' R_t a_t}{\|a_t\|^2} \right)^{\frac{1}{2}} + \log \left[\frac{G_r^2}{G_t^2} \right]$$

source component

$$G = \sqrt{a' R a}$$

⇒ Best sounding

vowels (reference)

Ci) (5)

take distance
from test data

T_{.....}

(iv) Steady frame

• for each frame calc

R_t, a_t, C_t

• Shift by ~80

• apply hamming window

• (long double)

Assignment Steps

(i) calc DC shift - with power
before recording

(ii) normalise - to ± 5000

(iii) recording & Succeeding frame
- Own Algo based on ZCR
& Energy

from all 5 ref ⇒ whichever min (is the same) ⇒ find recognition score

- NOTE: Don't write the distance calculation code

→ use excel sheet : Rot Cir, Ctr, & tomorrow's not
→ calc there

24/09/25

How to prove those 12 values are useful?

→ vowel recognition } - by calculating the distance!!

→ regenerate the waveform → Accuracy - 85%

LPC algorithm = backtracking search

Cepstral Coefficients

Cepstrum: The Fourier transform of a signal $x(n)$ is given by,

$$\text{gives } X(k) = \sum_{n=0}^{N-1} x(n) e^{-j2\pi k n} \quad 0 \leq k \leq N-1$$

fourier coefficients

Let, ~~number of samples~~

$$\hat{x}(k) = \log [X(k)] \quad - \text{ never try to take}$$

log[•]

to apply FFT⁻¹ on this → original signal P

∴ $\hat{x}(k) \leftarrow \text{real part of } X(k)$

for more detail CP to chapter

Cepstrum is defined as,

$$\hat{x}(n) = \frac{1}{N} \sum_{k=0}^{N-1} \hat{x}(k) * e^{j2\pi k n} \quad 0 \leq n \leq N-1$$

IDFT

well behaved

inverse discrete fourier transform

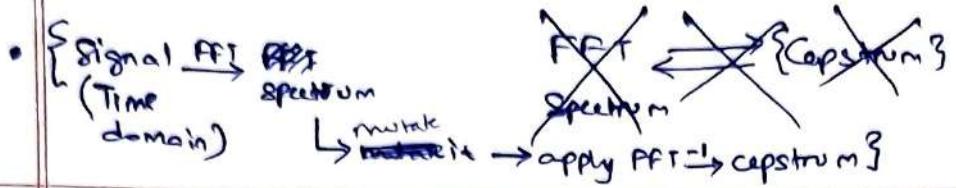
Cepstrum is defined by,

$$C(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log |X(e^{jw})| e^{jwn} dw$$

it's

analytic, etc. B.TP changing

- $\infty < n < \infty$



classmate

Date _____
Page _____

$X(k)$ approximately behaves like $X(K)$

Example 1 # raw signal : 10 kHz (1000 samples in 1 sec)
 each sample : 16 bits
 $\Rightarrow 160,000 \text{ bps}$

Capstral analysis $\leftarrow \{ \text{let's say } p = 10, 100 \text{ analysis points per sec} \}$
 16 bits data vector computation
 $\Rightarrow 16000 \text{ bps} - \text{data rate}$

NOTE: If we use original data $\Rightarrow 160000 \text{ bps}$

If we use capstral representation $\Rightarrow 16000 \text{ bps}$,
 data rate: i.e. 10:1

going to (i) \Rightarrow how to improve

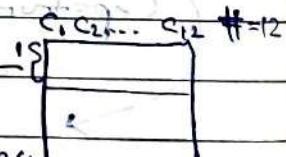
$$\pi = \frac{2\pi}{7} = \{ \text{Scalar quantization} \rightarrow 3.14 \text{ (approx)} \}$$

(representing decimal in small nos.): 7

$$[(x)]_{\text{qst}} = (x)$$

Vector Quantization

$\begin{bmatrix} C_1 \\ C_2 \\ \vdots \\ C_{12} \end{bmatrix}$ \rightarrow We would require a 'code book' of
 approximately 1024 entries (25 variants of 4 basic speech units)



Each row is called a

Code word / template?

size of codebook

= levels

{ # here recognition can be done again by comparing C.C & checking distance

Book: Grayson & (from) thy table \Rightarrow QI reporting associated

Grey for vector

Speech sound?

Quantization (either C1, C2) -

• Constituent of codebook

Pointed example 13

In this situation, there will be 100 vectors packed

$$\Rightarrow 1000 \text{ bits/sec} : \text{as each row by } 10 \text{ bits}$$

$\uparrow \quad \quad \quad \Rightarrow 10 \times 100$

Throughput achieved by

vector quantization

(no. of data bits
dealing)

hence \Rightarrow so codebook assigns 10 bits/vector

NOTE: here 1 vector = 1 frame
containing all c.c

Advantages:

- (i) less storage space is used
- (ii) reduced computation for determining similarity of spectral vectors
- (iii) discrete representation of speech signals

Drawbacks:

(i) spectral distortion is introduced $\{(\sum |E_i|)^2\}$ even

(ii) here quantization error is strictly greater than 0

(iii) storage of codebook vectors may not be non-trivial for large

sized codebooks (e.g. 10000)

Why \Rightarrow NOTE: most of important info in speech is in the spectral

CC reduce it to envelope (formant) not in every sample

$\frac{1}{10}$

CC analysis capture this envelope using just few coefficients

\uparrow per frame (1 rep)

comprised form representation

Defn: It was proposed by Gray in 1984, Makhoul et al 1985

Let \tilde{x} be R -dimensional vector whose components are real valued random values.

A vector \tilde{x} is mapped onto another R -dimensional vector \tilde{y}

$$\tilde{x} \rightarrow \tilde{y}$$

$$\text{where } \tilde{y} = q(\tilde{x})$$

That is \tilde{y} taken on one of a finite set of values,

$$y = \{\tilde{y}_i\} \quad i = 1, 2, \dots, K$$

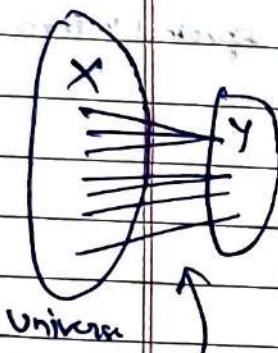
\uparrow no. of levels

The set of vectors \tilde{Y} is called the codebook & each \tilde{y}_i is called a codeword or a template.

The size K of the code book is referred to as the no. of levels.

When \tilde{x} is quantized to \tilde{y} , a quantization distortion measure,

$D(\tilde{x}, \tilde{y})$ can be defined between \tilde{x} & \tilde{y} .



The overall avg distortion is then represented by,

$$D = \lim_{M \rightarrow \infty} \frac{1}{M} \sum_{n=1}^M d(x(n), y(n))$$

A quantizer is said to be optimal (minimal distortion) if the overall distortion is minimized over all K level quantization.

→ To make such a codebook, the k -dimensional space of vector \tilde{x} is partitioned into K regions, i.e.,

$$\{C_i\}_{i=1, \dots, K}$$

with \tilde{y}_i being associated with each region C_i .

The quantizer then assigns the code vector \tilde{y}_i if \tilde{x} is in C_i ; that is, $q(\tilde{x}) = \tilde{y}_i$.

• Voronoi Diagram / Sections

→ just defn
← how
dist to
vry

Conditions necessary for optimality: {of a codevector being assigned}

i) nearest neighbour selection rule

$$q(\tilde{x}) = \tilde{y}_j \text{ iff } d(\tilde{x}, \tilde{y}_j) \leq d(\tilde{x}, \tilde{y}_i)$$

$$\forall i \neq j \quad \forall 1 \leq j \leq k$$

so quantize \tilde{x} to \tilde{y}_j

centre most point if there are no islands in the encoded structure

ii) Centroid Condition

\tilde{y}_i should be chosen to minimize the avg distortion in the region C_i , such a vector is called the centroid of the region C_i

If there are M_i codewectors in the region C_i then the centroid will be given by,

$$\tilde{y}_i = \frac{1}{M_i} \sum_{j=1}^{M_i} \tilde{x}_{ij} \quad \forall M_i \in C_i$$

(most frequent assignment after step 1)

Generalized Lloyd's Algorithm {k-means algorithm}

Let K be the no. of clusters {the K regions}

that is,

$$\{C_i\} \quad i = 1, 2, \dots, K \quad \text{C}_i \text{ universe have to be cut in } K \text{ regions}$$

Step 1: Initialisation

- Set $M = 0$ where m is the iterative index
- Choose a set of initial codereceters, $\{\tilde{y}_i(0), 1 \leq i \leq K\}$ using an adequate method

Step 2: Classification

Let $\{\tilde{x}_i(n)\}$

Take each $\tilde{x}_i(n)$ & put them into cluster $C_i(m)$ using the nearest neighbour selection rule, that is,

$$\forall i \in \{1, 2, \dots, K\} \quad \tilde{x}_i(n) \in C_i(m) \text{ iff } d(\tilde{x}_i(n), y_i(m)) \leq d(\tilde{x}_i(n), y_j(m)) \quad \forall j \neq i, 1 \leq j \leq K$$

Step 3: Code vector updating

$$m \leftarrow m + 1$$

Update the code vector of every cluster by computing the centroid of each cluster.

$$y_i(m) = \text{Centroid}(C_i(m)) \quad 1 \leq i \leq K$$

Compute the distortion measure, $D(m)$

Step 4: Termination?
 If, $D(m)$ at iteration m relative to $D(m-1)$, is below a certain threshold, stop, else go to step 2

Initial values of $y_i(0)$ and $C_i(0)$ are given

• What should be the starting centroids?

→ random vectors from universe - is this optimal? (No)

→ first k vectors in universe - not optimal as one type of sound

→ 3rd method {missed}

→ Judgment sampling - (?) {if we know the speech samples, we choose samples from each type of sound as our initial centroids}

K-Means Alternative - Starting point problem

• Linde-Buzo-Grey (LBG) Algorithm

→ no decision on Starting point {system takes care of it}

→ start with optimal centroids → then expand it by ensuring optimality

→ starting with code book of size = 1 is centroid of the whole universe

Split it into 2 centroids → suboptimal {run k-means again}

↓ to make it optimal

Split it again i.e. 4 centroids → {run k-means again}

It is also called - the binary split algorithm, or modified K-means algorithm.

The steps of the algorithm are as follows:

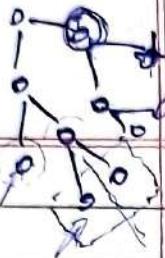
(i) Design a 1 vector codebook. This is the centroid of the entire training vectors

(ii) Double the size of the codebook by splitting each code vector \tilde{y}_n as follows:

$$\tilde{y}_n^+ = \tilde{y}_n (1 + \tilde{\epsilon})$$

$$\tilde{y}_n^- = \tilde{y}_n (1 - \tilde{\epsilon})$$



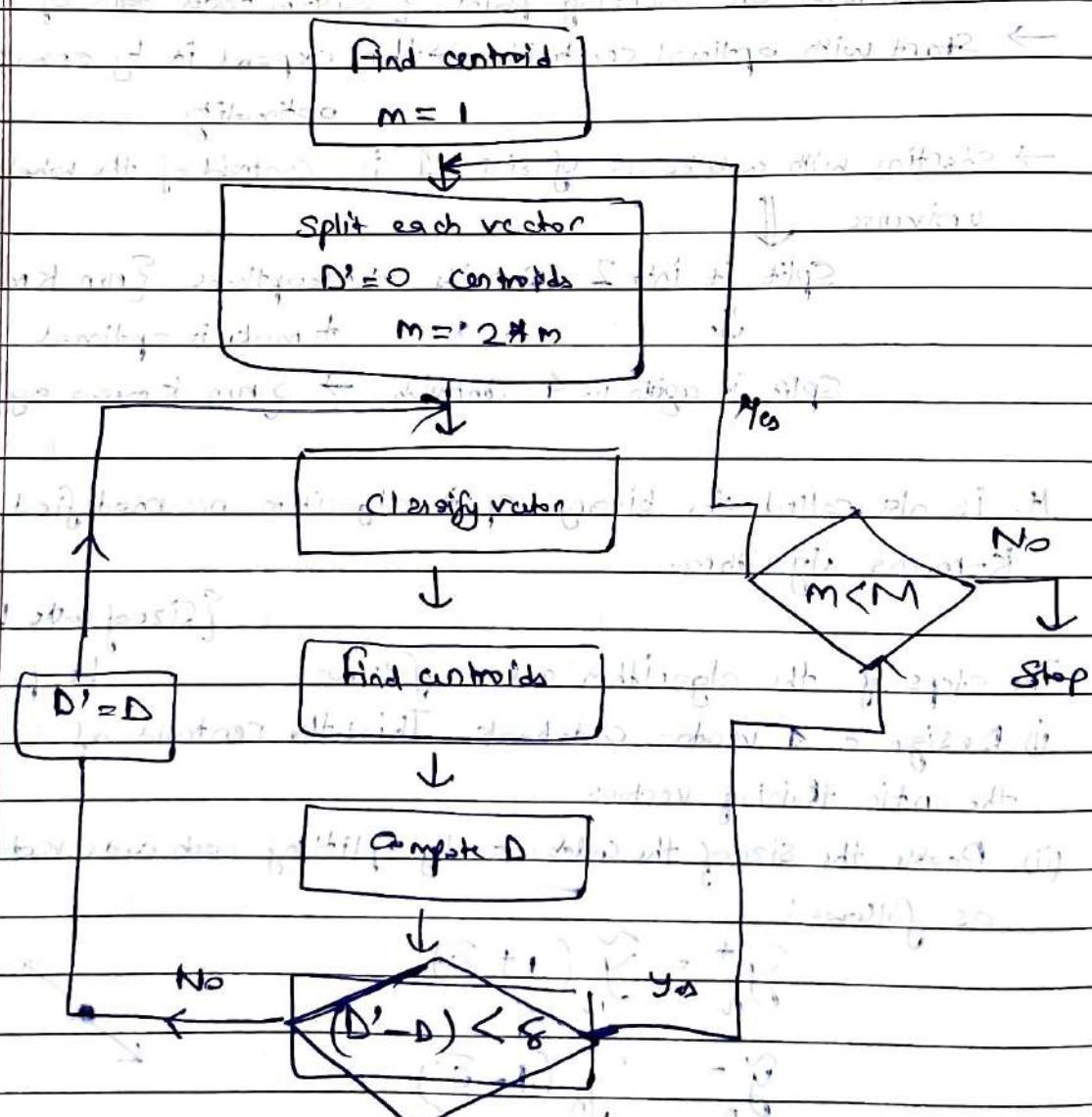


Where n varies from 1 to current size of the codebook
 $\epsilon \in$ is the splitting parameter. It is normally between
 $0.01 \leq \epsilon \leq 0.05$

$$0.01 \leq \epsilon \leq 0.05$$

(ii) use k-means algo to get the best set of centroids for the split codebooks i.e. codebooks of twice the size.

(iv) Iterate steps (ii) & (iii) until a codebook of size M (2^n) that is reached.



What is the weak point in this algo?

- The choice of epsilon (ϵ)
- All points go with 1 centroid & the other is with no points
 - here no new centroid can be computed for the other point as %

This is called empty cell problem

Can also occur in k-means algo? - but somehow settled

- Q How to solve this?
- Stop the algorithm for bigger codebook
 - throw away the vector which was having no points associated with it, and split the vector with high density?
 - run the ~~algo~~ # k-means again!
 - gather some more data when empty cell occurred
 - {increase the size}

- | | | | |
|---|------------------|--|-----------------|
| (i) Since at each split, | # General Points | (ii) | # 12 dimensions |
| size of codebook decreases | randomize which | ϵ wrt each dimension | |
| \Rightarrow Make ϵ logarithmic | | is many for $d_1 \rightarrow$ both pos | |
| & that it also dec per | | $d_2 \rightarrow$ poss neg | |
| split | | $d_3 \rightarrow$ both neg | |

Final set of centroids $\{C_1, C_2, \dots, C_m\}$ & std deviations of the clusters

$\{$ data drives the algo $\} \Rightarrow$ Centroid of data chosen of data

LBG Algorithm

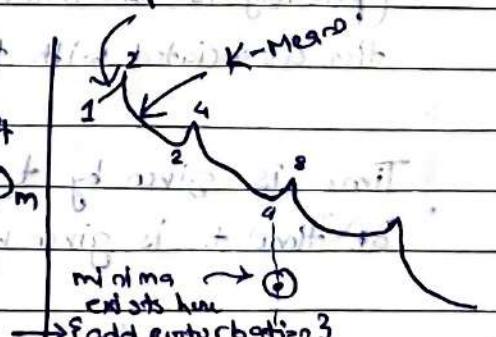
\hookrightarrow ϵ optimally chosen

\Rightarrow else more # iterations?

\Rightarrow omission error & rounding off

Dm

K-Means



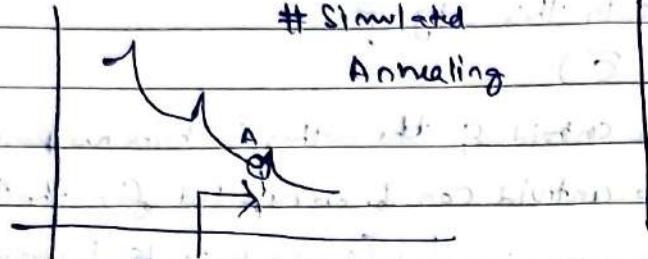
iteration

iteration

LBG

It helps in improving codebooks

Simulated Annealing



Annealing of Glass

Why does it happen?

Physical process - chemical process change

→ To get a point over here, repeat

Local optima

Known at A

→ by adding perturbation { 16th & 17th bit change }

Only one find and then move the next word (iii)

Next, next word and so on till the last is done

→ Make a codebook } Input

of size R

using K-Means

LBG

⇒ Match C_i with the codebooks

gives the index value of

the codebook entry which

gives the match perfectly!

→ per frame, 1 noisy frame to 1000 (ii)

For example, from file 3 match done to 1000

Statistical Analysis (contd.)

Discrete Markov Processes

Consider a system which may be described to be in one of any distinct N states, S_i ($i = 1, 2, \dots, N$)

At regular intervals the system undergoes a change of state, (possibly to the same state) according to a set of probabilities associated with the state.

Time is given by $t = 1, 2, 3, \dots$ and the actual state at time t is given by q_t ,

(i) Initial

Emitted state (iii)

Output (ii)

and (i)

$$P(q_t = s_j \mid q_{t-1} = s_i, q_{t-2} = s_k, \dots)$$

$$= P(q_t = s_j \mid q_{t-1} = s_i)$$

$$= a_{ij}$$

1st order markov

(S1S) \rightarrow (S2S) \rightarrow (S3S) \rightarrow ... process

$$(S1S) + (S2S) + \dots \quad (1 \leq i, j \leq N)$$

Properties:

(i) a_{ij} is a probability, hence $0 \leq a_{ij} \leq 1$

$$(ii) \sum_{j=1}^N a_{ij} = 1 \quad | \quad 0 \leq i \leq N$$

Let us consider a weather forecasting problem,

3 states (i) Raining

(ii) Cloudy

(iii) Sunny

Q What's the probability according to the given model that the weather for 8 consecutive days will be sunny - sunny - sunny - rainy - cloudy - rainy - sunny - cloudy - sunny?

Today is sunny & the model is.

$$A = \{a_{ij}\} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

Today

Cloudy

Rainy

Sunny

(i) $i+1 = S_1$ (ii) $i+1 = S_2$ (iii) $i+1 = S_3$

$\{\text{row sum} = 1\} \leftarrow \{\text{is stochastic matrix}\}$

This type of matrix

matrix

Probability it is sunny today

$$P(3) = 1$$

$$P(\sim \text{Mod})$$

$$= P(3, 3, 3, 1, 1, 3, 2, 3 \neq \text{Mod})$$

Probability it is sunny today is 3, $P(3) * P(3/3) * P(3/3) * P(1/3)$

Yesterday it was 3 $* P(1/1) * P(3/1) * P(2/3)$

This probability gives

$$* P(3/2)$$

$$= 1 * 0.8 * 0.8 * 0.1 * 0.4 * 0.3$$

$$* 0.1 * 0.2$$

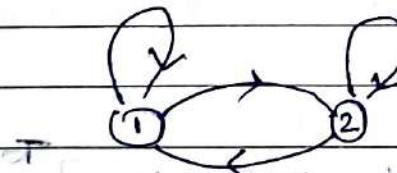
$$= 1.536 * 10^{-4}$$

08/10/25

- if 1 coin, $P(H) = \frac{H}{H+T}$

$$P(T) = T = 1 - P(H)$$

- for 2-coins, it will need which coins picked first info



probabilities (transition matrix)

$$(ii) A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

(ii) Characteristic probabilities $B = \begin{cases} P(H_1) = ? & P(H_2) = ? \\ P(T_1) = 1 - P(H_1) & P(T_2) = 1 - P(H_2) \end{cases}$

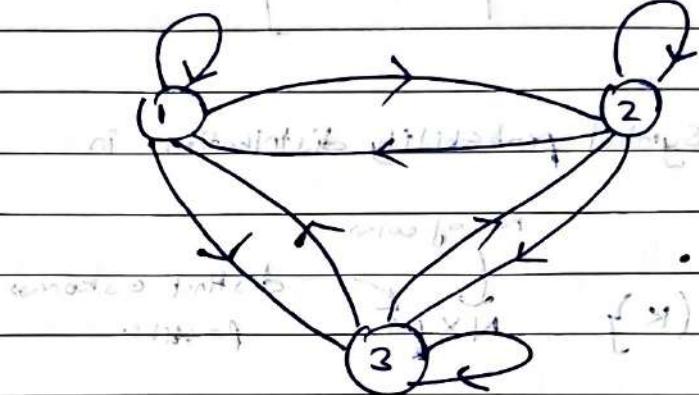
$$\begin{bmatrix} 0 & 1.0 \\ 1.0 & 0 \end{bmatrix}$$

Transition probability matrix

(iii) Initial probabilities (at $t=1$) $\Pi = (\Pi_1, \Pi_2)$

What if coins increased to 3?

$$\text{Initial Probabilities} \quad \Pi = [P_1 \ P_2 \ P_3] = [0.5 \ 0.25 \ 0.25]$$



• Fully connected graph

→ This can be increased to no. of coins (N) if needed

Hidden Markov Models (HMMs)

It is a doubly embedded Stochastic process with an underlying stochastic process that is not observable, (that is hidden).

(remm-
ber
wordy
word)

It can be observed through another set of stochastic processes that produce the seq of observations!

- i) How many wins? (taken)
- ii) Nature of coins (own)

Example: {either by ghost game or 2 more that's following}

Elements of an HMM → no. of coins

① N - no. of states, $S = \{S_1, S_2, S_3, \dots, S_N\}$

② At time $t \rightarrow q_t$ (state at t) $q_t \in S$

③ M - no. of distinct observation symbols per state

i.e. codebook size = power of 2 (64, 128, 256, ...)

for us: 128

1 - industry: 128 //
(CCITT std)

if diagonal values high \Rightarrow inertia of state?

③ A — The state transition probability distribution

Stochastic matrix $\rightarrow = [a_{ij}]$

$$= P[q_{t+1} = s_j \mid q_t = s_i] \quad 1 \leq i, j \leq N$$

④ Observation symbol probability distribution in state s_j

Stochastic matrix $\rightarrow B = \{b_j(k)\}_{N \times M}$ no. of coins \downarrow
 where $b_j(k)$ is the probability of observing k at time t provided $q_t = s_j$
 distinct outcomes possible

kind of saying

If C_1 is picked up $b_j(k) = P(V_k \text{ at time } t \mid q_t = s_j)$

What is the id of $b_j(k)$, answer to all of both using stochastic

Probability of getting

getting M/T which is \downarrow distinct symbols having $1 \leq k \leq M$

• Row sum of A & B are 1 for each row

2 stochastic matrix

initial state distribution is no more than 1 unit \Rightarrow ignore

⑤ Initial state distribution

$\Pi = \{\Pi_i\}$ where i belongs to $\{1, 2, 3, \dots, N\}$

$\Pi_i = P(q_1 = s_i) \quad i = 1, 2, 3, \dots, N$

state at time point 1

$$\sum_{i=1}^N \Pi_i = 1$$

We can write an HMM as,

$$\lambda = (A, B, \pi) \text{ Where } N \& M \text{ are implicit}$$

A is $N \times N$ matrix $\{$ already known $\}$

$$A = [a_{ij}]_{N \times N} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1N} \\ \vdots & & & \vdots \\ a_{N1} & a_{N2} & \dots & a_{NN} \end{bmatrix}$$

$$B = [b_j(k)]_{N \times M} = \begin{bmatrix} b_1(1) & b_1(2) & b_1(3) & \dots & b_1(M) \\ \vdots & & & & \vdots \\ b_N(1) & b_N(2) & \dots & \dots & b_N(M) \end{bmatrix}$$

Memory

09/10/25

Memory

- $N - A$

- $M - B$

- $N \times M - \text{matrix } A$

- $N \times M - \text{matrix } B$

- $N \in \Pi$

Ball and Urn Model

i) Picks an urn out of N based on prob

ii) Picks a ball from the urn it is with min. \rightarrow $N \in \Pi$

iii) Report the color of ball - {the observation}

NOTE: Π - initial prob to pick an urn

A - given current state, what prob to move to next?

B - what observations persist? - have 4 colour balls

per urn

Objective is: Determine how many no. of urns present

3 Problems of HMM

Pb1: Evaluation problem (Scoring Problem)

Given a model λ^i and a sequence of observations, how to{ Compute that the observation sequence was produced from
the model? Also it's called the scoring problem.

{ Kind of

Compute
pb3

How well a given model matches an observation sequence

$$\max_i P[\mathcal{O} / \lambda^i] \quad \begin{matrix} \text{: probability of obser-} \\ \text{vation Seq provided mo-} \\ \text{del } \lambda^i \end{matrix}$$

Pb2: Uncovering' Problem

{ We try to decipher the "hidden" part of the model that is
to find the "correct" state sequence.

Eas in

Which
State

the out-

Come
Come

from3

We try to use some optimality criterions to solve this pb

as best as possible.

Pb3: Resimulation Problem

We try to optimize the model parameters to best describe how
a given observation sequence comes about.The observation sequence used to adjust the model parameters
is called the training sequence.

Training sequence for training and adjustment of the model

- $W = 10$

Pb3 will make the model : Pb models

Pb2 will help in evaluating the model - testing : helps in redesigning the model

Pb1 will help us determine from which model it came!

↳ for deployment : only $\frac{1}{3}$ rd of the whole system

but in M1, we used the whole model during

deployment/training;

bigger size compared to M1/M2

Solution to Pb1

{forward Process}

M1 : We want to calculate probability of different things at each

$P(O_t | \lambda)$ probability of Observation seq given λ

(where O_t is the observation seq (going from) $t=1$ to T)

$$(O_1, O_2, \dots, O_T) = P(O_1, O_2, \dots, O_T | \lambda) = P(O_1) P(O_2 | O_1) \dots P(O_T | O_{T-1})$$

$M = 32$

depends on no. of frames

In order to find a solution we enumerate all possible state seq of length T , total no. of states $= N^T$

$T - \text{obs}$

Let form of the state be, s_t \rightarrow (x_1, s_t) \rightarrow $N - \text{states}$

$$\tilde{q} = (q_1, q_2, \dots, q_T) \quad \left\{ \begin{array}{l} \rightarrow \text{for a single observation} \\ (O_t) \end{array} \right.$$

$\rightarrow (x_1, q_1) \rightarrow \text{initial state}$

$$Z = (x_1, q_1)^T =$$

Now the probability of O given the state seq \tilde{q} can be written as,

$$P(O | \tilde{q}, \lambda) = \prod_{t=1}^T P(O_t | q_t, \lambda)$$

(here O_t should be independent)

$$P(\tilde{Q} / q, \lambda) = b_{q_1}(O_1) \cdot b_{q_2}(O_2) \cdots b_{q_T}(O_T)$$

Also, the probability of such a state seq \tilde{q} , is given by,

$$P(\tilde{q} / \lambda) = \prod_{t=1}^T a_{q_t, q_{t+1}} \cdots a_{q_T, q_1}$$

Now the joint probability of $\tilde{Q} \in \tilde{q}$,

$$\rightarrow P(\tilde{Q}, \tilde{q} / \lambda) = P(\tilde{Q} / \tilde{q}, \lambda) \cdot P(\tilde{q} / \lambda)$$

This is $= \prod_{t=1}^T b_{q_t}(O_t) a_{q_t, q_{t+1}} \cdots a_{q_T, q_1} b_{q_1}(O_T)$

Finally,

$P(\tilde{Q} / \lambda)$ is obtained by summing the joint probability
 $P(\tilde{Q}, \tilde{q} / \lambda)$ over all possible state seq \tilde{q} ,

$$\Rightarrow P(\tilde{Q} / \lambda) = \sum_{\text{all } \tilde{q}} P(\tilde{Q} / \tilde{q}, \lambda) \cdot P(\tilde{q} / \lambda) \cdots *$$

$$(\because \prod_{t=1}^T a_{q_t, q_{t+1}}) \cdot \prod_{t=1}^T b_{q_t}(O_t) = T(\text{Time complexity}) = 2T \cdot N^T = O(N^T)$$

Now how to do this operation.

$$N = 5$$

$$T = 100$$

It is not feasible for a real time implementation.

M2: The forward Procedure

Consider a forward variable $\alpha_t(i)$ as follows

$$\alpha_t(i) = P [O_1, O_2, O_3, \dots, O_t, q_t = S_i | \lambda] \\ \text{partial forward probability} \quad i = 1, \dots, N \\ t = 1, \dots, T$$

That is the probability of the partial observation sequence,

$O = (O_1, O_2, \dots, O_t)$, the state S_i at time t given the model λ .

16/10/25

We can solve for $\alpha_t(i)$ inductively as follows:

(i) initialisation

$$\alpha_1(i) = \pi_i b_i(O_1) \quad i = 1, 2, 3, \dots, N$$

(ii) induction

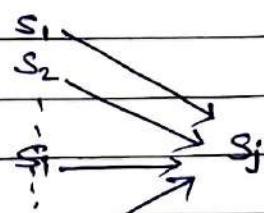
$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1})$$

$$1 \leq t \leq T-1$$

$$1 \leq j \leq N$$

(iii) termination

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$



Time complexity: $O(N^2 T)$ / per model

$$\Rightarrow \text{for } 10 \text{ models} \rightarrow 10 N^2 T$$

The Backward Procedure

Consider a backward variable,

$$\beta_t(i) = P(O_{t+1}, O_{t+2}, \dots; O_T | \underline{q_t} = s_{i,d})$$

~~It has already
been observed~~

9) Initialization:

$\beta_7(i) = 1$, $1 \leq i \leq N$.

ii) Innovation

$$\beta_t(i) = \sum_{j=1}^N (a_{ij} b_j(O_{t+1})) \beta_{t+1}(j)$$

just for interpretation

↳ probability of partial observation say if not then

for $Q_{t+1}, Q_{t+2}, \dots, Q_T$ at time t , $q_t = s_t$

Soln 2:

$$\gamma_t(i) = P(q_t = s_i / \Omega, \lambda)$$

That is the prob of being in state s_i at time t given observation seq Ω , model λ

Using the defn of fwd & bwd probabilities, we may

note,

$$\gamma_t(i) = P(q_t = s_i / \Omega, \lambda)$$

$$= P(q_t = s_i, \Omega / \lambda)$$

$$= P(\Omega / \lambda)$$

$$= \frac{\alpha_t(i) \beta_t(i)}{P(\Omega / \lambda)}$$

$$= \alpha_t(i) \beta_t(i)$$

$$\sum_{i=1}^N \alpha_t(i) \beta_t(i), \quad t=1, 2, 3, \dots, T$$

$\alpha_t(i)$ takes into account the prob of the observation seq $\Omega_1, \Omega_2, \dots, \Omega_t$ & state s_i at time t

(x, 0/b)

Wanted to find maxima of p(x)

arg

$B_t(i)$ takes into account observation O_{t+1}, O_{t+2}, \dots
given state s_t at time t

$$\text{Hence } \sum_{i=1}^N \gamma_t(i) = 1$$

Using $\gamma_t(i)$ we can solve for the individually
most likely state q_t at time t as follows:

$$q_t = \arg \max_{1 \leq i \leq N} [\gamma_t(i)], \quad 1 \leq t \leq T$$

↑ 2D matrix

Here transitions are not considered

e.g. $a_{15} = 0$ but we may get 5th
state as very likely

Some a_{ij} 's may be zero for some $i \neq j$ resulting
in an optimal state sequence which will be invalid.

One possible sol' is to modify the optimality criterion
for example: Solve for State sequence that maximises
the expected no. of correct pairs of states (q_t, q_{t+1})
or triple of states (q_t, q_{t+1}, q_{t+2})

$$\text{Hence } \sum_{i=1}^N p_i(i) = 1$$

~~Given state is of form~~
 ~~$B^t(i)$ takes int acc to observation O^t~~

We need to maximise p_b of,

$P(q/\Omega, \lambda)$ which is equivalent to maximising p_b of,

$$P(q, \Omega/\lambda)$$

$$P(\Omega/\lambda)$$

A formal technique exists & is based on dp method.

Viterbi Algorithm

To find the single best state seq q ,

$$Q = \{q_1, q_2, \dots, q_T\}$$

for a given observation seq Ω ,

$$\Omega = \{\Omega_1, \Omega_2, \dots, \Omega_T\}$$

We need to define the quantity,

Scalar quantity $\delta_t(i) \rightarrow \delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P[q_1, q_2, \dots, q_t = s_t, q, \Omega_1, \Omega_2, \dots, \Omega_t / \lambda]$

That is $\delta_t(i)$ is the best score (highest p_b) along a single path, at time t which accounts for the

$$\left[\frac{1}{(1-g)} \right]^{N-t} = \max_{1 \leq i \leq N} * P = \text{Max probability of } g \text{ given state}$$

~~prob <-- w
prob too <-- 1, 1. - w ... n~~

3. Transition: \rightarrow

The first t observations ξ_t ends in state S_t .

By induction we have,

$$\delta_{t+1}(j) = \max_i [\delta_t(i) a_{ij}] b_j(O_{t+1})$$

To actually retrieve the state s_{t+1} we need to keep track of the argument that maximised $\delta_{t+1}(j)$
i.e. i for each t & j

We do this via the array $\psi_t(j) \rightarrow \{ \text{keeps track of } i \text{ which are optimal} \}$

The complete procedure for finding the best state s_{t+1}
can be now stated as,

1. Initialisation: $\delta_1(i) = \Pi_i b_i(O_1), 1 \leq i \leq N$
 $\psi_1(i) = \emptyset$

2. Recursion: $\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(O_t)$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}]$$

$$\begin{cases} 2 \leq t \leq T \\ 1 \leq j \leq N \end{cases}$$

$$(1+2) f_0 \left[f_1 \circ (1)^T g \right] \max = (1)^{1+2} g$$

By induction we have,

The first + observation in end is state 1.

Next row good is λ if $P \downarrow \rightarrow \lambda$ not good
 $P \uparrow \rightarrow \lambda$ good

3. Termination:

Max probability \hat{P}^* = $\max_{1 \leq i \leq N} [\delta_T(i)]$

\rightarrow Seq & observation

(i.e. of complete path)* $\hat{Q}_T^* = \arg\max_{1 \leq i \leq N} [\delta_T(i)]$

most optimal state

at the last state / time T

4. State sequence (Path) backtracking :

\hat{Q}_t^* optimal state seq $\{ \hat{Q}_t^* = \psi_{t+1}(\hat{Q}_{t+1}^*) \}$

& it's probability available

$$t = T-1, T-2, \dots, 2, 1$$

This is similar to fwd procedure except for the maximisation step.

Input:

Observation seq

model

Output:

Globally optimal state seq.

Solution to Problem 3 : The Reestimation Problem

In this problem, we want to determine a set of parameters that is,

$$\lambda = (A, B, \pi)$$

to maximise the probability of occurrence of the observation seq given the model

There is no known way of solving this problem explicitly.
we can choose a set,

$$\lambda = (A, B, \pi) \text{ such that,}$$

$P(O|\lambda)$ probability of observation seq given λ is locally optimised

We use the Baum - Welch method (or equivalently the Expectation Modification Method Algorithm or Gradient Techniques (Hill Climbing Problem))



it can be used in any situation: where temporal data is used

The Solution: To solve this problem using iterative update & improvement of HMM parameters we define a variable,

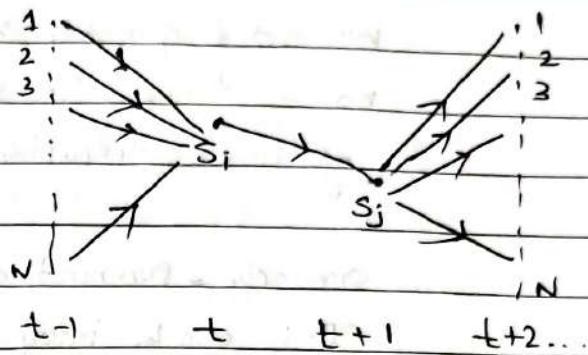
$\xi_t(i, j)$ {the probability of being in state s_i at time t and state s_j at time $t+1$ given the model λ & the observation seq i.e.,}

$$\xi_t(i, j) = P[q_t = s_i, q_{t+1} = s_j / O, \lambda]$$

Using the dyn of fwd S, bkd pb, we can write,

$\xi_t(i, j)$ as follows,

$$\cancel{\xi_t(i, j) =}$$



$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(\Omega/\lambda)}$$

$$= \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

$$\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

The denominator makes $\xi_t(i, j)$ into a probability measure.

$$\gamma_t(i) = \text{product of } \xi_t(i, j) \text{ where } j = 1, 2, \dots, N$$

Pb of observation seq ending in State s_i at time t given $\Omega \xi \lambda$

We can now relate $\gamma_t(i)$ & $\xi_t(i, j)$ as follows,

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j) \quad t = 1, 2, \dots, T$$

(# Paradigm Shift)

If we sum $\gamma_t(i)$ over t , the time index from 1 to $T-1$,

We get a quantity which can be interpreted as the expected no. of times the state S_i is visited or expected no. of times transitions are made from S_i .

Similarly, summation of $\gamma_t(i,j)$ over time t , from 1 to $T-1$ can be interpreted as,

expected no. of transitions from state S_i to S_j

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected no. of transitions from state } S_i \text{ to } S_j$$

$$\sum_{t=1}^{T-1} \gamma_t(i,j) = \text{expected no. of transitions from state } S_i \text{ to } S_j$$

Using this formulae, we give a reasonable set of re-estimation formulae for Π , A , B as follows:

$$\bar{\Pi}_i = \text{Expected freq in state } S_i \text{ at time } t=1$$

$$= \gamma_1(i), \quad i=1, 2, \dots, N$$

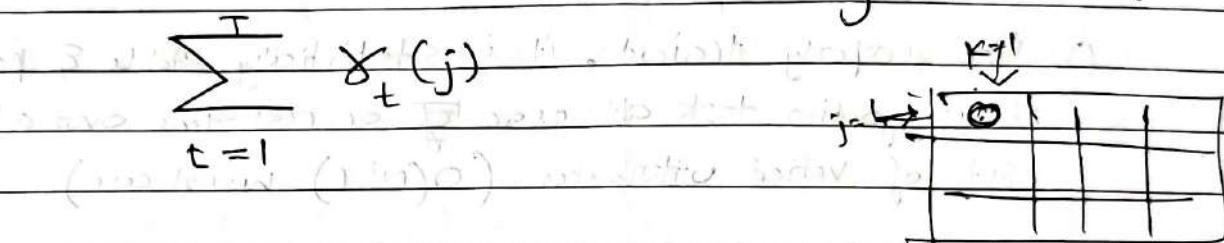
$$\bar{A}_{ij} = \frac{\text{expected no. of transitions from state } S_i \text{ to } S_j}{\text{expected no. of transitions from state } S_i}$$

$$= \frac{\sum_{t=1}^{T-1} \gamma_t(i,j)}{\sum_{t=1}^{T-1} \gamma_t(i)}, \quad i, j = 1, 2, \dots, N$$

$\bar{b}_{jk}(\leftarrow)$ = expected no. of times in state s_j and observing symbol v_k

expected no. of times in state s_j

$$= \sum_{t=1}^T \gamma_t(j) \quad (v_k: output of vector quantized) \\ t \text{ s.t } O_t = v_k \quad K = 1, \dots, M \\ j = 1, \dots, N$$



If we define $\lambda = (\Lambda, \beta, \pi)$ as the current model, $t=1, \dots, T$

model $\tilde{\lambda}$ compute new model $\bar{\lambda} = (\bar{\Lambda}, \bar{\beta}, \bar{\pi})$ $O_t = v_k$

by using the above set of formulae it is proved
that, either

(i) $\bar{\lambda}$ defines a critical point of the likelihood function
in which case, $\lambda = \bar{\lambda}$

(ii) $\bar{\lambda}$ is more likely than model λ is the sense that,

$$P(\Omega/\bar{\lambda}) \geq P(\Omega/\lambda)$$

What should be starting λ ?

(i) Random

(ii) Set it as: $\begin{bmatrix} 1 \\ N \end{bmatrix}$ for Λ

$\begin{bmatrix} 1 \\ M \end{bmatrix}$ for β

$\begin{bmatrix} 1 \\ N \end{bmatrix}$ for π

NOTE: until $\lambda = \bar{\lambda}$,

repeat the process

i.e if $P(\Omega/\bar{\lambda}) \geq P(\Omega/\lambda)$

set $\lambda = \bar{\lambda}$ & repeat

no information
model
(unbiased)

fully connected state machine

The Three problems in HMM - 4 algorithms

Advantages Of HMM

- It is very accurate
- Time alignment is done automatically & optimally
- It is mathematically tractable & is flexible in modelling real life phenomena
- Once properly trained, it is statistically stable & performs the recognition task at near ~~of~~ or real time over a large set of varied utterances ($O(N^2 T)$ worst case)
- It can be used in speaker-independent case

Disadvantages

- The no. of states have to be carefully selected as once finalized it cannot be changed for that model
- Training of HMMs require huge amt of training data over a large no. of speakers or utterances
- Choice of (starting model) is critical for the proper convergence of the model

Digit Recognition System

$N=15$

$N=7-9$ for sentence

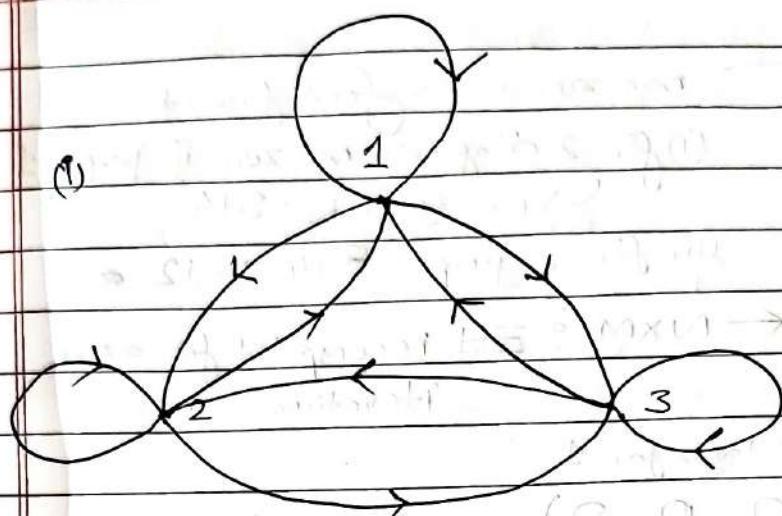
& phrase based detection

- Usually a no. of starting HMMs are used & the model that gives the best performance (given by $P(\lambda | \Omega)$) is chosen

Types of HMM

Model

- Based on state topology



(Ergodic Model)

Fully Connected Model

Challenge: converging ← • Most flexible that can
 model identification

→ n built
 (for image
 recognition
 system)

(ii)

$$A = \begin{bmatrix} [1/3] \end{bmatrix}_{3 \times 3} \quad \pi = (1/3, 1/3, 1/3)$$

$$B = \begin{bmatrix} [1/4] \end{bmatrix}_{4 \times 4}$$

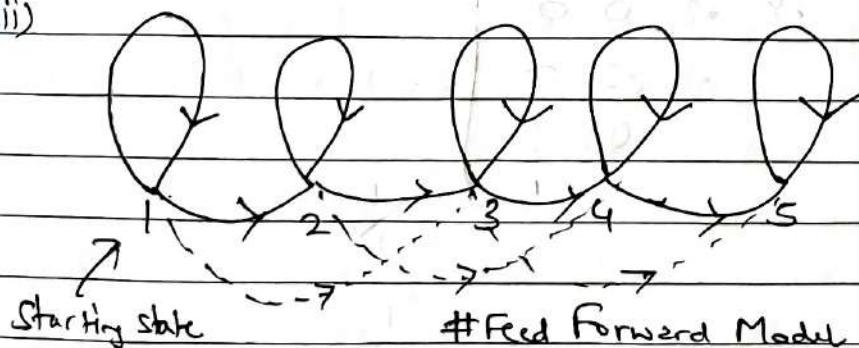
$$N=3$$

$$M=4$$

(Unbiased Model)

= input? only observation seq! which is vector quantization

(iii)

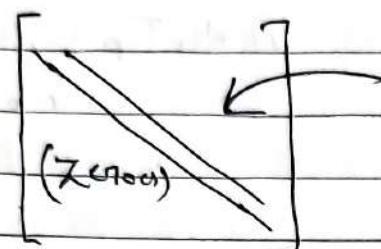

 Relevant in
 speech

Feed Forward Model (Moving forward)

(Baldi's Model) - (for 1 step moving fwd)

NOTE: Beta
 ↑
 If 2 steps skipped ⇒ 2 skipping F.F.M

Both will be in different state due to model constraint



non zero as feed forward

(i) for 2 diagonals non zero if jumps 1

i.e. g status: 5+4

(ii) for 2 jumps: $5+4+3=12$

$B \leftarrow N \times M$: ~~all~~ recomputed for every iterations

Choosing first Model

(Bakis Model) as it starts from 1

$$\pi = (1, 0, 0, 0, 0)$$

$\Rightarrow P(Q/A) \uparrow$

unbiased Model

$$A = \begin{bmatrix} 0.5 & 0.5 & 0 & 0 & 0 \\ 0 & 0.5 & 0.5 & 0 & 0 \\ 0 & 0 & 0.5 & 0.5 & 0 \\ 0 & 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$[Q/A] = A$$

$$b_{ij} = \frac{1}{M} \sum_{k=1}^M \pi_{ik} j$$

Adding bias (Intentional test) (Intuition Model)

$$A = \begin{bmatrix} 0.8 & 0.2 & 0 & 0 & 0 \\ 0.8 & 0.2 & 0 & 0 & 0 \\ 0.8 & 0.2 & 0 & 0 & 0 \\ 0.8 & 0.2 & 0 & 0 & 0 \\ 0.8 & 0.2 & 0 & 0 & 1 \end{bmatrix}$$

Architecture of Speech Recognition System

(# Fill in the blanks
of source code)!

Soln 1: $\underline{\lambda}, \lambda \rightarrow P(\underline{\lambda}, \lambda)$

Soln 2: $\underline{\lambda}, \lambda \rightarrow p^*, q_t^*, t=1, \dots, T$

Soln 3: $\underline{\lambda}, \lambda \rightarrow \bar{\lambda} \quad P(\underline{\lambda}/\bar{\lambda}) > P(\underline{\lambda}, \lambda)$

(# ALL VARS GLOBAL)

Step 1: Make a vocabulary: w words - 10 digits

Step 2: # Training

- First utterance of zero is the

Training seq

$M = 32, N = 5$

Code book size

Weak law of large numbers

- Req 29 samples

at least + follow normal distribution

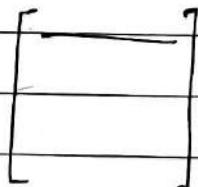
to build the model

Step 3: Use trial / starting (λ)

⇒ keep 30 for

→ use the solution to problem 2 for initial λ for training and→ print q^* (p^* will be very small) 10 for testing

(1 2 3 4 5 5 ... 5)

Step 4: Compute Soln to problem 1 - $\lambda \underline{\lambda}, p$ Step 5: use Soln to prob 3 → get new $\lambda \rightarrow \bar{\lambda}$

go back to Step 3 for estimating how

good the model is

 $M \times 12$
=Check if better (print q^*, p^*)write old model to a text file; $\lambda = \bar{\lambda}$ Repeat Step 2: 30 models !!! (30 λ for digits)

after we see convergence

in Step 3 //

• for 30 models, generate 1 avg model for O

{ → we get stochastic model - as all stochastic matrices

use this starting model & repeat

→ new avg find (run it 3-times)



get the avg → i.e. final Model

run 1 more time ←

→ pick one of the 30 models

which has the highest !!

probability

repeat for 10 digits — (20 models) final!

Text: FOC → run on 20 models → get 20 prob



picks highest &

report !!

(observation: P# ←)

should increase) # Adjustment of op matrix

• After computation of A, check if row sum is 1

(if not) add the fraction value which makes it 1
to higher pb

• for B matrix, to get non zero values in every row

{ on every iteration → take it as float value

if val < 10^{-20} → make it as 10^{-20}

→ not stochastic matrix anymore

⇒ do row sum

⇒ remove the extra value from highest pb !!

o For 30 models, generate 1 avg model for O

(\rightarrow we get stochastic model - as all stochastic matrices)

use this ^{as} starting model & repeat

\rightarrow new avg find (run it 3-times)



get the avg \rightarrow i.e. final model

run 1 more time \leftarrow

\rightarrow pick one of the 30 models

which has the highest!!

probability

repeat for 10 digits — (20 models) final

Test: rec \rightarrow run on 20 models \rightarrow get 20 prob

picks highest &

(observation: P^*) \leftarrow

report!!

Should increase) # Adjustment of o/p matrix

• After computation of A, check if row sum is 1

(if not) \rightarrow add the fraction value which makes it 1
to higher pb

• for B matrix, to get non zero values in every row

at every iteration \rightarrow take it as floor value

if val $< 10^{-20}$ \rightarrow mark it as 10^{-20}

\downarrow not stochastic matrix anymore

\Rightarrow do row sum

\Rightarrow remove the extra value from highest pb!!

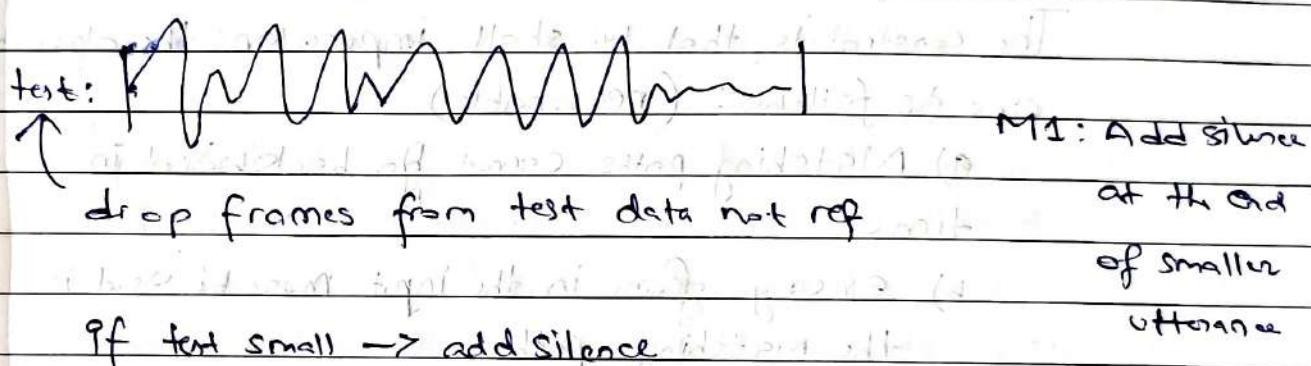
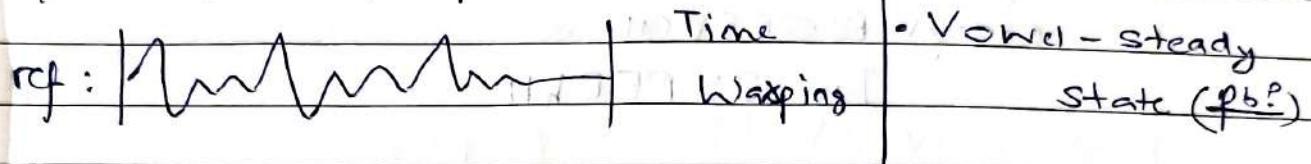
Word Matching

Text Based - DTIF } we can use
TTIF }

Speech: not linear

padding can be done to make both same

M2 # Matching 2 unequal speech utterances? Linear



(?) Struts points - vowels - matching in both & adjust the frames

Dynamic Time Warping (DTW) - 1970's - 80's - Vishnu Atal

It involves two concepts:

- computation of the features
- distances

Some form of metric has to be used in order to obtain a matched path. There are 2 types of distances:

(a) local distance: it is a computational distance between ~~between~~ a feature in one signal with the feature of another signal

→ Taxicab / Manhattan

etc.

(b) Global distance: it is the overall computational distance between an entire signal & another signal of possibly different length

2 algorithms

(i) Symmetric DTW

example .Re: SPEECH

Tat: SPEECH

The constraints that we shall impose on the algorithm are as follows: (reasonable).

a) Matching paths cannot go backward in time.

b) every frame in the input must be used in the matching path

c) local distance scores are combined by adding them to give a global score

o We will have a matrix, Time-Time Matrix

X-axis : input pred curr-61

y-axis: ref

→ Choose from

$$(i-1, j-1) \otimes (i-1, j) \Rightarrow (i, j-1)$$

$\rightarrow D_{ij}$) is the global distance

update (i, j) and local distance

at (i,j) is given by

$$1 - \text{Res}_i \neq j+1$$

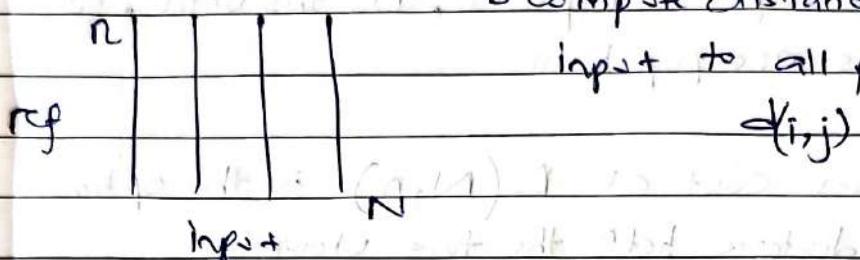
- no jumps allowed
- no skipping of cells

- Movement shown by arrows

$$D(i,j) = \min [D(i-1, j-1), D(i-1, j), D(i, j-1)] + d(i,j) \quad *$$

Distance Matrix ~~(N X N)~~ ($N \times n$) - input to the algorithm

compute distance from all point of
input to all point of ref



The algorithm to compute the least global cost

- Load only first 2 columns of TT-Matrix
 $\hookrightarrow n \times 2$ points

(i) We need to start at $C(1,1) (1,1)$ and go upto $D(N,n)$

Step 1 (ii) Compute. $D(1,1) = d(1,1)$

(iii) populate predecessor column : $D(j, 1) = D(i, 1)$
 $\quad \quad \quad + d(i, 1)$
 $\quad \quad \quad i > 2$

$$D(1, j) = D(1, j-1) + d(1, j) \quad j \geq 2$$

now,

for $D(2, j)$ = first fill $D(2, 1)$

$$D(2, j) = \min (D(1, j), D(1, j)) \quad (*)$$

Store curr in pred & load the next col in curr
repeat at 2 onwards

Step(3)

Step (2) populate current column Q_1 , calculate

at (i, j) the local distance = $d(ij) + \min_{\text{global}} \text{cost at either } (i-1, j-1) \text{ or } (i-1, j) \text{ or } (i, j-1)$

Step (B) Corr column is now assigned to pred coln
Eg, the steps are repeated. Do this until all columns are populated

Step (5) The global cost at $D(N,n)$ is the optimal least distance b/w the two utterances

(ii) Asymmetrical DTH shapes the voltage plots.

There are 3 main points allotted

(we can \leftarrow { a skip 1 jump } in ref } reason of
here many! { no skip on the ~~ref~~ input is allowed } } asymmetrical
for sentence
recap)

Not Computable		X	X				
		X	X				
		X	X				
		X	X				
		X	✓				
		(X)	(X)				
				✓			

$$(1,2)(i,j) \cdot (i,j) = (i,j)(1,2) = (i,j)$$

$$-\left((i+1) \cdot 2 + (i(i+1) \cdot 2) \right) \sin \theta = (i, 2) A$$

and also from the top? \Rightarrow big or no \Rightarrow

$$\alpha_{t+1}(j) = \sum_{i=1}^N \alpha_t(i) a_{ij} b_j(O_{t+1})$$

$$\beta_t(j) = \sum_{i=1}^N a_{ij} b_j(O_t) \beta_{t+1}(j)$$

$$\delta_t(i) = \max_j (\delta_{t-1}(j) a_{ji}) b_i(O_t)$$

$$\psi_t(i) = \max_j (\delta_{t-1}(j) a_{ji})$$

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

$$p^* = \max_i \delta_T(i)$$

$$q^*_T = \arg \max_i \delta_T(i)$$

$$\pi_t = \frac{\alpha_t(i) \beta_t(i)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)}$$

$$q^*_t = \psi_{t+1}(q^*_{t+1})$$

$$b_j(k) = \sum_{t=1}^{T-1} \frac{\gamma_t(j)}{Q_t} = R$$

$$\xi_t(i, j) = \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

$$\sum \sum$$

$$\pi_i = \gamma_1(i)$$

$$a_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$