



Lecture 01 (11 Aug)

Speech : legal sequence of legal sounds produced by humans. understandable by other humans

Broad goal of speech processing

↳ Human-machine interaction as natural as possible

ASR → automatic Speech Recognition
↳ speech → text

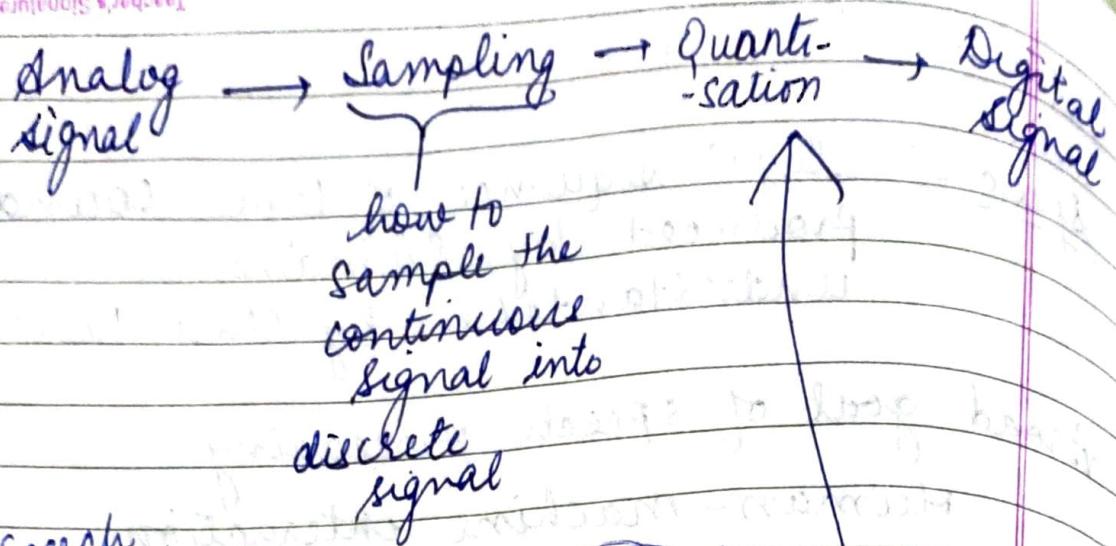
Applications of Speech processing

- ① ASR
- ② Language recognition
- ③ Speaker recognition
- ④ Speech synthesis — Text to speech (TTS)
e.g. GPS

⑤ Speech analysis
↳ Event detection
Speech to features
useful extract features

⑥ Speech Coding : since the wireless channel bandwidth is limited

so we



2 FM, WHY? to avoid aliasing (loss of info)

max frequency

$\text{FM} = 4 \text{ KHz}$

Compressed speech

each sample represented using 8 bits

if $\approx 8 \text{ kHz samples} \times 8 \text{ bits/sample}$

\Rightarrow bitrate 64 Kbits/sec

$\Rightarrow 64 \text{ Kbps bandwidth}$

but since speech has a lot of redundancy we can reduce the bitrate ~~bandwidth~~ to ($\text{eg. } 12.2 / 13.3 \text{ Kbps}$)

| and then GSM

reconstruct

back from 12.2 to 64

$30 + 12$

42×8^4

28×10^3

10^5

this is called speech coding

(@)

Now, we can accommodate 5 users instead of 1 — because we reduce bitrate to $1/5^{\text{th}}$.

5G uses Enhanced Voice Service (EVS)
Up to 3G/4G — AMR coder

Adapt to multiple coders
codes the speech signal depending on the quality of speech
If it is very imp. frame it uses high bitrate.

Bitrate — (5.7 to 12.2 kbps)

↓
unwanted speech

→ Speech Coding uses LP analysis

⑦ Speech Enhancement

↪ reduce noise

3 kinds of noise

a. Bg noise

b. Reverberation

c. multispeaker noise

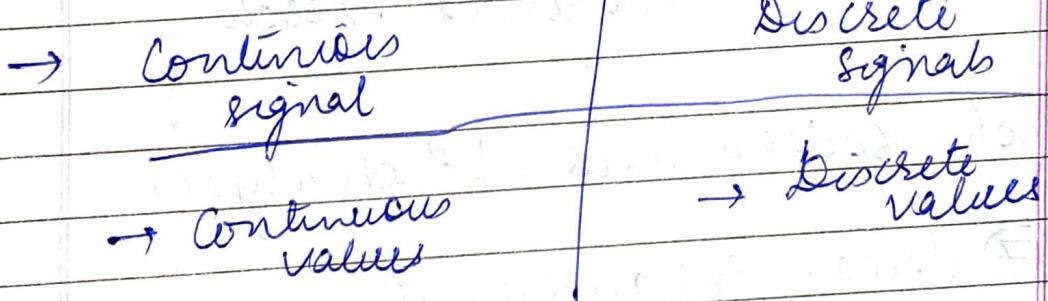
$f_s = 30 \text{ kHz}$

- Speech recognition - HMM
(baseline)
- Speaker recognition - Variants of gaussian models
(baseline)

Now, some NN is baseline
eg. ASR - RNN or Seq2Seq
Speaker - DNN + attention / multiheaded attention

Lecture 02 (14 Aug)

→ signals - expressed as funcn of independent variable



Analog to Digital

① Sampling (Analog to discrete)

a. Analog signal \times uniform impulse \rightarrow discrete signal

→ Sampling rate matters when reconstructing the signal back

→ Sampling $f_s > \text{high pass } f_1 \text{ to } \infty$

lowpass sampling

0 to max frequency
(f_m)

bandpass sampling

$f_1 \text{ to } f_2$

sampling rate = f_s

$$f_s \geq 2 f_m$$

→ ~~Q~~ time domain $\xrightarrow{\text{FT}}$ frequency domain

→ a. spectrum in frequency domain \times spectrum of impulse

Spectrum becomes repetitive with f_s

so if $f_s < 2f_m$ then aliasing happens and hence signal lose \star
PTO X2

and if $f_s \geq 2f_m$ even though they repeat, they won't touch each other

Nyquist Criteria

$$f_s \geq 2 f_m$$

sampling frequency ≥ 2 (higher frequency content)

A continuous time signal can be completely represented in its samples and recovered back if the sampling frequency is twice of highest frequency

In speech we take,
min. sampling rate = 8 kHz
because we assume our sounds have frequency upto zero to 4 kHz.

so max freq = 4 kHz (f_m)

so $f_s \geq$ as 8 kHz minimum (f_s)

Some sounds have frequencies upto 8 kHz so ideal $f_s = 16$ kHz
but due to bandwidth limit often 8 kHz is used.

HD quality sound use high sampling rate
 i.e. more samples are used to represent the signal.

② Quantisation (discrete to digital)

no. of bits used to represent a sample

if more no. of bits is used -
 i.e. more bitrate - the quantisation error reduces

more Bitrate = high sampling rate

and
 more no. of bits

information is preserved hence quality is good.

we need low quality also to fit to a given bandwidth.

LTI (Linear Time Invariant)

Linear - homogeneity & superposition
Time - invariant

Homogeneity

$$a x_1 + b x_2 \rightarrow a y_1 + b y_2$$

Superposition

$$x_1 \rightarrow y_1$$

$$x_2 \rightarrow y_2$$

$$x_1 + x_2 \rightarrow y_1 + y_2$$

Impulse response - output of impulse

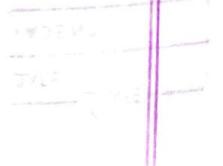
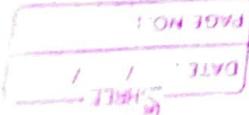
~~Output~~ Impulse response for LTI is convolution
of input \times impulse response.

- ⑪ Discrete time \rightarrow convolution sum
- Continuous time \rightarrow convolution integral

Fourier series - for periodic

Fourier Transform - non periodic &
periodic

\rightarrow To know characteristics of system/signal
we multiply it by an impulse
response (train of impulses).





Aliasing

when aliasing occurs — i.e. when there is an overlap ~~in~~ in signal low frequencies will get affected which are the voiced regions.

② In discrete domain

$$\text{output} = \text{conv}(\text{input}, \text{system response})$$

\downarrow

$$= x[n] * h[n]$$

SSP

Class : 21st Aug : Lecture 04

- voiced consonant
- unvoiced consonant
- vowels

	energy	duration	harmonic uniformity random beat
lower	low	high	unvoiced
lowest	low	highest	voiced
highest	highest	lowest	vowels

→ Formants — in similar range for vowels

→ spectrum section — FFT of the spectrum

→ Linear Prediction

→ LPC — smoothing of spectrum

→ in FFT — peak picking is difficult while it isn't so in LPC.

→ LPA fails in unvoiced consonants due to random (non-uniformity) & in aspirated & other sound but since speech is dominated by

→ vowels → vocal fold vibration (

rate: Pitch

VC &
vowels
it is
feasible

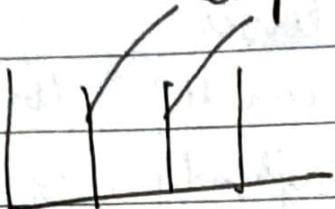
7 EGG Electroglottograph

○xension in waveform.

Teacher's Signature

→ Impulse - open

Regular Intervals - open & close
Impulse like (epoch locations)



→ Diff b/w impulses - pitch period

→ Impulse train ← assuming that voiced sounds are like impulse train

→ Glottal closure instant / instant of significant excitation

→ In diff EGG - the negative peaks are epoch locations

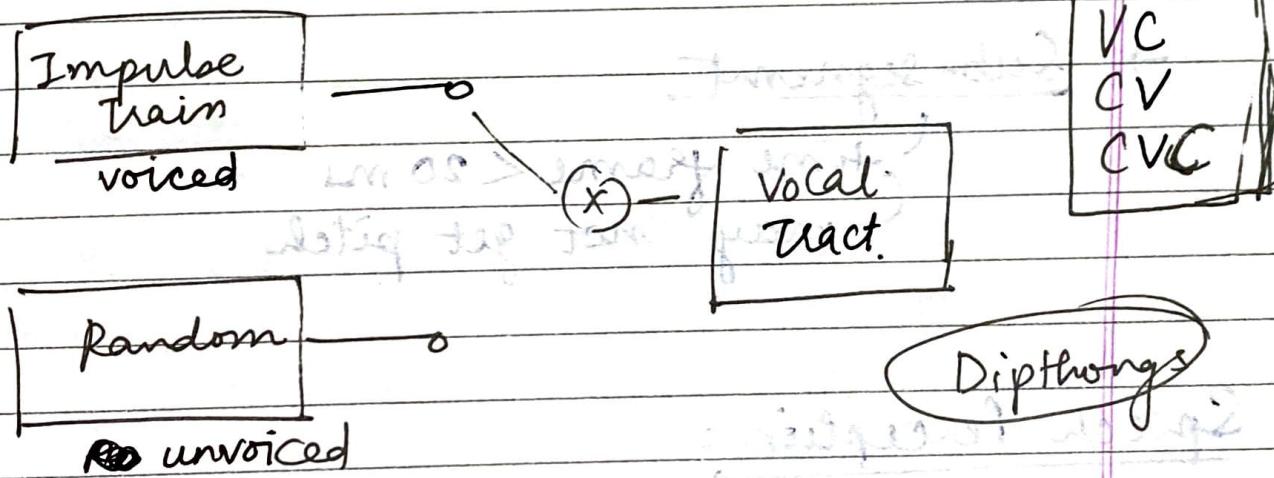
Epoch

→ LP - predict current sample because there is periodicity in picture



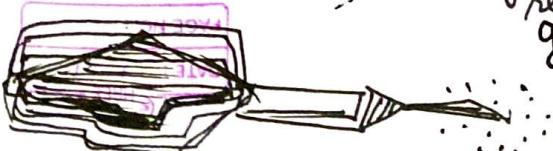
- Frame size - 20 ms - 2 to 3 pitch cycles
- $\{P = 8 \text{ or } 10\}$ in LP.
- z-transform — Study.
- Residual — excitation
- all pole filter because if we invert the "forward filtering" equation — there are all poles.

→



- Plosive — MOA — if phone is blocking the air

- Spectral details — f_0 } white formants
- Formant tracking } f_1
- Formant tracking } f_2
- yellow red green } mark range (Pg 13 Table 5)



Lecture 05 (25th Aug)Speech analysis

extract features to build

speech
systemsmost speech sys
use segmental

what features to use to automate
phonetic transcription? And how to
get these features?

Ans.

- a. Information from frequency bands
- b. frequency bands
- b. spectral energy bands

~~(b) Related Energy~~

Speech analysis techniques - depending on frame size

related to vocal tract
use of spectral features or frequency related features

① Segmental analysis

(around 20-50 ms)

→ 20ms corresponds to 2-3 pitch cycle
& signal won't change much i.e.
assume that stationary signal

related to vocal tract
use of spectral features or frequency related features

② Subsegmental (around 5ms)

frame size < 20ms → may have redundancy (i.e. excitation source feature have redundancy compared to vocal tract feature)

use → extract excitation source related features like epoch-related features

related to vocal tract
use of spectral features or frequency related features

③ supra-segmental (imp to human but not to speech systems)

frame size > 100ms

use of spectral features or frequency related features

extract energy patterns
duration patterns
pitch patterns } natural variation

Class : 28 Aug (Lecture 06)

→ Segmental technique is popular ($TF = 20 \text{ ms}$)
 wth spectral features

→ Supra segment
prosodic features
 (pitch ~~in energy contours~~,
 energy & duration of wave)
 ↓
 related to
 human voice/emotion

→ Sub-segment excitation features
 (time frame $< 20 \text{ ms}$)
 (may not get pitch)

Speech Perception:

as Voiced / Unvoiced detection

outer ear
middle ear

inner ear

neuro muscular commands
 converts to neuro commands

→ speech understanding also depends on language / context understanding

Voiced/Unvoiced detection — is used in most speech systems + first step - before extracting segment features

- Applications :
- speech encoding
 - voiced activity detection — activate VAD
 - AMR (Adaptive Multi-Rate) — Alex
 - Emotion/Language/Speaker identification

- zero-crossing is less in voiced.
- energy is more in voiced
-

Voice activity

-
- more the bit rate — better the coding

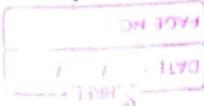
(AMR) →

Adapt to multistate code — supports multiple bit rates

4.75 to 12.2 — most popular
 (lowest bit rate) (highest BR)

→ Frame size — 20ms

- if voiced — use highest bit rate
- if silence — use lowest bit rate



→ AMR uses VAD

Q Why removing unvoiced / silence part
 → Extract features ^{only} from voiced speech
 while performing speaker recog / emotn
 etc.?

→ But for speech recognition — silence and unvoiced also important

Ans. Only voiced regions have information about the excitation source / vocal tract characteristics.

Silence will be common across all speakers — it is redundant informat.

In silence & unvoiced speech the vocal tract system is not active. And they may confuse the system very prosody info isn't useful for speech systems. Prosody may change for the same speaker.

Frame-level / Block processing of Speech

Used mostly for speech processing

Voiced/Unvoiced Detect

Basic approaches \rightarrow pros & cons and formula:

① Easiest: Energy (for one (100 ms) frame or 3-4 (multiple) frames)

$\rightarrow 20/50/100 \text{ ms}$

\rightarrow Compute energy or log energy

\rightarrow put threshold \leftarrow greater than this → voiced
0 → silence
less → unvoiced

Pros

\rightarrow easiest & simple

\rightarrow if clean speech,
 \times if we put hard

easy

Cons

\rightarrow noise can't it affect will give wrong results

\rightarrow we aren't exploiting any characteristics of voiced/unvoiced dynamics

\rightarrow if putting a threshold is difficult \because we cannot put a static threshold.

\rightarrow some voiced consonants may get removed if not enough data is available

② Zero crossing - no. of times ~~on~~^{signal} crosses zero
 ↓
 more for unvoiced

Cons : if noise ~~isnt~~, it may look like noise & noise will also have more zero crossing

Pros : - putting up a threshold may be slightly easier but fails in case of noise
 - easy & simple

→ Spectral energy is more for lower frequencies. for voiced & vice versa.

③ ~~Spectral energy~~
 → ~~Lower Energy~~
 → ~~Total band energy~~
 more for voiced
 less for unvoiced

- more efficient till now
- Spectral energy may sometimes have an issue of threshold

(4) Periodicity / Quasi-stationary

voiced — periodicity \oplus more
unvoiced — " less

Q How to check whether signal in frame is periodic or not?

Ans. Auto correlation - formulat $(\sum_{n=1}^N s(n)s(n))$

Correlation of two signals is high when the signals are same.

In Auto correlation ~~formulat~~,

if we apply auto correlation ~~formulat~~ on signal (on 20 ms frame)

then the correlation coefficient -

{2nd peak & 1st peak} \rightarrow almost similar

\because because i.e. correlatⁿ b/w the two is voiced

of periodicity

but in unvoiced — there will be no peaks in auto correlatⁿ

deviatⁿ b/w first peak & second peak will be much higher because no correlatⁿ b/w signal

→ can use spectral energy in frequency domain
and autocorrelation in time domain.

→ Excitation for unvoiced is like white noise so ideally autocorrelatⁿ should be an impulse funcⁿ for that, but there are few informatⁿ that is left in unvoiced — hence we won't get an impulse.

→ Sometimes even noise can have periodicity so method ③ & ④ may fail there — so we need better methods for voiced / unvoiced detection.

How to handle then?

→ LP analysis

→ glottal closure instants

Q Can we use epoch information in voiced / unvoiced detection?

Ans Locations & strengths of epochs will be periodic & highly good in case of voiced speech

& random (not pattern) for unvoiced speech.

→ also a good approach even in case of noise (i.e. no epochs in noise)

Lecture 07 (1st September)

- Pitch: Rate of vocal fold vibration
- Impulse like excitations in voiced sounds
- The periodicity of those impulses is pitch period. $= 1/T \approx 5\text{ms to } 9\text{ms}$.
- Children > Female > Male
(7-9) (5-7)
- No concept of pitch in ~~silence or~~ unvoiced region.

① Approaches to detect pitch from speech waveform

① Zero crossing — very erroneous
— won't work if noise is dominant

② Autocorrelation

③ Epoch locations

→ successive differences of epochs will be pitch period

④ Cepstrum

Q. How can we use autocorrelation for pitch detection?

Ans a. Segment size = ~~20ms - 50ms~~ (Sampling rate
per ms) \downarrow
~~8 samples~~

b. Difference first

Distance b/w two successive peaks (dominant)
after autocorrelation is applied is the pitch period.

Q. Why does it give pitch?

Ans. When the moving frame comes in same configuration i.e. when the cycles match — the correlation will be higher.

Cons: getting second dominant peak can be difficult ~~&~~ in case of large no. of spurious peaks.
if the segment size is not taken carefully

- We get pitch only in voiced relation
- You can put a threshold on the ratio of first & second dominant peaks to determine voiced/unvoiced.

Ratio is less — voiced
u more — unvoiced

Short Term Fourier Transform Analysis

Need for STFT

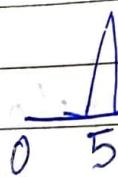
Q Signal $\in x\text{Hz}$, ~~FT~~ will give

Ans an impulse at ω . i.e. ~~it~~ it indicates the frequency of the signal.
So, the Fourier transform spectrum gives the frequency of the signal.

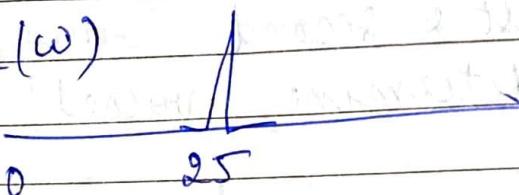
$$\begin{cases} \text{eg} \\ \text{a. } x_1(t) = \cos(2\pi \cdot 5 \cdot t) & - 5\text{Hz} \\ \text{b. } x_2(t) = \cos(2\pi \cdot 25 \cdot t) & - 25\text{Hz} \\ \text{c. } x_3(t) = \cos(2\pi \cdot 50 \cdot t) & - 50\text{Hz} \\ \text{d. } x_4(t) = x_1(t) + x_2(t) + x_3(t) \end{cases}$$

Stationary
signal
↓
do not
change w.r.t
time

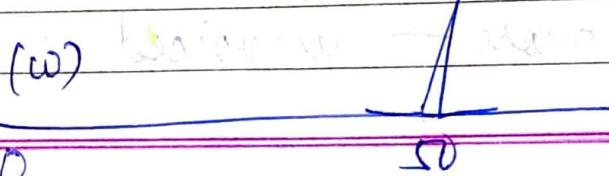
a! $\bullet X_1(\omega) = \text{FT}(x_1(t))$



b! $X_2(\omega)$



c! $X_3(\omega)$



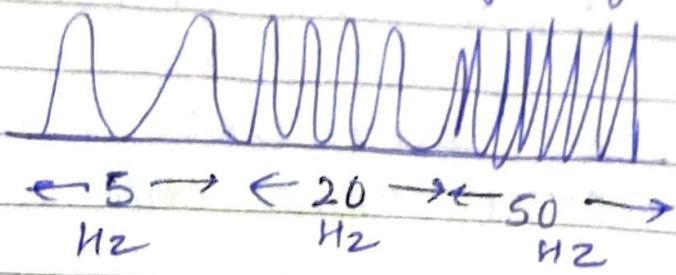
d! $X_4(\omega)$



non-sta

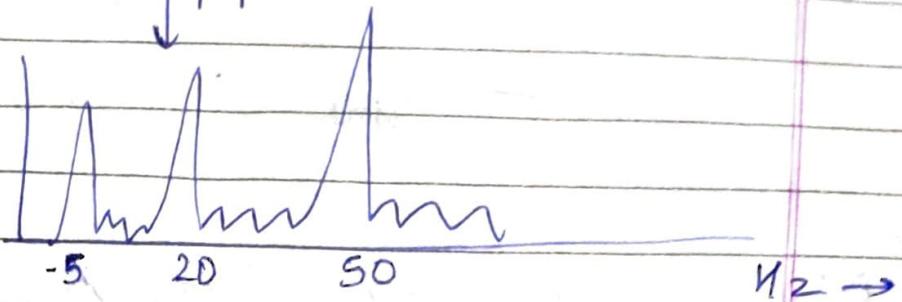
→ if we take a non-stationary signal

e. $x_5(t)$:



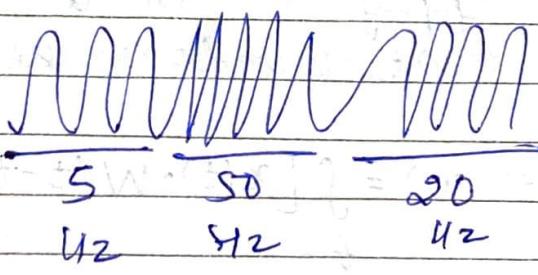
↓ FT

$X_5(\omega) =$



So, FT gives us knowledge of what frequencies exist but no information about where these frequencies are located in time.

So, for eg. if we change $x_5(t)$ to $x_6(t) \rightarrow$



↓ FT

will give same

$X_5(\omega)$ spectrum

Drawback:

So, we cannot classify/distinguish b/w non-stationary signals using FT. Since for both $x_5(t)$ & $x_6(t)$ it gives same spectrum.

FT

L converts time to frequency

Pros: all frequency info is known

Cons: loses time info

signal:

Amp

time

FT of amp

frequency

So, FT works well if signal is stationary
and speech is non-stationary
so I need both time & frequency
analysis

→ So, we use STFT — FT for short segments
and find spectrogram i.e. time vs
frequency plot/spectrum.

→ In STFT, we apply FT on a windowed
signal and that window keeps on moving

$$\text{STFT}_x^w(t, w) = \int_{-\infty}^{\infty} [x(t) \cdot W(t-t')] e^{-j\omega t} dt$$

→ For $x_5(t)$ there will be only one
magnitude spectrum ~~if~~ if we apply FT.

But, in STFT (t vs Hz) — the spectrogram

will vary because it depends on
the window size

Q Ideal length of window = 20ms (160 samples)

Q How to choose the ideal window length?

Ans:

Since FT works well for stationary signals — within the window the signal should be stationary. So we take 20ms window size — since we assume the signal is stationary.

Advantage STFT

Disadv. of STFT

→ a kind of uncertainty/tradeoff is present b/w time & frequency, since the two are opposite to each other. i.e. we have time resolution (when time info is seen well) and frequency resolution (where frequencies well seen)

→ For non-stationary signals we require Joint time-frequency analysis

→ Output of STFT is spectrogram and it is not unique — changes in

① window size

② window shape ↗ Hanning (bell shape)

Hanning (bell shape),

can be changed in spectrum section in wave surfer

→ The output of STFT that comes as a spectrogram is the convolution of signal spectrum and window spectrum because we are windowing the signal.

→ Window Shape

→ bell shape window → these windows have less sidelobe leakage. Hence, the window won't affect much of the signal spectrum.

→ but if we take rectangle window the spectrum of signal gets affected.

Q If we take a small window, which one is good - time resolution or frequency resolution and why?

Ans. Time Resolution → when we take narrow window

Frequency Resolutn - wide window because we will get all frequencies

The spectrogram flattens for high window size and does not capture any time variations

$$2^x = 32768 \quad x =$$

and for small window size — the spectrogram shows more time variations vs freq.

Q To get frequency info in FT what transform should be used?

~~At point DFT~~ n-DFT (Discrete FT)
where 'n' should be in powers of 2.

Q To compute DFT fast we use FFT.

So. n should be a value in powers of 2 greater than the no. of samples

→ n is the number of points where frequency values are computed and those points are joined

if $n <$ no. of samples — we lose information

so for 20 ms frame \approx 8 kHz sampling rate FFT points should be 256 point or 512 point DFT.

DFT
(min)

256 64 128 21

(DFT is applying FT on discrete signals)
and FFT is used to perform fast DFT.

→ Spectrogram is 3D plot -
 Time vs Frequency vs Amplitude
 (Intensity)

→ We can also use spectrogram for
 Voiced / Unvoiced detection -

ratio of high freq
low freq.

more for unvoiced (\therefore dominated by high freq)
 less for voiced (\therefore dominated by low freq)

→ Formants tell us where spectrogram is dark

→ So in 512 point DFT all points other than the 160 samples are masked zero — this is called zero-padding

Q How does changing the window size affect the joint time frequency resolution

if large window — same spectrum.

you can see all freq. \oplus nt but not when i.e. the time info.

Page No.:	1
Date:	Shree

So FR ↑
 TR ↓



Aliasing

when aliasing occurs - i.e. when there is an overlap in signal low frequencies will get affected which are the voiced regions.

② In discrete domain

$$\text{output} = \text{convulat}^n(\text{input}, \text{system response}) \\ = x[n] * h[n]$$

Lecture 8 (18 Aug) 4th Sept) - STFT for ~~short~~ time scale modification

- if we resample -

- (a) - discard 1/2 sample every 5 samples
 - pitch cycle decreases
 - pitch rises

- (b) undersampling - ~~increases~~ add zeros in b/w
 - pitch cycle increases
 - hence pitch rises

Q what should be window size?

if we do resampling in frequency domain - why pitch stays same

Resampling in time - we take

domain

sample acc. to time

However,

in spectrogram — resampling happens
only in 20 ms

so we have \times 20ms
frames

so we discard ~~or~~ / add
frames

In 20 ms — 2-3 pitch cycles so
the pitch is preserved
and speaker charach. also.

Q why is it cl. Phase Vocoder

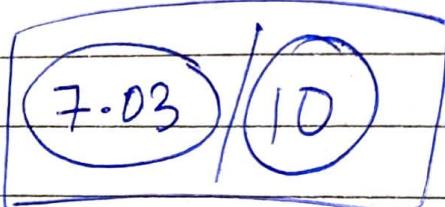
- Phase Vocoder is used for time scale & pitch modification
- increasing/sampling rate — will make the audio fast / slow
- Talking time — speeding up (sample it in time domain)
 - ↓
 - discard some samples
 - slow down — undersampling

Quiz

- Q1. Wider window TR↑ - False (+1)
- Q2. Speech Coding for cellular - True (+1)
- Q3. Female pitch and bird - False (0) (+1)
- Q4. Synthesiser - ① — (+1)
- Q5. Sampling Theorem → $\frac{-1}{3}$
L limits the min. freq requirement.
- Q6. 270 kbps $\frac{-1}{3}$

$$8 - \frac{2}{3}$$

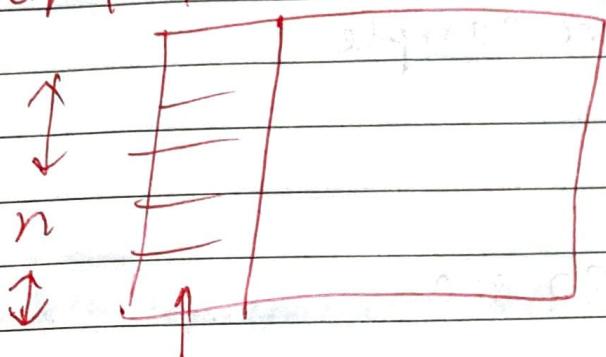
8-



SSP A2

buffer — x into frames of size n
 $\hookrightarrow x.size()$ — total #samples
 $n \rightarrow$ # samples in a frame

output $\lceil \frac{L}{n} \rceil = \text{no. of frames}$



one frame

94274 — 160 window size

1 0 0 5 0 0 0 3 0 3 0 5 0 0

(10-2+1)

$$\begin{aligned} & 94274 - 160 + 1 \\ & = \underline{\underline{94115}} \end{aligned}$$

94274 —

for $i=0$ to $i < 94274$; $i+=2$

1 sec 8000 samples
 $\frac{94274}{8000 \text{ sec}}$? 94274 sam

$$8 \times 0.000125 \times x = 5.30$$

~~8x~~

2 msec in one frame

$$\frac{1}{8000} \text{ time per sample}$$

$$\frac{20 \times 500 \times x}{1000} = 5.30 \text{ sec}$$

$$= \frac{530}{2}$$

11-8

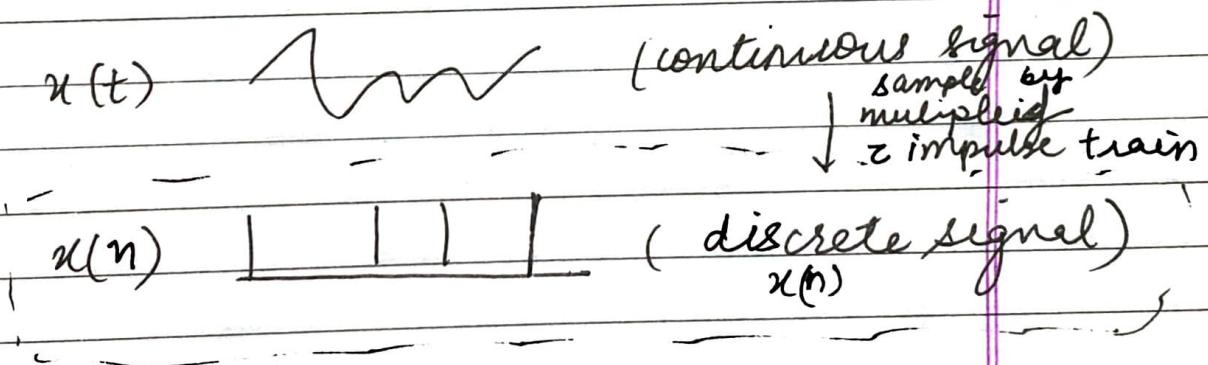
Q Circular convolution for DFT. Why?

$$\begin{array}{c}
 160 \\
 20 \\
 20 \\
 80 \\
 \hline
 3x1 \\
 100
 \end{array}$$

Class : 8th September 2020

Linear Time Invariant System

Impulse Train



Linearity condit'

$$a x_1(t) + b(x_2(t)) \rightarrow a y_1(t) + b(y_2(t))$$

$$y_1 x_1 \rightarrow y_1$$

Three	Four	Date:	Page:
Three	Four		

$$y_2 x_2 \rightarrow y_2$$

Linear

$y = mx \quad \checkmark \text{ linear}$

$y = mx^2 \quad \times$

?

 $y = mx + c \quad \times$

Time Invariance

$x(t) \rightarrow y(t)$

$x(t-1) \rightarrow y(t-1)$

if input

Class 11th September

→ Z-transform

→ Impulse response &

Discrete signals are represented in the form
of z-transform.

→ Motivation: very easy for discrete system analysis and also for some of the signals FT may not exist but z-transform exists. Why?

We can represent 'Z' as $A e^{j\phi}$

if $A = 1$, then it is ~~DFT~~ DTFT only.

because even $x[n]$ not stable

by multiplying it into $\cancel{A} A^{-n}$.

will become

A^{-jn} will become.

This makes it as a stable signal

and then we can get $x[z]$

Q Given a signal, how to get z-transform

Q Circular convolution for DFT. why?

- Signal is periodic
- We discretize the input signal in time domain - because of convolution property
HOW?
- Discretizing by taking impulse train and multiplying it signal.
- + because of this multiplication - the spectrum is the convolution of the impulse response & signal response.

because of which the output becomes periodic and to avoid aliasing we want

$$f_s \geq 2 f_m$$

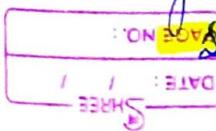
But,

Even after sampling in time domain, we cannot compute DFFT or compute (because for that we need to discretize in frequency domain).

Now since the signal is periodic we cannot apply linear convolution so we apply circular convolution.

At the same time, we have to take 2 signals (constraint in circular conv.)

should be of same length) → & apply convolution



→ In n -point DFT

$\text{if } n >$ also - we use 0-padding
 it is ✓ technique (adding zeros to make same length/ n points)
 this increases frequency resolution

Mel-Frequency / Cepstral Analysis

Spectral Features

- Spectral features extracted using segmental analysis
- what do these features represent?
- modifying spectral features → cepstral features

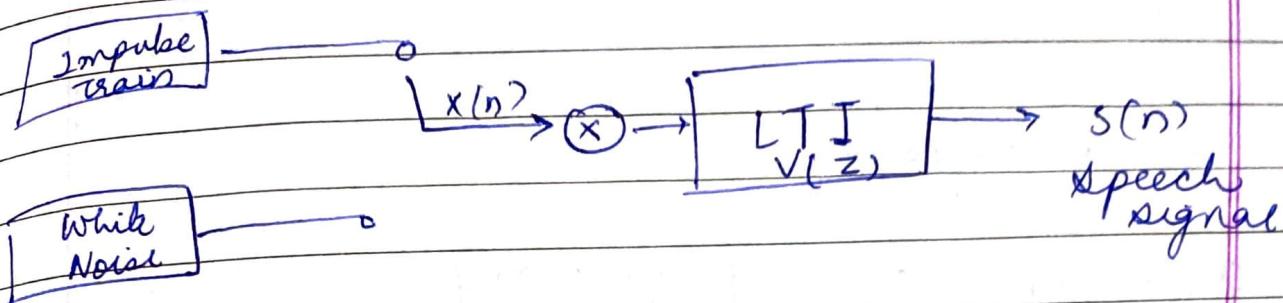
Spectro

SPECTRUM → CEPS TRUM

Cepstrum is inverse of Spectrum

Spectrogram

LTI



$$s(n) = x(n) \text{ convolution } v(n)$$

So, if we observe, we want information related to the vocal tract because it is responsible for changing sounds & it is unique for every human & it is that acts as resonator — hence it is required / more useful.

Unfortunately, the signal we receive is $s(n)$ which is convolution of $x(n)$ & $v(n)$

So we need to separate the two — using LP analysis

Cepstrum = Inverse FFT (log(spectrum))

$$\text{FFT}(S(k)) = X(k) * V(k)$$

$$\log(\text{FFT}(S(k))) = \log(X(k), V(k))$$

Cepstral Analysis

after FFT — frequency

~~peaks~~ p
 spectral env. peaks - formants
 spectral env - spectrum of CP quotients
 smooth spectrum of FFT.

How to achieve separation?

of $X(n) * V(n)$

i.e. given

$\log X[k]$ get $\log H[k] \oplus$, $\log E[k]$

$$\underbrace{\log X[k]}_{\text{spectrum}} = \underbrace{\log H[k]}_{\text{env. slope}} + \underbrace{\log E[k]}_{\text{details}}$$

80 IFFT ($\log X[k]$)

will separate the 2 into
 high freq & low freq
 details (excitation)
 $\log H[k]$ envelope (VT feat)
 $\log E[k]$

$\&$ First 13/15 values of low frequency ($\log E[k]$)
 represent vocal tract

- diff for each frame in these 13 values
 show how vocal tract is changing

Mel-frequency Analysis

mel-scale \longleftrightarrow linear scale

$$m = 2595 \log_{10} \left(1 + \frac{f}{700} \right)$$

Spectrum \rightarrow mel-filters \rightarrow Cepstral analysis
 $\frac{1}{MFCC}$

LP Analysis (15-09)

exploits redundancy
redundancy of speech signal

$$\sum_{i=0}^p a(i) y(n-i) = \sum_{j=0}^q b(j) x(n-j)$$

Apply z-transform,

$$\frac{N(z)}{D(z)} = H(z) = \frac{\sum_{j=0}^q b(j)z^{-j}}{\sum_{i=0}^p a(i)z^{-i}}$$

↑ roots - zeros
↑ roots - poles
↑ or denominator

When only numerator — all-zero system
 (— $D(z)$ is unity)
 or Finite Impulse Response system
 or moving average (MA)

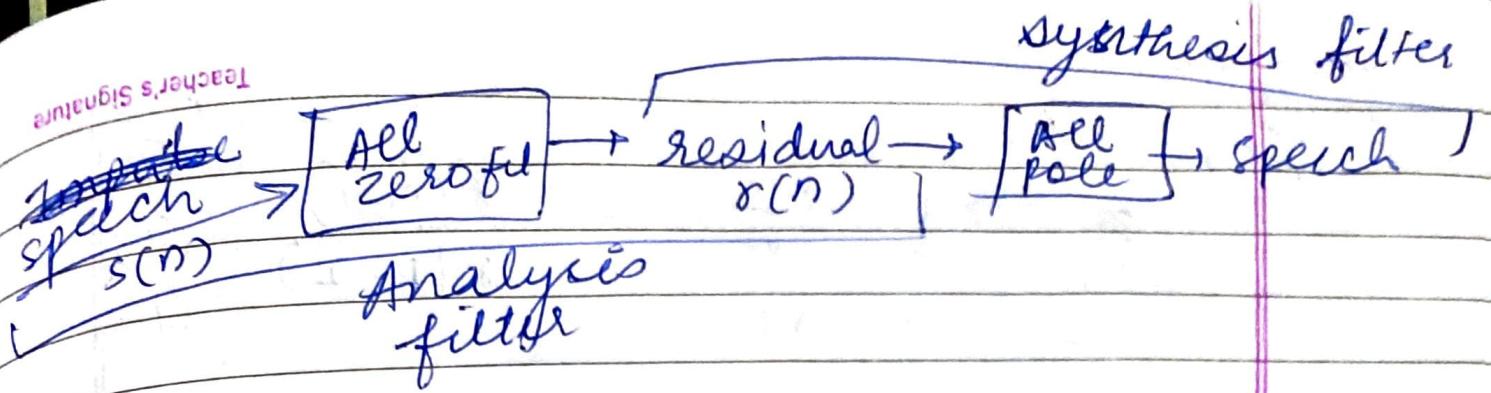
When $\frac{1}{\text{Denom}} = \text{all-pole system}$
 N(z) — constant

Infinite Impulse Response

(Antagonistic)
AR

vocal tract

When both — ARMA



$$\text{All zero filter} = \frac{1}{H(z)} = 1 + \sum_{k=1}^P a_k z^{-k}$$

In AR Model,

$$\hat{y}(n) = -\sum_{i=1}^P a(i) y(n-i)$$

represent
excitation

$$e(n) = y(n) - \hat{y}(n)$$

$$e(n) = y(n) - \sum_{i=1}^P a(i) y(n-i) \quad \text{--- (2)}$$

Solve (2) by auto correlation approach

Total error =

$$E = \sum_{n=-\infty}^{\infty} e_n^2$$

$$\frac{\partial E}{\partial a(i)} = 0; \quad 1 \leq i \leq P$$

$$\Rightarrow \sum_{i=1}^P a_i r(k-i) = r(k) \quad i \leq k \leq P$$

where $r(k) = \sum_{n=k}^{N-1} s_n s_{n-k}$

So,

$$\begin{matrix} \gamma_0 & \gamma_1 & \gamma_2 & - & \gamma(p-1) \\ \gamma_1 & \gamma_0 & \gamma_1 & - & \gamma(p-2) \\ \gamma_2 & \gamma_1 & \gamma_0 & - & \vdots \\ \vdots & & & & \gamma(0) \\ \gamma(p-1) & & & & a_p \end{matrix}$$

R
 L auto correlat' matrix
 A

Solve

$$\begin{aligned} RA &= R' \\ A &= R^{-1} R' \\ \text{Lp coeff} &\downarrow \end{aligned}$$

Prediction Order

no. of prev values you should take depend on analysis bandwidth

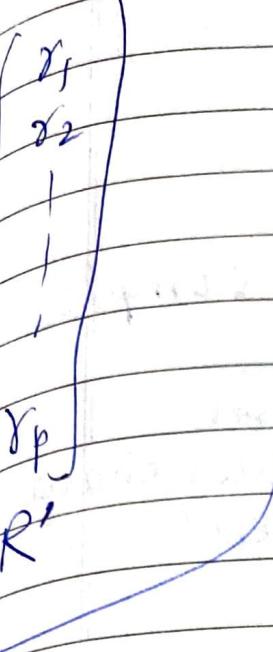
$$P = \frac{2(BW)}{1000} + C$$

→ 1 foeman per kHz of BW

and

1 foeman requires 2 complex conjugate poles

so, for every foeman 2 predictor coeff. → 2 ~~foeman~~ coeff per kHz of BW



greater P \rightarrow
more formants

~~VUV~~

- Error more where
more excitation

- errors - more for unvoiced

formant extraction

\rightarrow LP ~~residual~~ spectrum \rightarrow formants
LP ~~residual~~ \rightarrow VT features (ak)
error gives \rightarrow excitation
features

AK-LPCC - ~~graph~~

~~Dpitch~~

Auto correlate of LP residual

diff b/w 2 peaks

Q To get 5 formant for LP spectrum
\$ what should be LP order?

1 formant - 2 complex conjugates
 \Rightarrow 2 order

$$\underline{5 \times 2 = 10}$$

Q why cepstrum?

L to separate VT from excitation

Q Why LP analysis spectrum is smooth?

VT as all pole resonator
 assuming are resonator has only 5 responses ($P/2$) - so accordingly only 5 frequency are highlighted.

We get spectrum from LP coeff
 not directly from signal

We removed excitation
 so no zig-zag only vocal tract is highlighted

Is MFCC enough for ML model?

Yes.

In which scenario MFCC may not be enough.

e.g. emotional recognition because it requires excitatory features like strength of excitatⁿ

If training & testing condition is not matched

then LPCC is better
or MFCC - ans = MFCC

LP Coeff

using autocorrelatⁿ

minimizing diff of seqⁿ of original & predicted signal

Inverse filter is all pole filter.



No false

Residual → speech } all pole filter.



25th September - Epoch Extract

Method

① HE (LP Residual) - Peaks

② ZFF - directly from speech signal

} ① Instantaneous frequency
which comes from
analytical signal

} ② Property of Impulse (Epoch Local)

doesn't use these concepts

Q what is FFT (Impulse)

L constant (DC signal)

Impulse will have influence
on all frequencies

→ Instantaneous frequency represent
a change wherever there is
an impulse

→ filter signal at diff frequency to
get plot.

(depends on VT influence on
diff frequencies)

so if VT influence less

on 580Hz

Plot will show Epoch Local better

25, 29, 20

Teacher's Signature

Because impulse will have influence on all feed but VT can have on some.

Since VT doesn't influence low frequency (it starts at 100 Hz etc) so choose / filter signal at low freq.

ZPF

→ ZPF is a low pass filter (at 100 Hz)

→ Ideal ZPF is incubator
i.e. add the samples

(①) Difference i.e. integrating
signal to remove time varying low frequencies
① → pre-emphasis

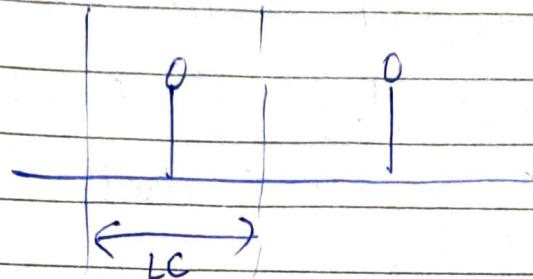
→ 12 times ZPF = 4 times integrated

② Trend Removal
remove mean from each value
around M values.

③ negative to positive zero-crossings
gives the epoch location
in mean - sub signal (mss)

Evaluation approach

Larynx Cycle - $\frac{1}{2}$ pitch period on both sides of epoch
(LC)



so within Larynx cycle is nothing is

⊕ mt \Rightarrow epoch "deat" missed

or is more than one

\Rightarrow false alarm.

when only one in LC - ✓

→ Can be used for V and VV detectn.

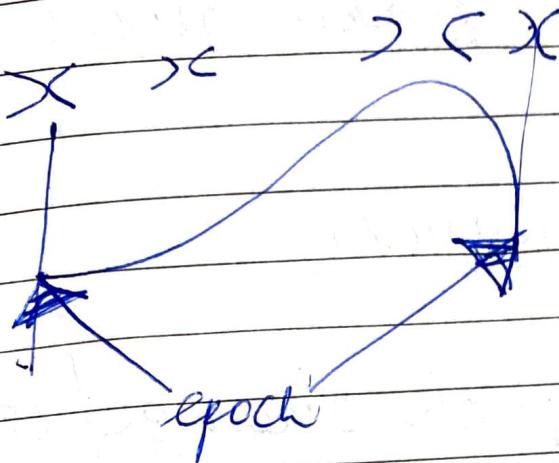
all peaks of impulse not same as
strength of excitatⁿ is diff.

is given by slope of MSS

sharply - strength ↑

slow slope - low strength

GVV - Glottal Volume Velocity
excitation feature



How to get GVV

LP residual \rightarrow [Integ ratio] \rightarrow GVV

\rightarrow how fast opening/closing
vocal folds

\rightarrow useful for emotion, singing speech

\rightarrow works well even if noise

\rightarrow CMU Arctic Database

Noise

↳ reverberant

↳ bg noise

↳ babble / multi-speaker noise

(29m-09) Prosody features

↳ suprasegmental

- ① Intonation pattern → pitch change / contour
- ② Duration → length of syllables - rate of speaking
- ③ Intensity → energy contour

② Duration is difficult to extract

Issues in prosody extract?

- ① what is high for someone may be normal for other
- ② for each ind. we need to know normal voice first to understand if there is a change
- ③ Duration is difficult



Applications

- ① TTS
- ② Lang. Ident.
- ③ Error corr.
- ④ Speaker recog
- ⑤ Speech recog

→ extract features for whole regions
and use them as anchor points

Teacher's Signature

6-9 Oct missing

13th OctoberVOP - vowel onset point

↳ vowel regions — ROBUST
 feature extraction

Indian L — syllable

↳ has vowel nuclei

helps in

↳ speech recognition

↳ prosody modifcatⁿ

↳ changing duratⁿ of speech

↳ sketch / squeeze only
 vowels

→ vowels act as anchor points

→ vowel onset ← vowel offset
 ↓ vowel region

→ 3 approaches =

↳ ~~VOP~~ source = excitation source - [residual
 complementary] - spectral peaks (10 peak values) are summed
 to each other - time informatⁿ (modulatⁿ spectrum)

all three

show significant
 change
 at VOP

time domain
 envelop of
 signal

- ① Signal
- ② LP Residual
- ③ Hilbert Envelope of LP residual

assumpt'n - significant change in envelop — VOP — peaks

- ④ but spurious peaks
- eliminate — using - first order diff
& keep a threshold

peak & valley values are
normalised

- ⑤ convolut'n \cong first order Gaussian differentiator

$L = 500$ (length)

$\sigma = 200$ (sigma)

- ⑥ Peaks in the final signal are VOP

-ve peak — V offset point

16th October (VOP)

- Spectral info is better for coded speech
- Changes =
 - instead of n peak values find epoch locat's using ZFF
 - Spectral energy around ω_L is high
 - within 2 epochs, 30% of cycle is estimated spectrum — vowel/spectral energy will be high

AI & ML Trends

- speaker, languag, emotion — simple
- ASR, TTS — advanced
- data ↗ diff accent
diff gender, age group

GMM (Gaussian mixture Modelling)

lot of data

classes - 23 (languages)

① Assume some gaussians and fit the features into these gaussians. More gaussians is good because they capture more information provided we have enough data.

② → GMM is universal background model

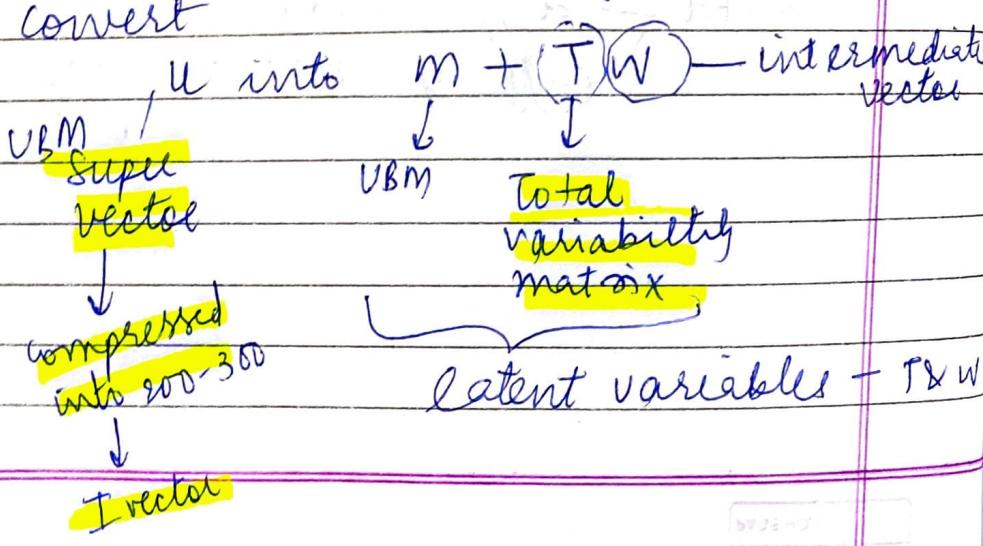
where:

VBM will be built & huge data & adopt that model to the target data (1024 mixtures)

We are only interested in speaker dependent component of the VBM vector

③ So use Joint factor analysis

convert



④ I-vectors are used as a model or also as a feature vector.
Then use SVM & DNN for classification

Now, I-vector has been replaced -
DNN & attention

| discriminative models at large
weights of
connection are optimised
that is the model

| GMM is generative - for each class a diff GMM

but here in DNN - only one network that classifies all classes.

| fail when training & testing conditions are diff

| work well with large data

- Architecture
- Layers
 - Neurons
 - learning Rate etc

} Parameters

DNN & attendⁿ

- frame level decision
- Utterance level decision

e.g. Lang. Ident.

- 10 lgs

attention mechanism

L take decision at utterance level & not at frame level
 context vector

→ frames are given to attendⁿ model
 & a weighted average of the hidden layers give context vector

→ So attendⁿ weights are low in silence region

→ Multi-headed attendⁿ

LID → explicit - STT - then find L
 implicit - feature - classif

ASR (Automatic Speech Recognition)

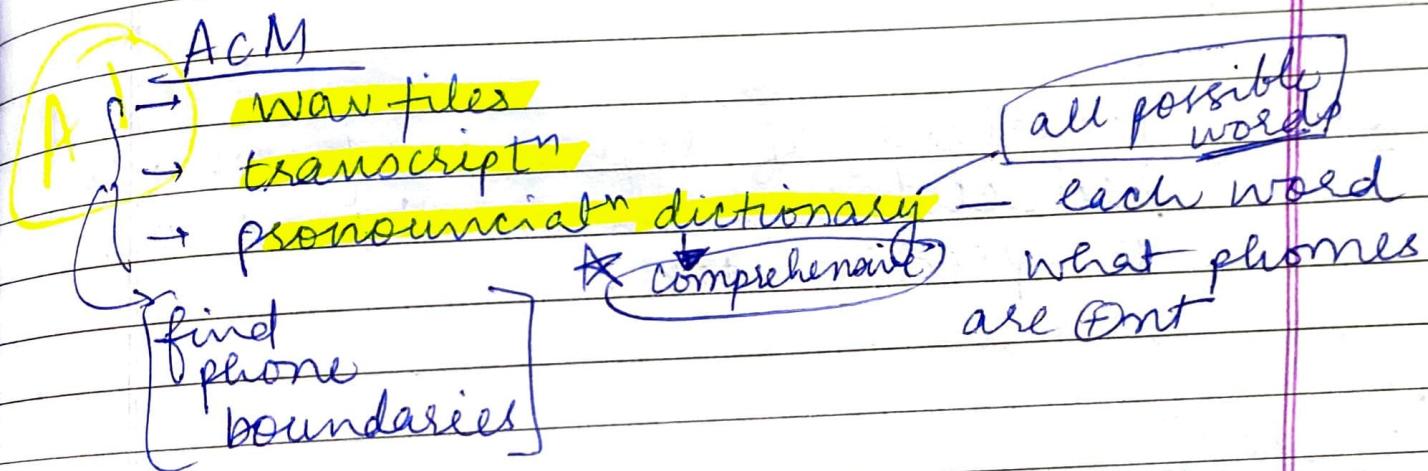
- Speech to Text
- Language Model → A conistic Model

Y - MFCC

W - sequence of words

$$\text{STT (ASR)} = P(W|Y)$$

$$P(W|Y) = \frac{P(Y|W)P(W)}{P(Y)} \xrightarrow{\text{lang M}} \xrightarrow{\text{constant}}$$



L multiple pronunc — mentioned
(depending on context)

20th October

Problems =

① Machine phoneme - word - sentence
listens

so errors may creep in

② Data - should nullify the accent / age group
gender / speaking style
variation

③ Fragmentation - chop wav files into
fragments
& manually transcri

LM - probability of sequence of phones
reduces / cuts down the errors
to predict text.

Evaluation

Word error Rate

① Substitution

② Insertion - extra word

③ Deletion

$$\frac{S + I + D}{N} = \frac{S + D + I}{S + D + C}$$

Dynamic Time Working

Better - HMM (Hidden Markov Model)

A6 } data < 1000 hours

else HMM gets saturated

Acoustic Models

- A3 ① GMM- HMM (most popular) captures sequential info
- ② DNN- HMM

→ ASR are speaker & domain specific

→ 3 diff ASR ~~can't~~ be modeled:

- A7 - Read speech - reading
- Extempore - lecture/talks
- Conversation - phone / meeting

So one ASR may not work for all

HMM - model the temporal variability

A8 → each phoneme will have a HMM

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

→ three states - beg
middle
end

→ depending on context, position ~~&~~ duration
pronunciation of a phoneme changes

Hidden states:

- [we don't know where the phoneme is starting / ending]
- [assume 3 hidden states in the productⁿ of a phoneme]

H Markov process

- [what is the prob. that I am in a particular state for time t=0]
 - [what is the prob. that I will remain in the same state or go to some other state.]
- (state probabilities)

in speech HMM

- [we cannot go back]
- [only left to right topology]

A = state transition prob. ($a_{ij} = i \rightarrow j$)

B = symbol probability ($b_{ik} = p(i^{th} \text{ symbol is state})$)

↳ GMM

π = initial probability

Three Basic Problems

① Evaluation Problem

$$P(O|\lambda)$$

given model & output

$P(\text{model output})$

forward & backward algo

② Decoding Problem

give $m \& o$
most likely
state sequence

vitte-Bi

③ Learning Problem

parameters
to get high
prob of generating
sequences &
model

limitation:

- limited parameters to optimise
- (19) - so will not change & increasing data — can't optimise much
- on the other hand — NN — have a lot of parameters

Q Sequential info is used in which applications :

- a. ~~LID~~ LID
- b. speech recog
- c. ~~Speaker Independent sys~~ Text

Q What is A in HMM

A - matrix

what are elements of a_{ij}

e.g. Transition from i to j
prob \rightarrow

Q What is symbol prob.

$b_{ijk} \rightarrow$ prob of symbol k in state j

N
no. of states

M
no. of symbols

optimised
based on experiments

Q What is input to AM
L MFCC

Q What is input to LM
L transcripts

Q Transcripts dict
L phonetic descriptⁿ of each word

Q HMM vs DNN
↓

one
HMM
for
one
phoneme
(generative)

A 10

one model
for entire
classification

(discriminative)

One model classifies
all phonemes

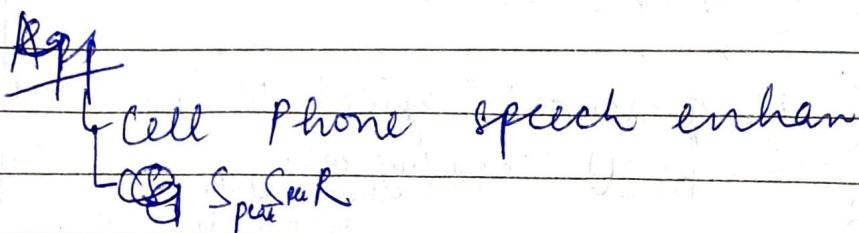
October 2017

Note 3 (Speech Enhancement)

↳ Once speech is degraded with noise - how to enhance speech & improve the quality so that intelligibility res.

Speaker Recognition

↳ Naturalness & intelligibility



3 kinds of degradation

① Bg noise - AC noise $s(n) + d(n)$

② Reverberation - Reverberant/reflected speech $s(n) * h(n)$

③ Competing speaker - $s_1(n) + s_2(n)$

↳ most difficult since both req. interfering speech

* not req. have the same characteristics

PROS = effective

Approach

problem - need for explicit modeling of degraded component

Major ① Spectral Processing - estimate the noise spectrum & remove it from degraded speech spectrum

② Temporal processing - in time domain
Enhances the high SNR region

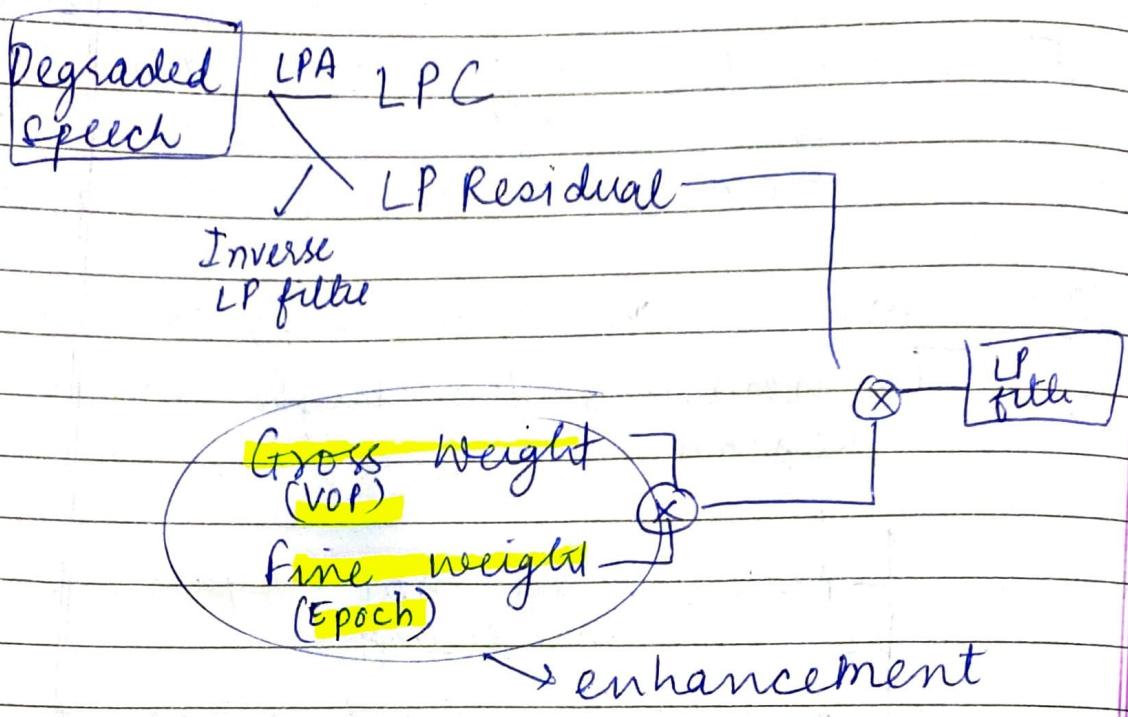
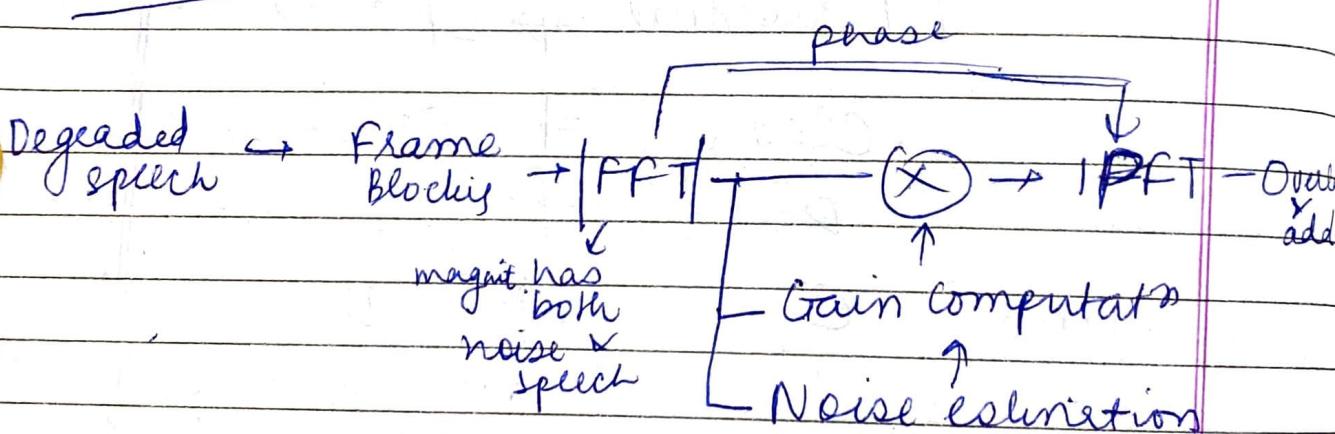
- └ PROS = may not require explicit modeling
- └ Demerit: Ineffectiveness

LP residual \Rightarrow LPC \rightarrow speech

$$\begin{array}{ccc} \text{Clean} & + & \text{Clean} \rightarrow \text{Clean} \\ \text{D} & + & \text{D} \rightarrow \text{D} \\ \text{C} & \neq & \text{D} \end{array}$$

so, the enhanced (LP residual)
+ degraded LPC

enhance

TemporalSpectral

approach: Spectral subtraction

→ Diff noise have diff spectrum

Music Noise

↳ speech may be \oplus in noise so when we do spectral subtraction - it creates some zeros in the spectrum

Temporal (non-speech regions become zero)

Gross Level (0-1 functⁿ)

- ↳ get window functⁿ
- ↳ plot 10 peak values (as in VOP detect)

So 3 evidences -

- residual sum of peaks in DFT spectrum
- Hilbert envelop
- modulation spectrum

are combined to enhance

→ Non-linear Mapping

functⁿ on the ~~ratio~~ ~~normed~~ sum

- ↳ making things zero in the low SNR (low energy regions)

UNIQUE IS TAACHERS

Finite Weight Final

- Find epoch using zero crossing
- then take 10 values around epoch
- * create a finite weight func

Final

- Combine Gross level fn & finite weight fn
- Multiply LP residual & the two

Method: combination of the two

- use enhance LP residual
- ~~Temporally Processed speech~~ is used to estimate the noise spectrum.
- Noise Estimation is done using (in spec proc)
the noisy speech only
and that will be applied on temporally enhanced speech

NOV 6

TTS

- ↳ natural & intelligible
- ↳ need only one speaker data

Front End —— Back End

① Text Analysis

- ↳ Text Normalisation
- ↳ Linguistic Analysis

Phonetic Analysis

Grapheme to Phoneme

- Prosodic Analysis

- ↳ Pitch & Duration (depending on context)

Back-End

↳ Speech synthesis

↳ approach:

↳ concatenated synthesis
(unit selection synthesis)

↳ phoneme level

or

↳ syllable level

↓
Select right

unit * concat-
enate

Optimal Text Selectⁿ

- ↳ Hand collect valid sentences
- ↳ select sentences that cover all sounds of a lang.
- in no. of combinations & no. of times (in beg, middle, end) — all variations of sound

Problem — emotionless

- ↳ always need a backend database
- time consuming

Adv

- ↳ more natural & intelligible

Parametric Synthesis

Using HTML

- ↳ Problem — naturalness (less)

- ↳ save parameters instead of the database

Pros → reduce size

- ↳ lower development cost

- ↳ flexibility — modify parameters to generate ~~content~~

↳ attractive — diff voices because

↳ e.g. storytelling

Cons

- ↳ less intelligibility
- ↳ robotic - less natural

↳ Speaker is fixed - same speaker
so we get bored and less attractive

concatenate % con

WaveNet

Q How to create emotional speech

TTS is evaluated - Mean Opinion Score

NOV 10Course

- Speech Basics
 - Productⁿ
 - LTI model of speech
 - Modelling
 - V/U/V detectⁿ, pitch detectⁿ
- Features Ex
 - Spectral - seg
 - Excitation - subseg
 - Prosody - supra
- Speech Application

STS (Bahubhashak) - reduce lg barrier

NPTEL

directly - ~~ASR~~

ASR + TTS

↓
 loss - emotion / characteristics / prosody

NPTEL is naturally recorded
 Eg. - filling words, correctn,
 pit
 repetition

Problems in Indian L

- ① low resource
- ② Morphologically rich
- ③ ~~See~~ free-order language

NN may help when structure is same from source L to target L

So we follow ASR + TTS.

① ASR (Word Error Rate)

② MT (fragmentation effects MT) so post processing is req.

③ TTS

↳ naturalness

* intelligibility

Problems in ASR

↳ time alignment

↳ lip syncing

↳ use same filling words?

↳ use incorrect text also

↳ one L - less words }
other L - more words }