

EC864

SPEECH AND AUDIO PROCESSING

MANVITH PRABHU
211EC228

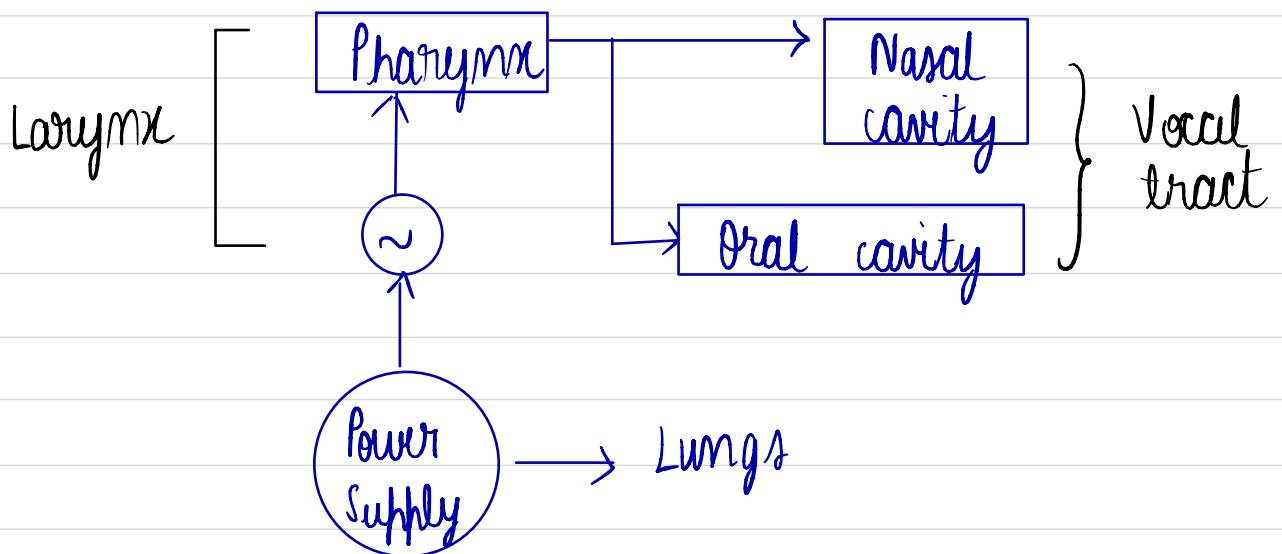
SPEECH

Speech generation in Humans:

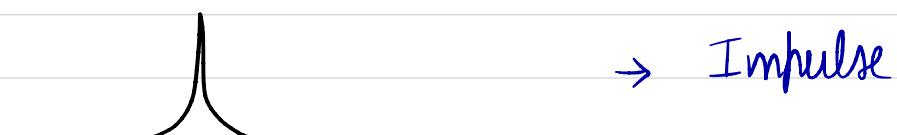
Brain → Lungs → Vocal folds
↓

Pressure signal ← Vocal tract

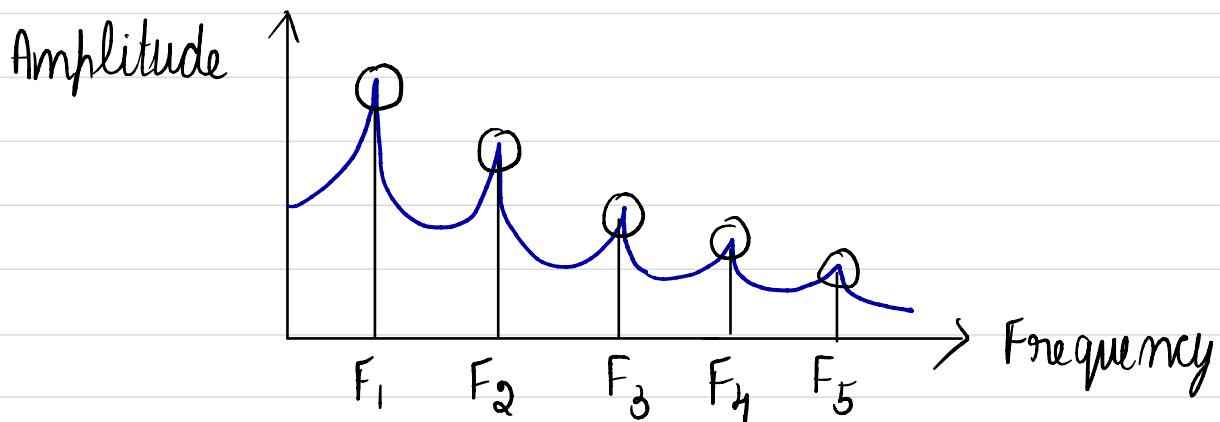
Block diagram



Types of sounds



NOTE: ECG is used to measure resistance between vocal folds \rightarrow gives a waveform.



The peaks in frequency determine the vowel

- \rightarrow On an average the frequency of voice of men is of range $60 - 100 \text{ Hz}$
- \rightarrow For women $\rightarrow 90 - 220 \text{ Hz}$
- \rightarrow For children $\rightarrow 250 - 500 \text{ Hz}$
(Fundamental frequency)

Fundamental frequency

Pitch / Fundamental Period : The time duration of one glottal cycle as the pitch period.

$$\text{Fundamental frequency} = \frac{1}{\text{Pitch period}}$$

NOTE: Pitch is often used to describe the subjectively perceived "height" of a complex musical sound even when no single fundamental frequency exists.

Glottal source band characteristics

Jitter: The variation in the pitch period across different glottal cycles

Shimmer: The amplitude variation in airflow velocity across glottal cycles.

Vocal Tract (VT)

→ Comprises oral cavity (from larynx to lips) and nasal passage

→ Average length of oral tract (adult male):
~17 cm.

→ For adult females, it is comparatively shorter

→ The pressure signal at the output of vocal folds (during voicing): time-varying buzz like sound.

→ VT: spectrally "colour" the source which is important for making perceptually distinct speech sounds.

Spectral shaping (formant frequency)

- Under certain conditions, the relation between glottal airflow velocity and vocal tract velocity output can be approximated by a linear filter with resonances.
- In speech science, these resonances known as formant frequencies / formants
- Formants change with different vocal tract configurations.

Categorisation of sound of source

- 1) Vocal sounds: Speech sounds generated with a quasi-periodic glottal source.
- 2) Fricatives / unvoiced sounds: These sounds are generated by forming a constriction at some point in the vocal tract, (usually towards the mouth end) and forcing air through the constriction at a high enough velocity to produce turbulence.
 - This creates a broad-spectrum noise source to excite the vocal tract.
e.g: "sh" in should -

3> Plosive sounds: Generated from complete closure (usually VT), building up pressure behind the closure.

Eg: "ch" in chase

Types of Phonemes

→ Every language has a set of distinctive sounds or phonemes.

Eg: American English has around 42 phonemes.

- 1> Vowels
- 2> Diphthongs
- 3> Semivowels
- 4> Consonants.

• Each of these classes can be further subdivided

1> Vowels: These are produced by exciting a fixed VT with quasi-periodic pulses of air cause by vibration of vocal chords

2> Diphthongs: It is a gliding monosyllabic speech item that starts at or near the articulatory position of one vowel and moves towards the position for another.

Eg: /eɪ/ (as in bay),
/aɪ/ (as in boy)

/ɔɪ/ (as in boy)
/ju/ (as in you)

/əʊ/ (as in boat)
/aʊ/ (as in how)

3) Semivowels:

- The group of sounds consisting of /w/, /l/, /r/, and /y/ are called semivowels because of their vowel like nature.
- Generally characterized by a gliding transition in vocal tract area function between adjacent phonemes
- Acoustic characteristics of these sounds are strongly influenced by the context in which they occur.
- These sounds can be described as transitional, vowel-like sounds, and hence are similar in nature to the vowels and diphthongs.
Eg: "w" in we

4) Nasals:

- Nasal consonants /m/, /n/ and /ŋ/ are produced with glottal excitation and vocal tract totally constricted at some point along oral passage.
- The velum is lowered so that air flows through the nasal tract, with sound being radiated at the nostrils.

Categorisation of speech sounds

- Phonemes arise from a combination of vocal fold & vocal tract articulatory features
- Articulatory features include:
 - 1> Vocal folds are vibrating or open
 - 2> The tongue position, height & it's constriction
 - 3> Vlum state (nasal sound or not)

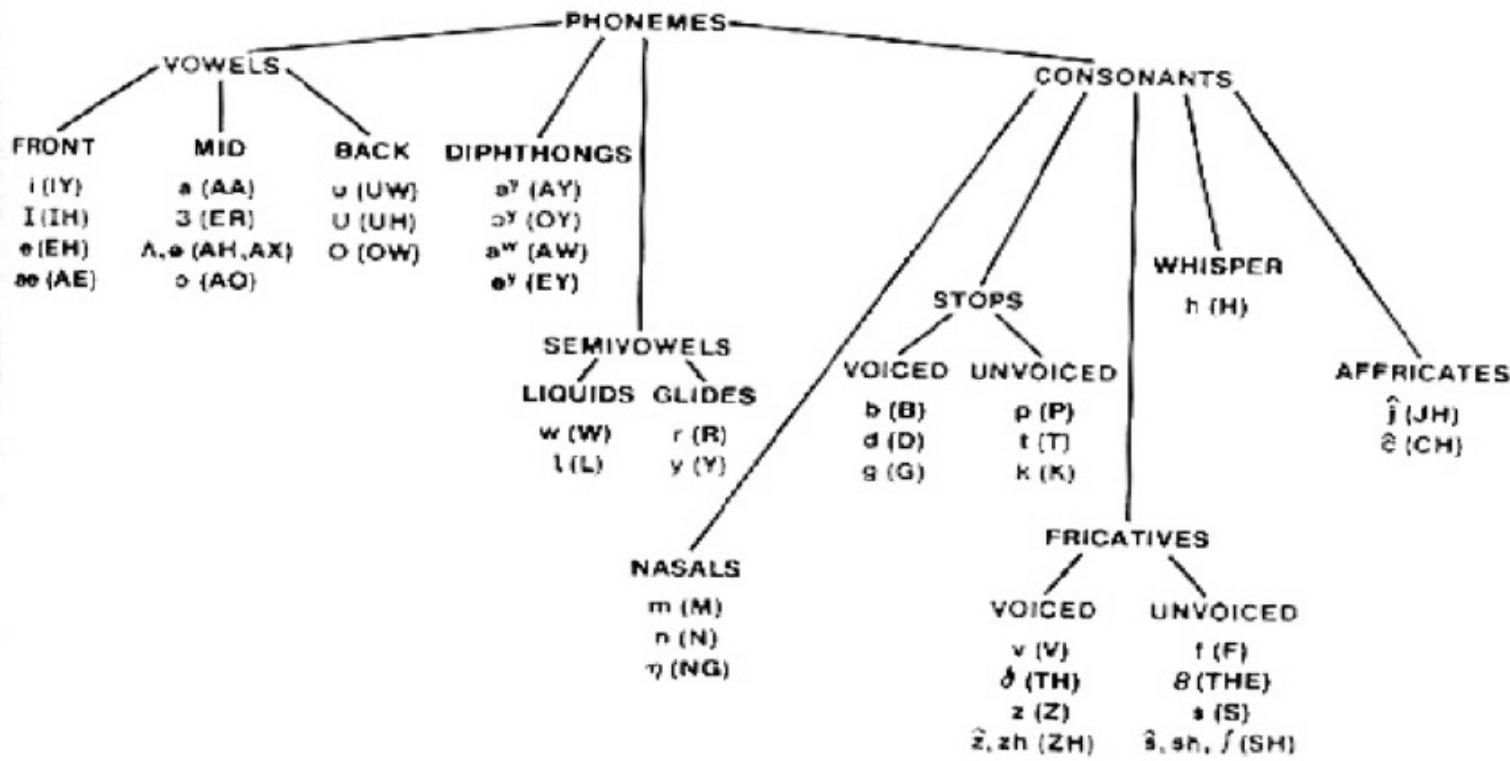
Vowels:

VT shape = f (tongue, jaws, lips & vlum)

For non-nasalized vowels, the nasal passage is not occupied with oral cavity

Tongue: Primary determinant of VT shape

Place of articulation: front, center & back of oral cavity



Consonants:

1> Nasals: Quasi-periodic airflow puff from the vibrating vocal folds.

System: The velum is lowered, oral tract being constricted totally, and the air flows mainly through nasal cavity.

Nasal sounds radiates at nostrils

Example: /m/- oral tract constricts at lips

Nasals have lower frequency compared to vowels

2) Fricatives :

Types : voiced & unvoiced

1) Source in unvoiced : Vocal folds are relaxed & not vibrating. Noise is generated by turbulent flow of air at some point of constriction (narrower than with vowels) along the oral tract.

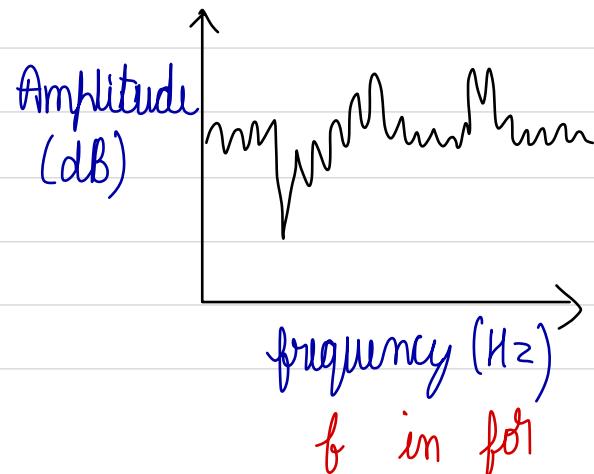
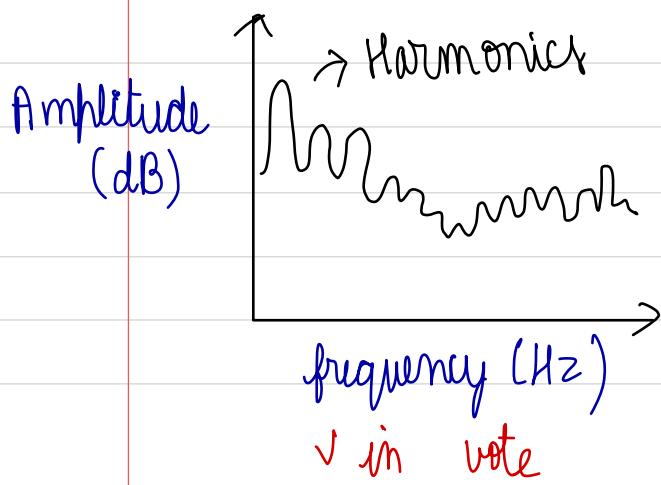
System : Location of the constriction by the tongue at the back, centre / front of the oral tract as well as teeth / lips.

Eg : f - (for) , h (he) , s (see)

2) voiced : Source has a similar noise source and system characteristics to unvoiced fricatives with simultaneous

Eg : v in vote

Wideband spectrum : unvoiced : noisy spectrum
voiced : both noise and harmonics



3> Plosive sounds

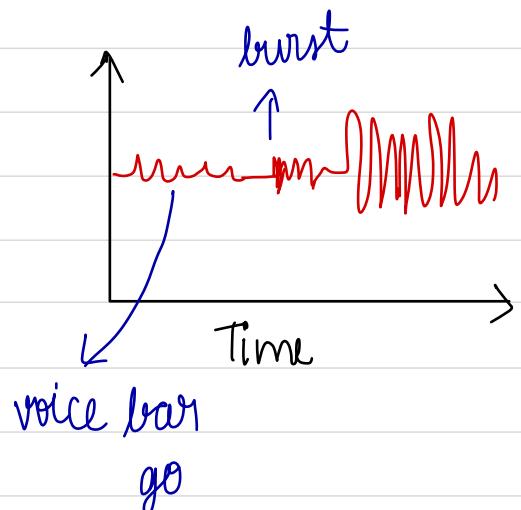
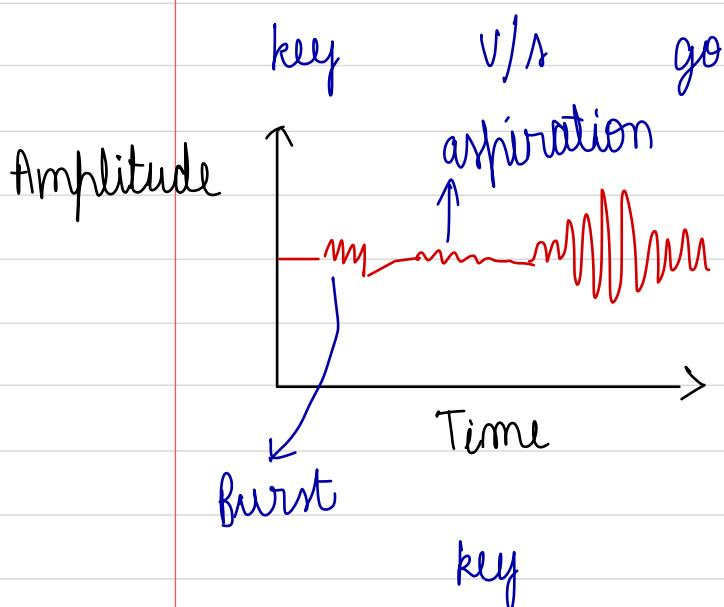
Unvoiced Plosive: • A burst is generated at the release of the buildup of pressure behind a total constriction in the oral tract.

- This burst is idealized as an impulse source
 - No vibration of vocal folds.
- Eg: p (pay) , t (to) , k (key)

Voiced Plosives: Pressure is built behind an oral tract constriction, but vocal folds can also vibrate.

Presence of voice bar

Eg: b (be) , d (day) , g (go)



Diphthong: They have 2 formants.

continuant sounds: eg: Vowels, Fricatives, Nasals
Produced by a fixed VT configuration
excited by the appropriate source

Non-continuant source: e.g.: Diphthongs, Semivowels, stops, affricates.

Sound v/s vibration

POA	UVVA	UVA	VUA	VA
MOA	χ, ψ (χ)	ψ, ϗ [χ ^h]	χ, ϗ [χ]	ψ, ϗ [χ ^h]
Palatal	ɛ, ē [t̪̚]	ɔ, ɔ̚ [t̪̚ ^h]	ɛ, ɔ̚ [d̪̚]	ɔ̚, ɔ̚ [d̪̚ ^h]
Alveolar	ʒ, ʃ̚ [t̪̚]	ɸ, ɸ̚ [t̪̚ ^h]	ɸ, ɸ̚ [d̪̚]	ɸ̚, ɸ̚ [d̪̚ ^h]
Dental	ʒ, ʃ̚ [t̪̚]	ɸ̚, ɸ̚ [t̪̚ ^h]	ɸ̚, ɸ̚ [d̪̚]	ɸ̚, ɸ̚ [d̪̚ ^h]
Bilabial	پ, پ̚ [p̪̚]	پ̚, پ̚ [p̪̚ ^h]	پ̚, پ̚ [b̪̚]	پ̚, پ̚ [b̪̚ ^h]

The figure shows a UV-vis spectrum with Absorbance on the y-axis and Wavelength (nm) on the x-axis. A sharp peak is visible at approximately 350 nm, and a broad absorption band is centered around 550 nm.

VUA: 

VIA: 

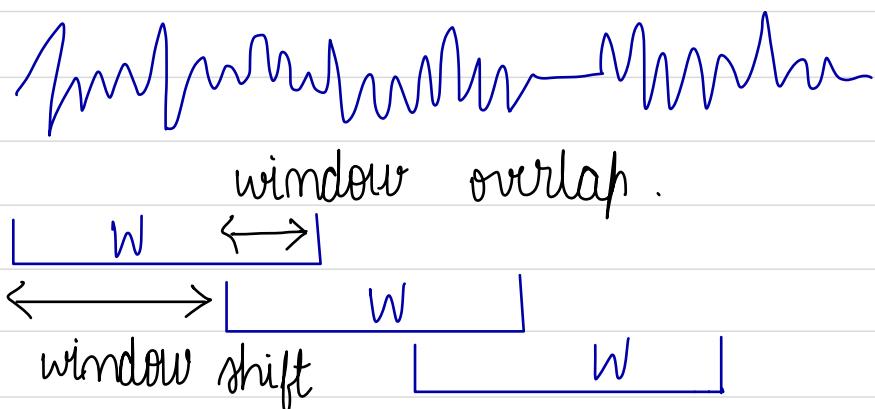
stationary signal: A signal which can be mathematically defined.

• speech signal is non stationary and quasi periodic.

• So to analyze speech signal, we break speech signal to smaller chunks which is called window.

→ windows can be non overlapping and overlapping. Usually overlapping window size is preferred to avoid problems during reconstruction -

• If window size is w , usually window shift is usually $w/2$



Energy of stationary signals:

For a discrete signal $x[n]$, its energy is given by : $E[n] = \sum_{n=-\infty}^{\infty} x^2[n]$

$$\text{If frame size } = N \text{ then energy of frame} \\ E[n] = \sum_{m=n}^{n+N-1} x^2[m]$$

Eg: For a window of $W = 25 \text{ ms}$ and sampling frequency, $f_s = 16 \text{ kHz}$. The no. of samples in the window is ?

$$\text{Ans} \quad \text{No. of samples} = \frac{25 \times 16000}{1000} = 400 \text{ samples.}$$

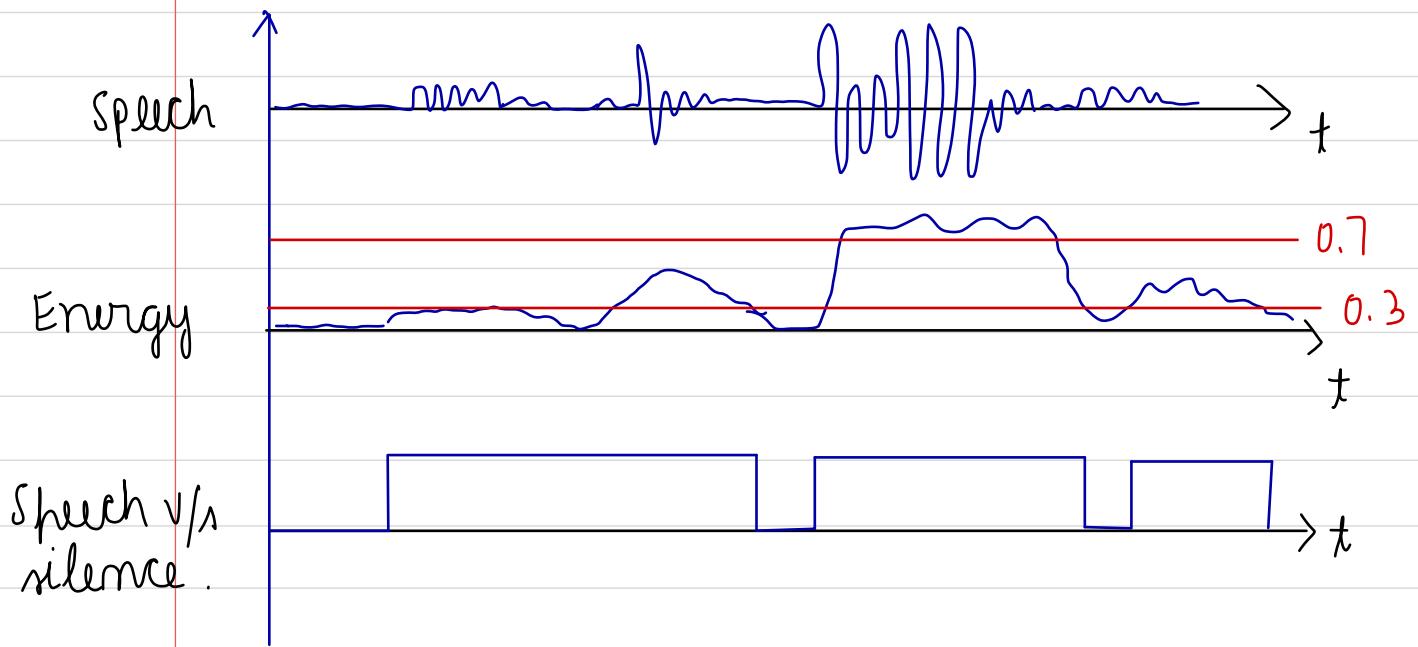
\therefore Frame size , $N = 400$ samples .

Rectangular window

$$w_n(n) = \begin{cases} 1 & , 0 \leq n \leq N \\ 0 & , \text{otherwise} \end{cases}$$

$$E[n] = \sum_{m=-\infty}^{\infty} \{ x[m] \cdot w[n-m] \}^2$$

For overlapping window $n = N/2$ shift is generally used.

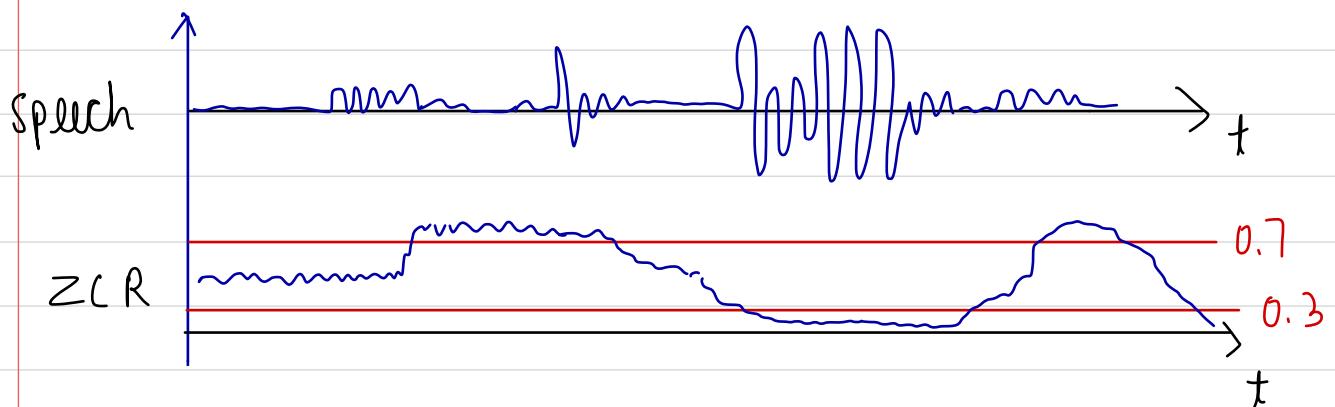


Zero crossing rate (ZCR)

$$ZCR = \frac{\text{No. of zero crossing}}{\text{Total duration}}$$

Zero crossing = No. of sign changes.

Also $(ZCR)_{\text{unvoiced}} > (ZCR)_{\text{silence}} > (ZCR)_{\text{voiced}}$



NOTE: Sigmoid function:

$$\text{Sgm}(x) = \begin{cases} 1 & x(n) > 0 \\ -1 & x(n) < 0 \end{cases}$$

For stationary signal:

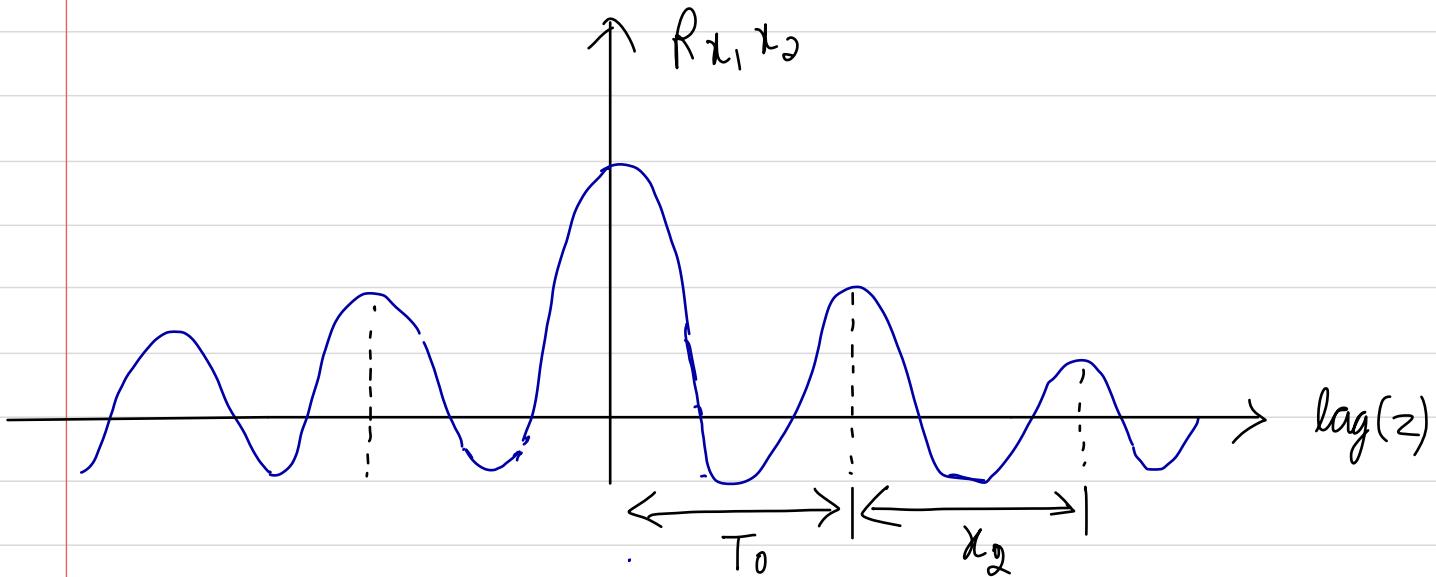
$$ZC = \sum_{m=-\infty}^{\infty} \frac{|\operatorname{sgn}[x(m)] - \operatorname{sgn}[x(m-1)]|}{2}$$

For total M samples:

$$ZCR = \frac{ZC}{M} = \frac{1}{M} \sum_{m=-\infty}^{\infty} \frac{|\operatorname{sgn}[x(m)] - \operatorname{sgn}[x(m-1)]|}{2}$$

For a windowed signal:

$$ZCR(n) = \sum_{m=n-N/2}^{n+N/2} \frac{1}{2} |\operatorname{sgn}[x(m) \cdot w(n-m)] - \operatorname{sgn}[x(m-1) \cdot w(n-m+1)]|$$



Convolution:

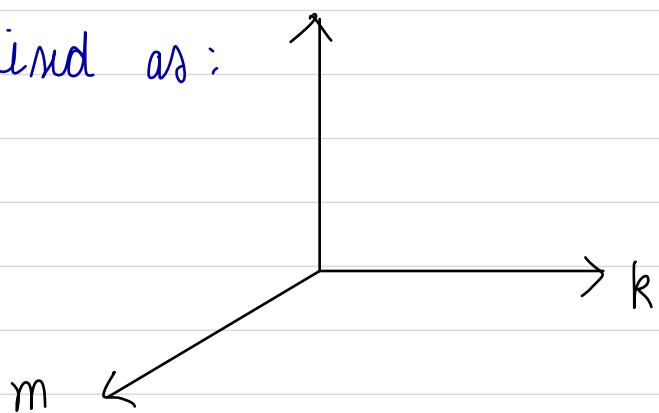
$$x_1[n] * x_2[n] = \sum_{m=-\infty}^{\infty} x_1[m] x_2[n-m]$$

Auto Correlation: $R_{x x}[k] = \sum_{n=-\infty}^{\infty} x[n] w[m-n] \cdot x[n+k] w[m-(n+k)]$

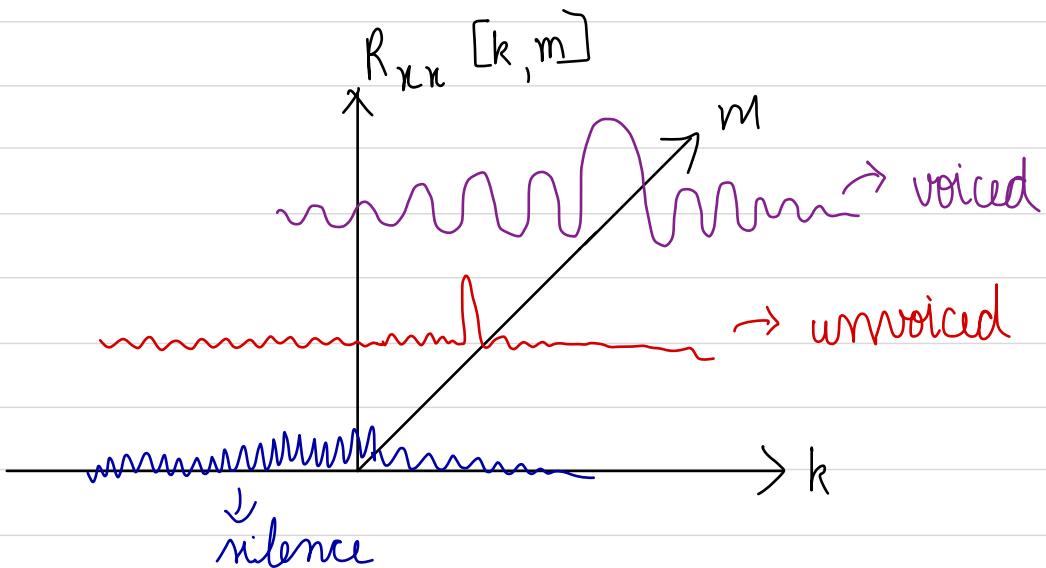
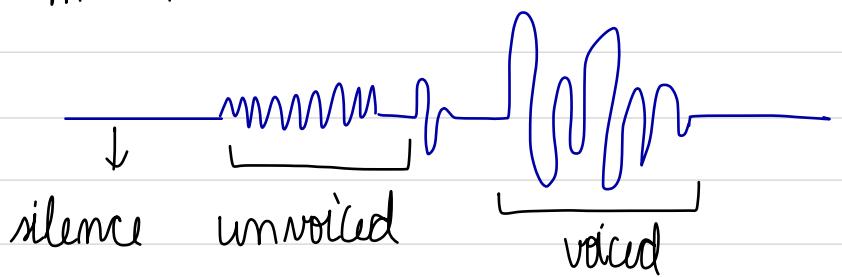
↓
non stationary signal

This can be visualized as:

$$R_{xx} [k, m]$$



Consider signal:



Short time Fourier transform (STFT)

$$X(\omega) = \sum_{n=-\infty}^{\infty} x(n) e^{-j\omega n}$$

$$X(\omega, m) = \sum_{n=-\infty}^{\infty} x(n) w(m-n) e^{-j\omega n}$$

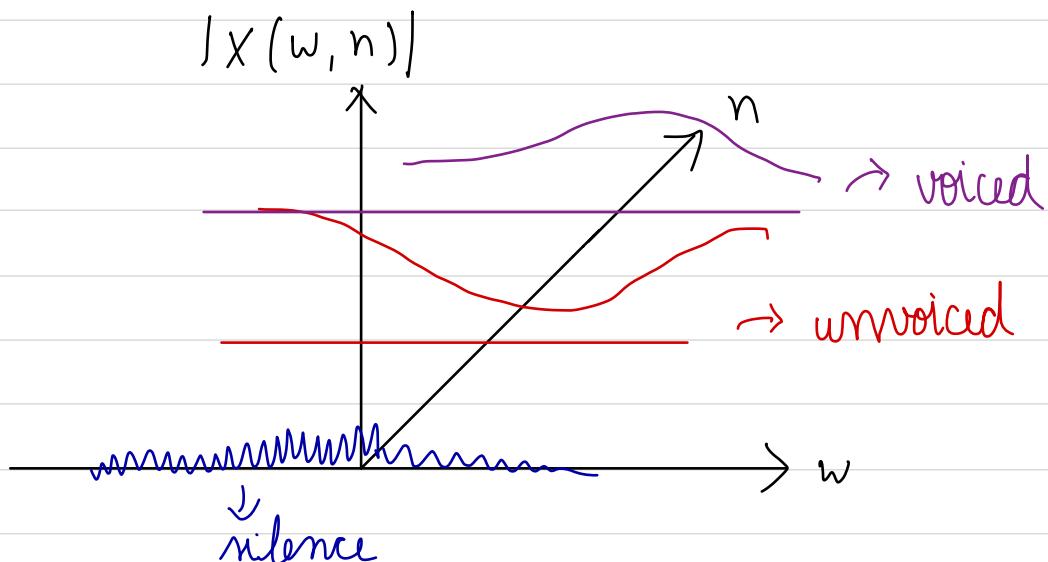
$$X(\omega_0, n) = \sum_{m=-\infty}^{\infty} x(m) w(n-m) e^{-j\omega_0 m}$$



Let $p = n - m$, $p \rightarrow -\infty$ to ∞

$$\begin{aligned}
 X(\omega_0, n) &= \sum_{p=-\infty}^{\infty} x(n-p) w(p) e^{-j\omega_0(n-p)} \\
 &= e^{-j\omega_0 n} \sum_{p=-\infty}^{\infty} x(n-p) w(p) e^{j\omega_0 p} \\
 &= e^{-j\omega_0 n} \sum_{p=-\infty}^{\infty} x(n-p) h(p)
 \end{aligned}$$

$h(p)$ is impulse response. So $w(p) e^{j\omega_0 p}$ acts as a impulse response.

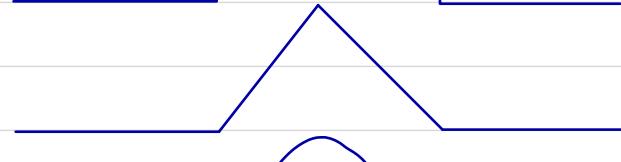


Types of windows

Rectangular



Triangular



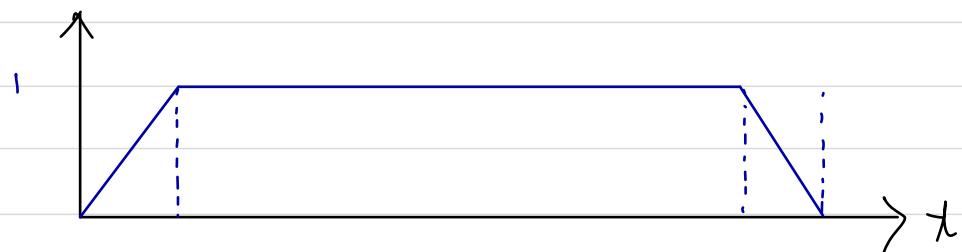
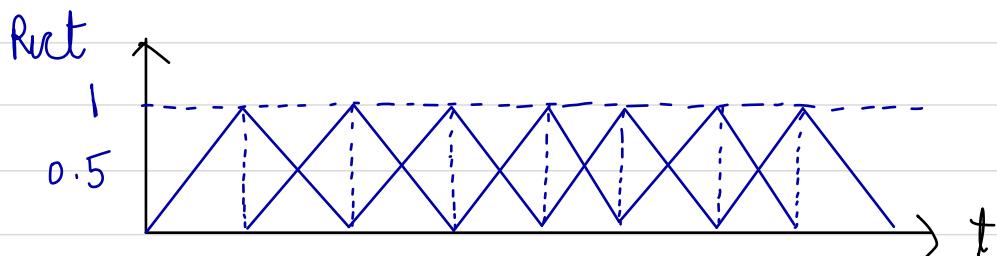
Hamming



Hamming



Rect



STFT Spectral analysis

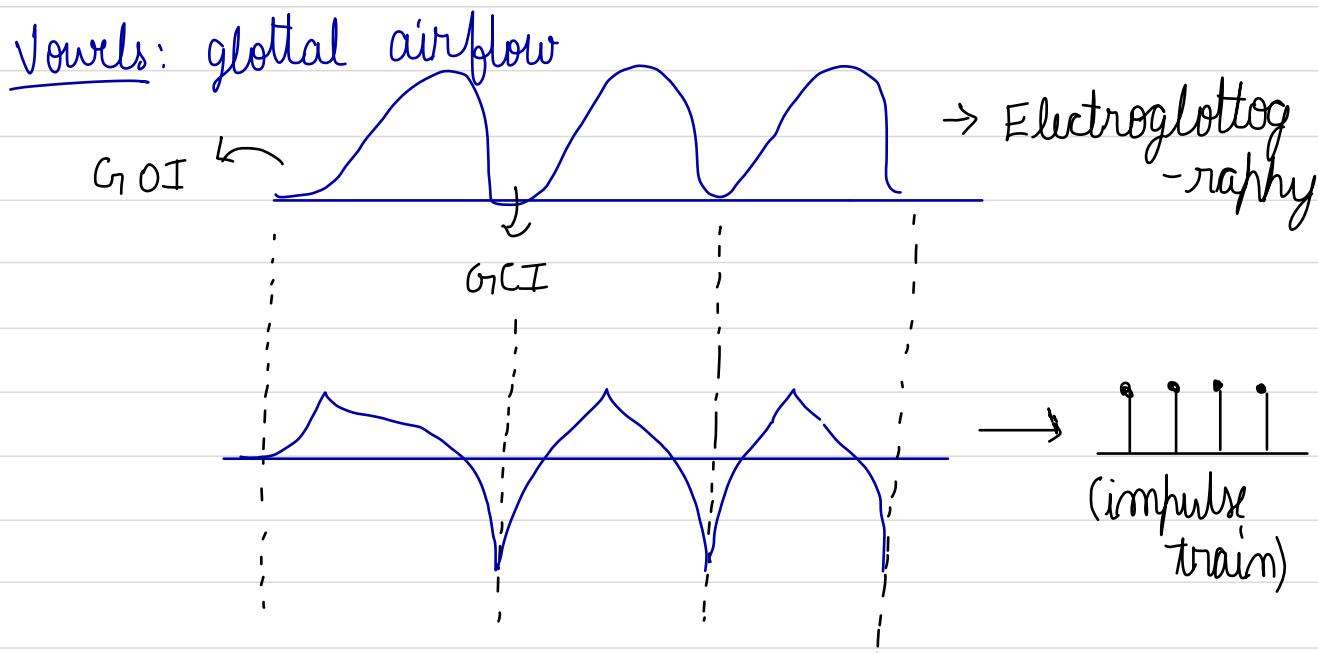
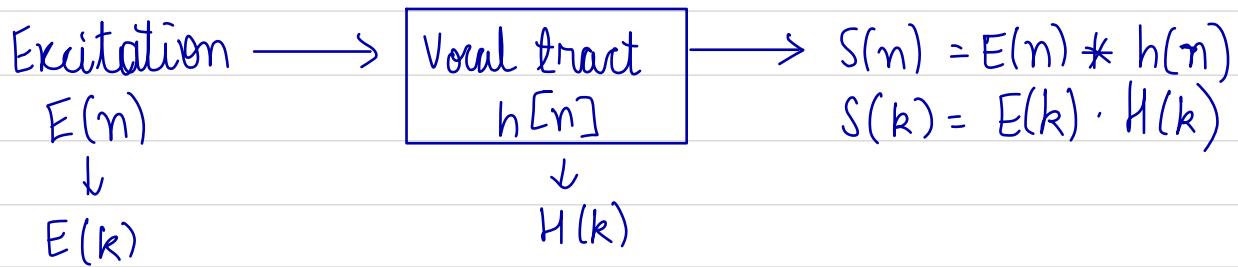
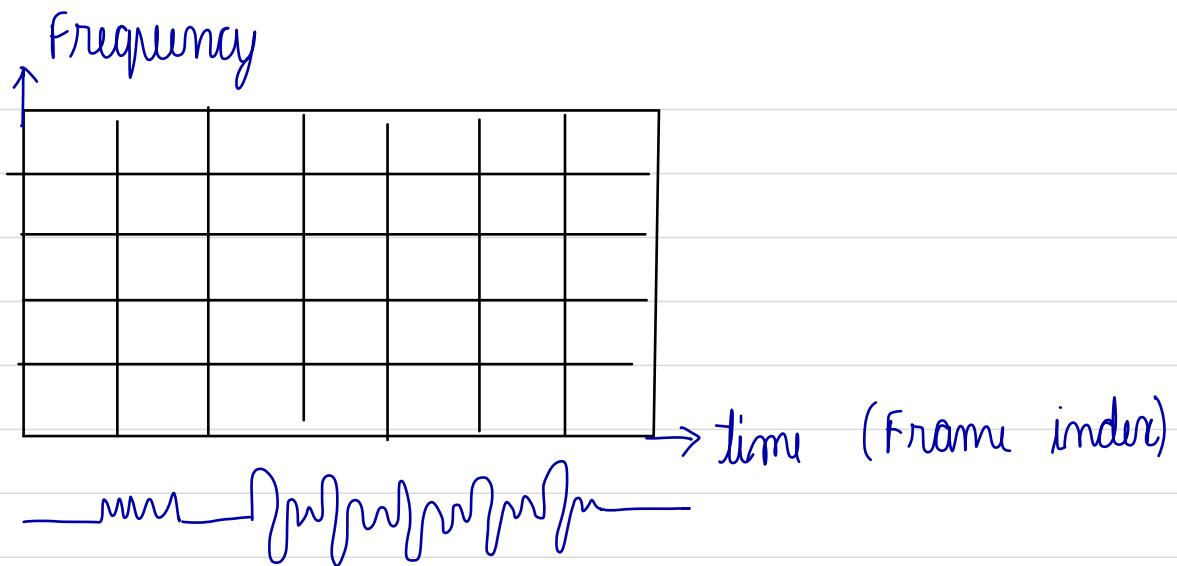
$x(t)$ Sampling \rightarrow
Time domain

$x[n]$

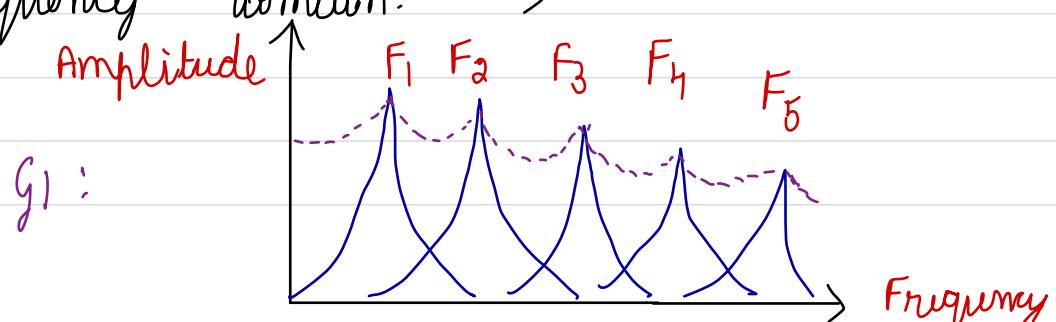
↓
Aliasing in frequency
domain

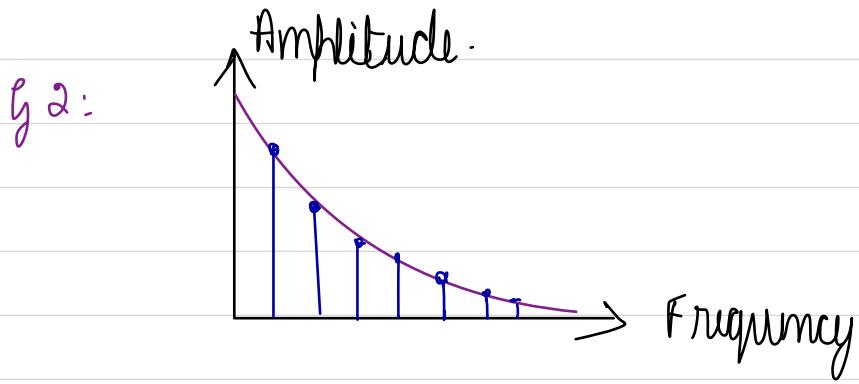
$x(t)$ $\xrightarrow{\text{DTFT}}$ $X(w)$ $\xrightarrow{\text{DFT}}$ $X(k)$

↓
Aliasing in Time
domain

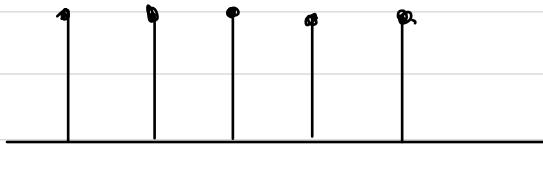
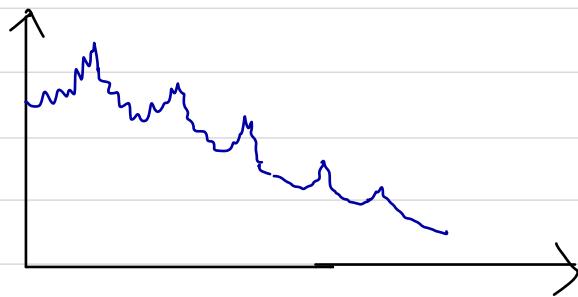


G_{OI} = glottal opening instance
 G_{CI} = glottal closing instance
 f_m frequency domain: \rightarrow

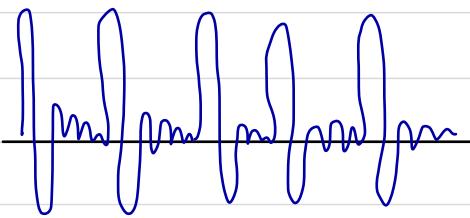




When $y_1 \times y_2$:

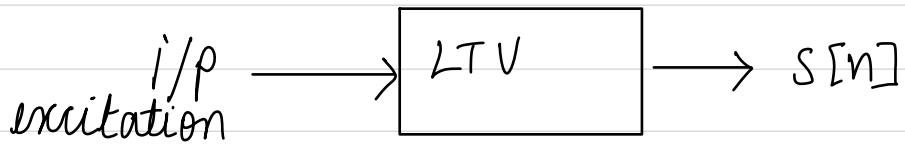


\rightarrow Impulse train



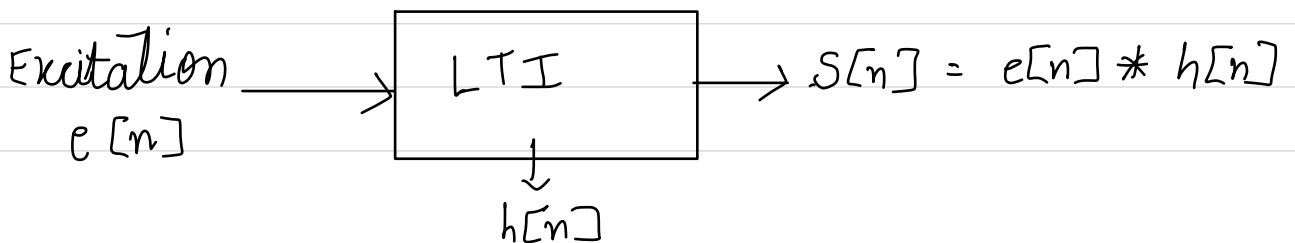
Glottal pattern

Homomorphic speech analysis

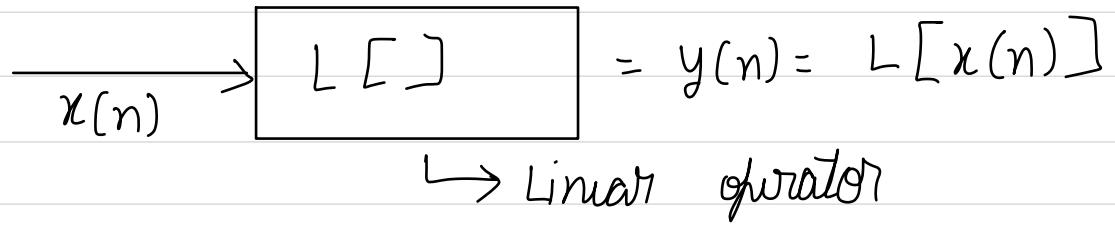
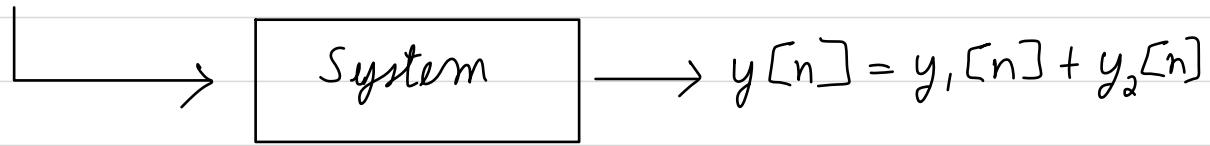


↳ slowly varies with time
Linear time variant

Linear time invariant.



$$x[n] = x_1[n] + x_2[n]$$



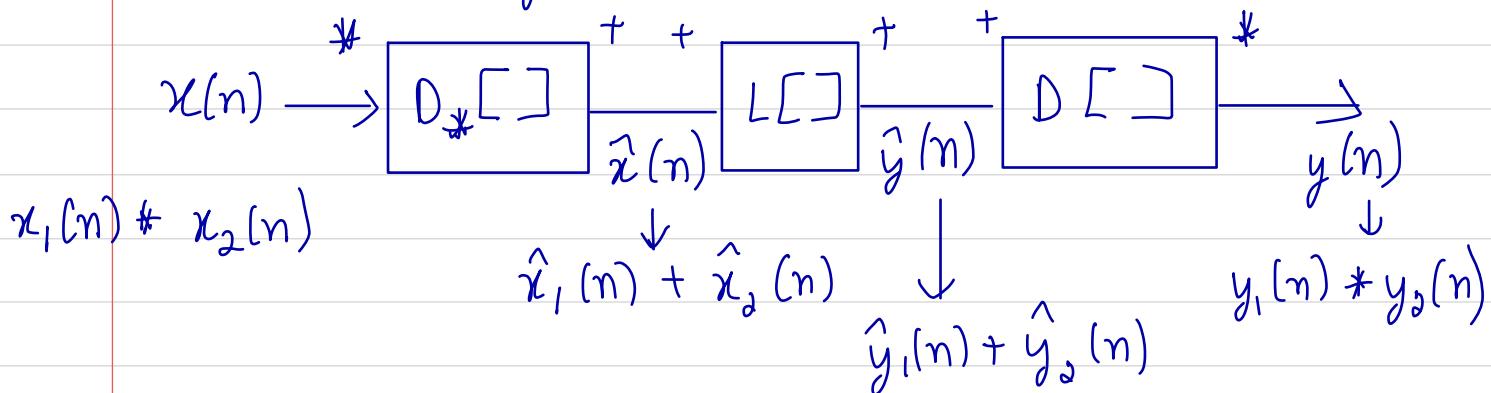
$$\begin{aligned} L(x[n]) &= L(x_1[n] + x_2[n]) \\ &= L(x_1[n]) + L(x_2[n]) \\ &= y_1(n) + y_2(n) \end{aligned}$$

$$y[n] = \sum_{k=-\infty}^{\infty} x(k) h(n-k) = h(n) * x(n)$$

$$H[x(n)] = H[x_1(n)] * H[x_2(n)]$$

$$= H[x_1(n)] * H[x_2(n)]$$

$$= y_1(n) * y_2(n)$$

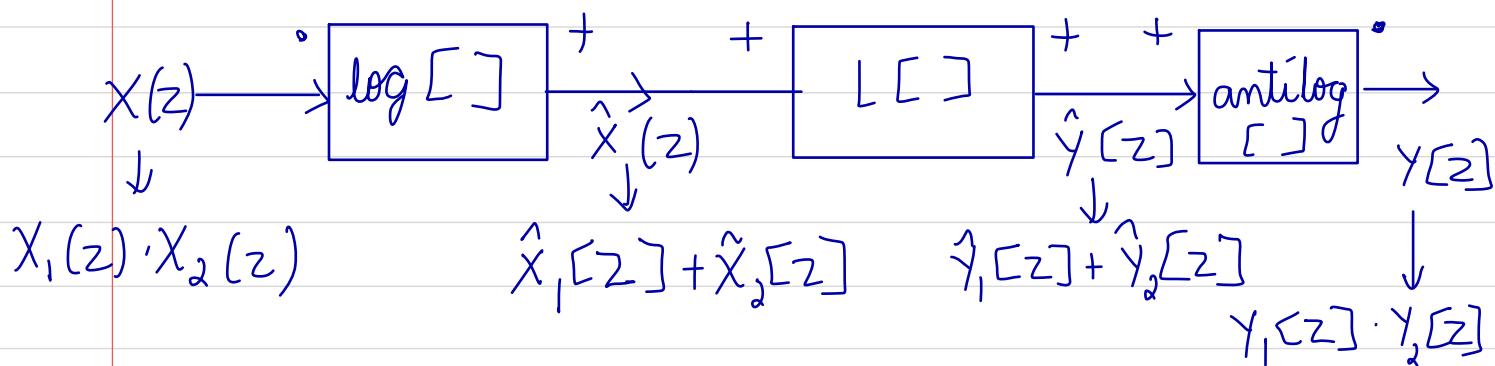


$$\begin{aligned}
 D_*(x(n)) &= D_*[x_1(n) * x_2(n)] \\
 &= D_*[x_1(n)] + D_*[x_2(n)] \\
 &= \hat{x}_1(n) + \hat{x}_2(n)
 \end{aligned}$$

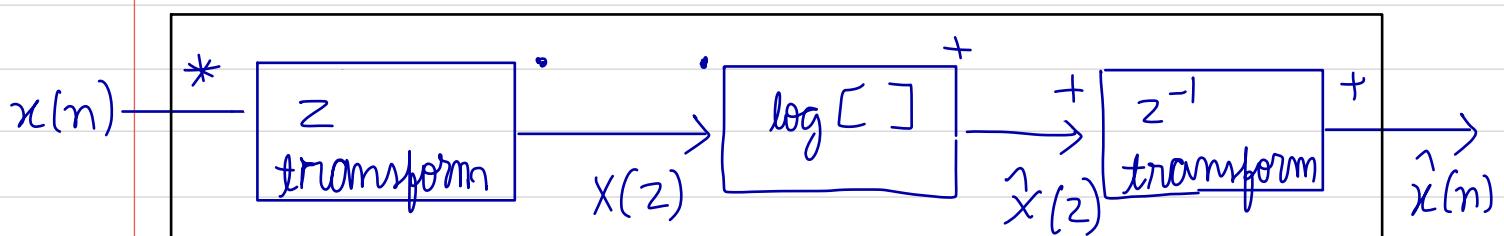
eg: $a * b \xrightarrow{\text{DFT}} A[k] \cdot B[k]$

$$\log(A[k]) + \log(B[k])$$

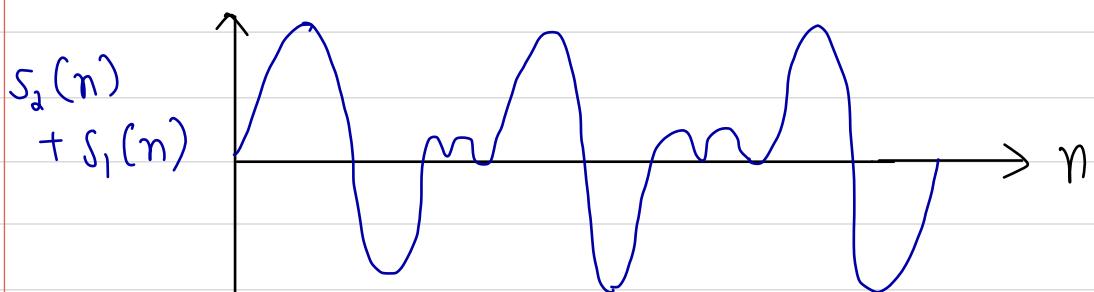
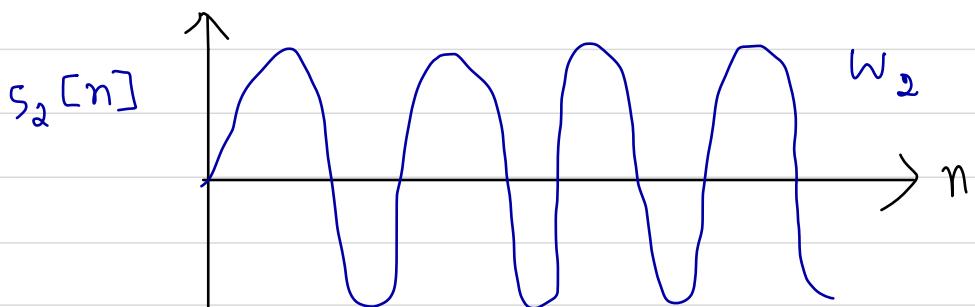
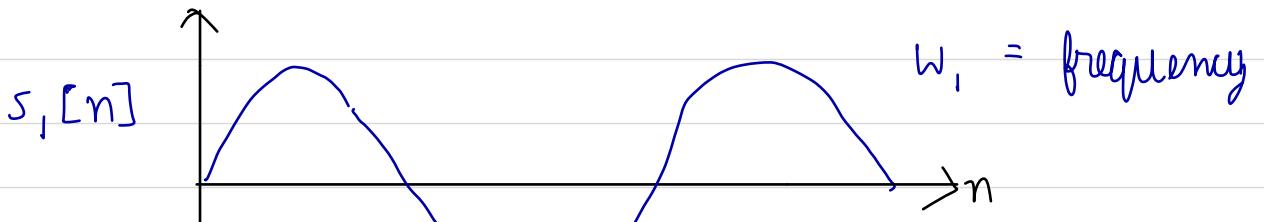
Frequency representation w.r.t convolution:



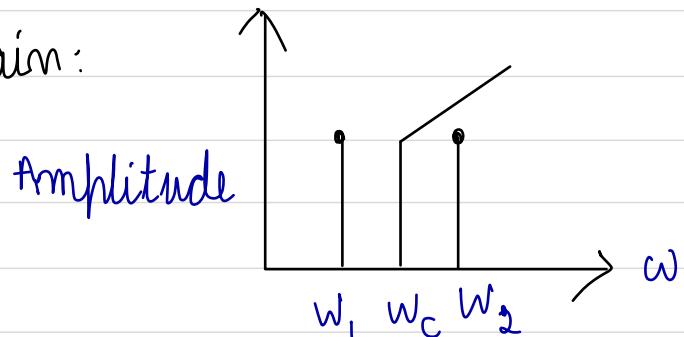
$$D_*^{-1}[\hat{y}(n)] = D_*^{-1}[\hat{y}_1(n) + \hat{y}_2(n)]$$



Homomorphic deconvolution



In frequency domain:

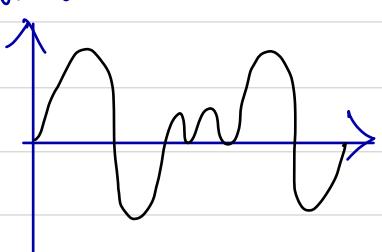


$$s[n] = e[n] * h[n]$$

$$S(\omega) = E(\omega) \cdot H(\omega)$$

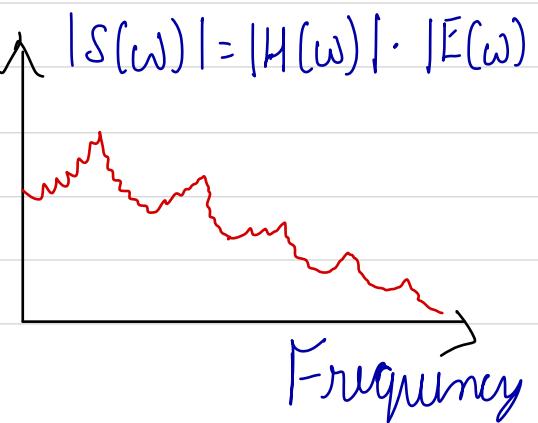
$$= \log |E(\omega)| + \log |H(\omega)|$$

voiced



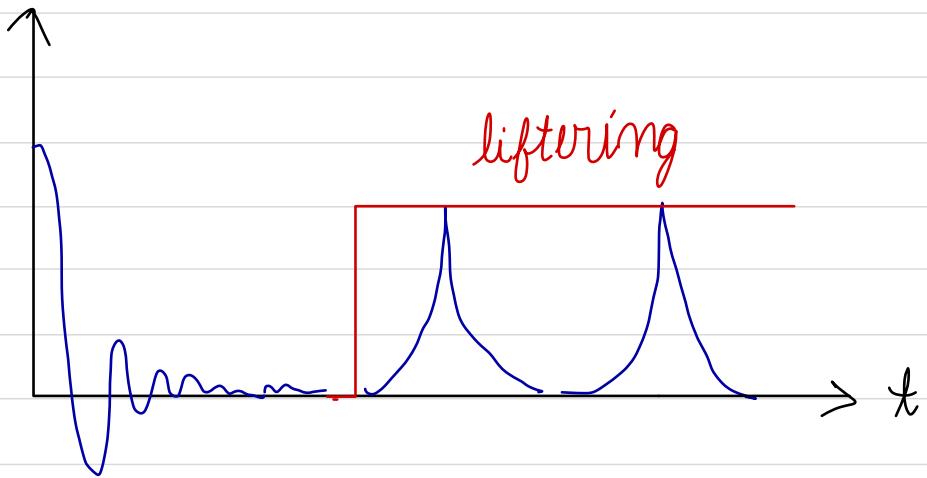
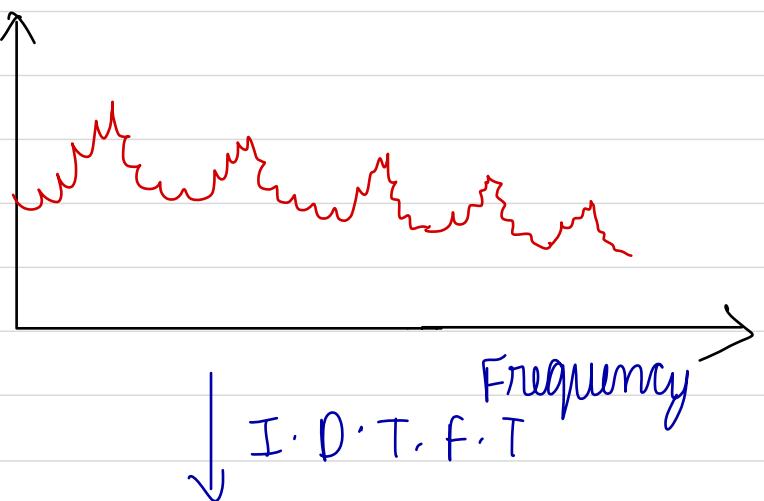
D.T.F.T

$$|S(\omega)| = |H(\omega)| \cdot |E(\omega)|$$



taking $\log |s(\omega)|$

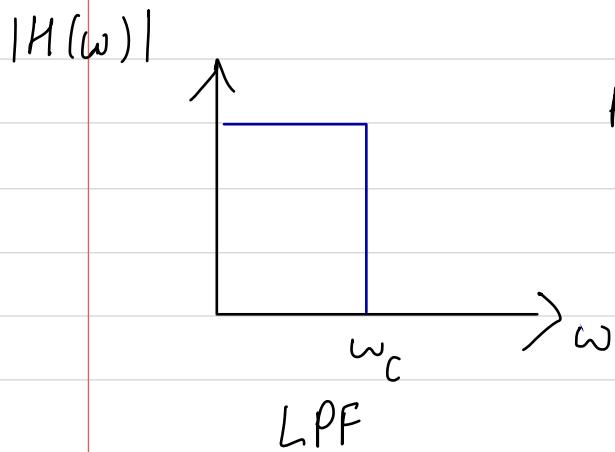
$$\log |s(\omega)|$$



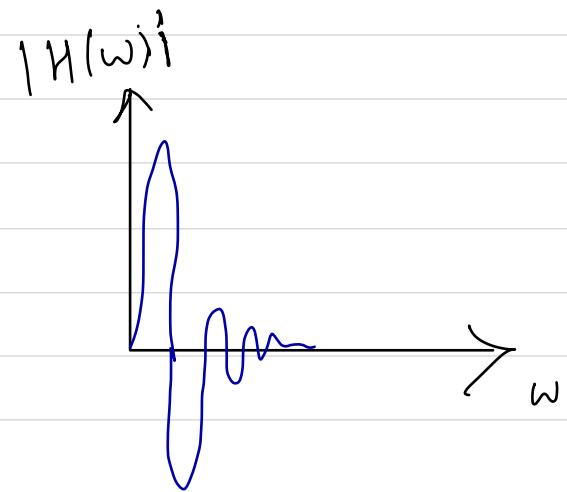
Cepstrum:

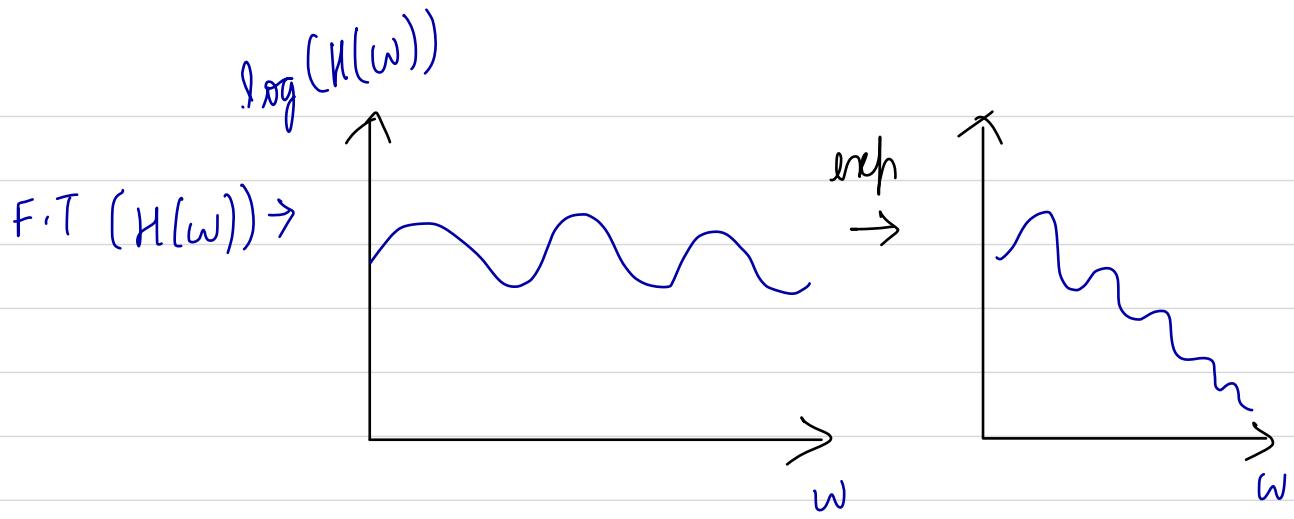
$$c(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log |s(\omega)| \cdot e^{j\omega n} d\omega$$

Real part of cepstrum

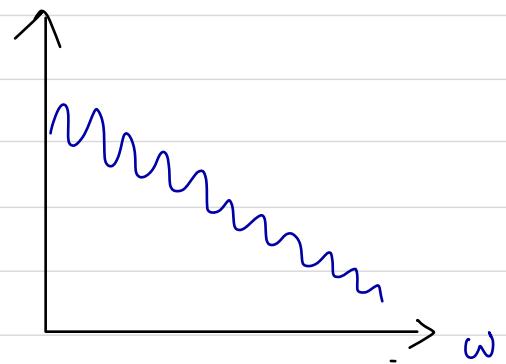
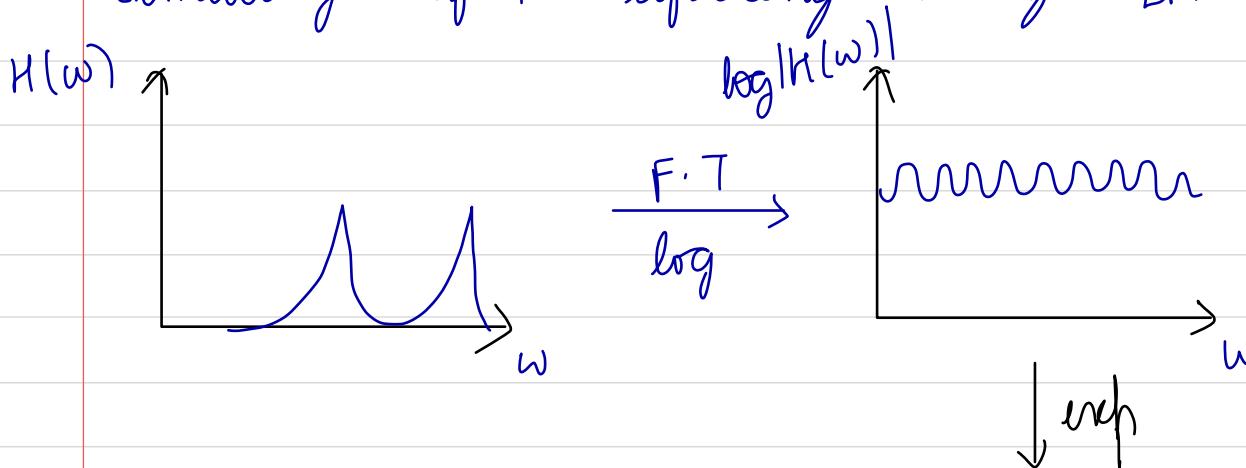


After
liftering
using
H.P.F

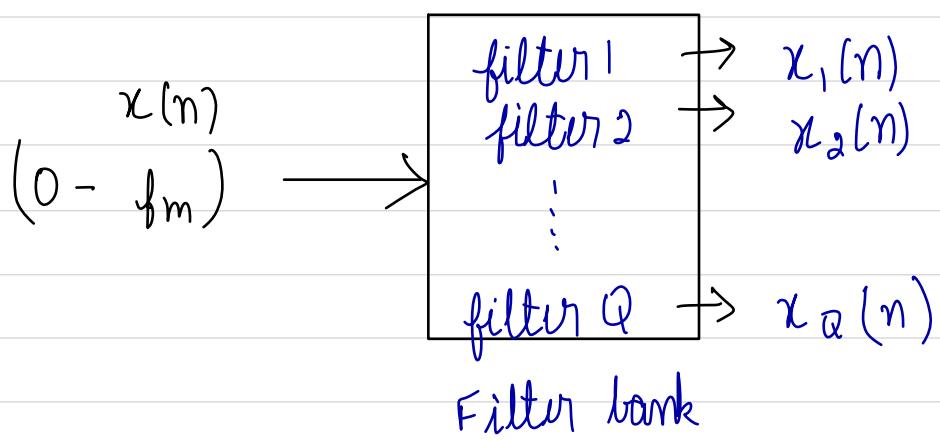




Similarly after lifting using LPF



Filter bank analysis



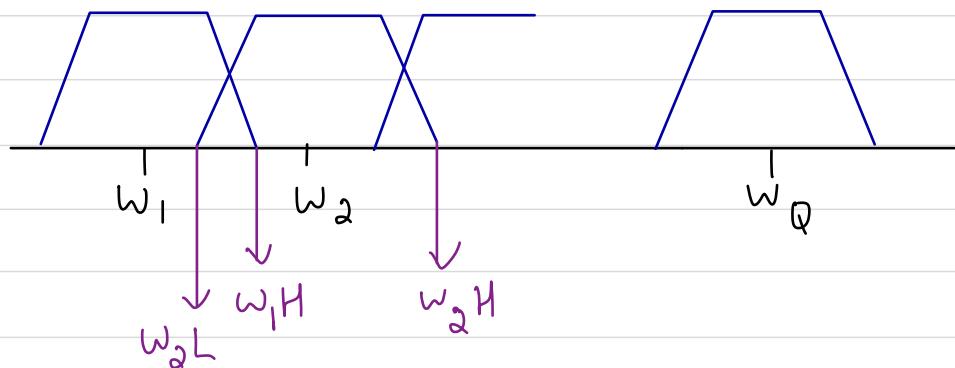
Lifting: Process of filtering out components outside the filter used. (choosing signal outside the filter)

Q = Number of bandpass filters.

$$Q = \frac{f_m}{B \cdot W} \quad \text{or} \quad \frac{f_s}{2 \times BW}$$

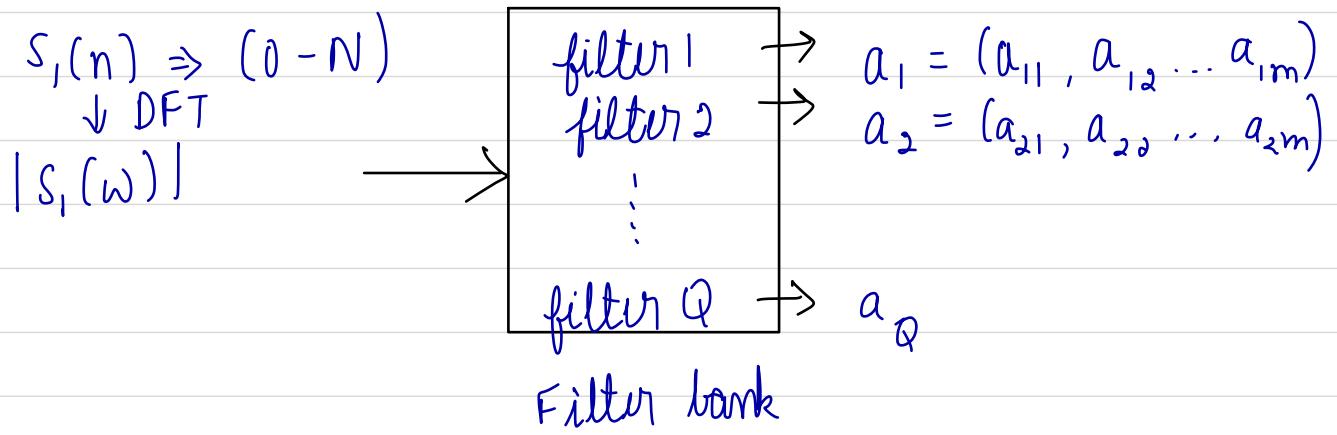
Eg: if $f_s = 16000\text{Hz}$ & $BW = 100\text{Hz}$

$$Q = \frac{16000}{2 \times 100} = 80 //$$



For non overlapping filters : $\omega_1H = \omega_2L$,
 $\omega_2H = \omega_3L$ & so on

NOTE: No. of frames = $1 + \frac{\text{len(signal)} - \text{len(frame size)}}{\text{len(hop size)}}$



Energy:

$$\sum_{i=1}^m a_{1i}^2$$

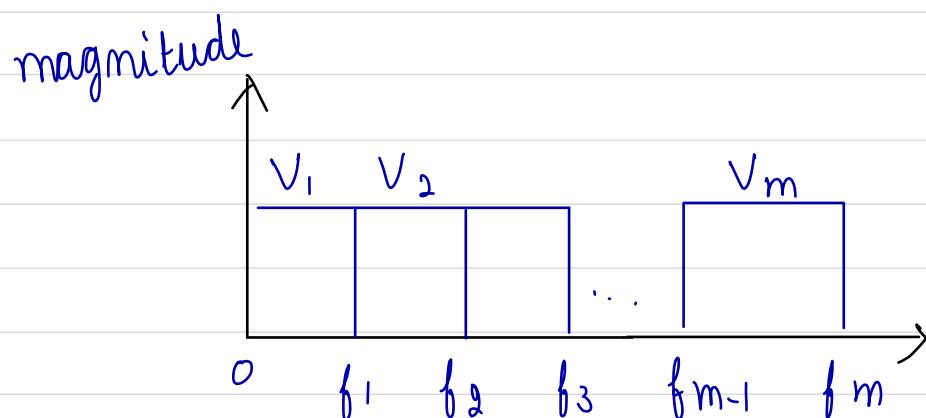
$$\sum_{i=1}^m a_{2i}^2$$

$$\vdots$$

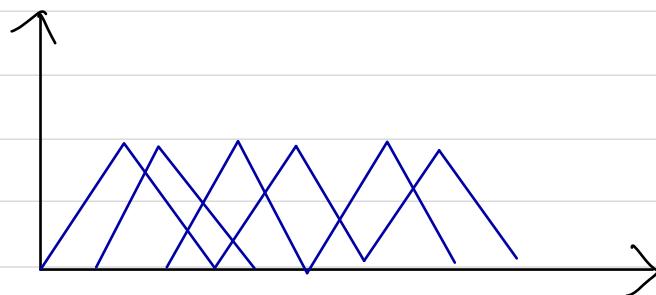
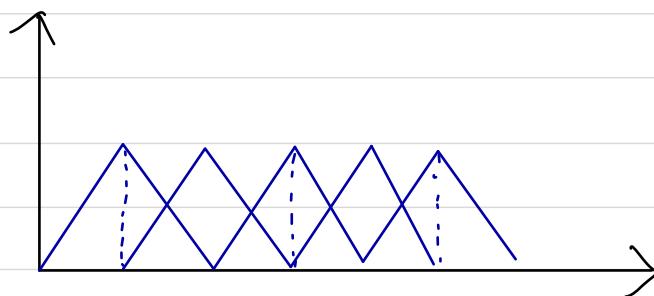
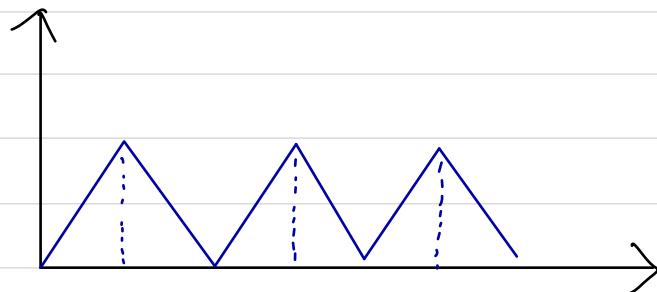
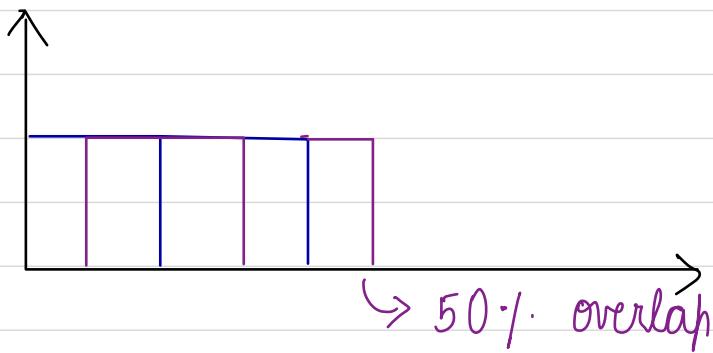
$$\sum_{i=1}^m a_{Qi}^2$$

$\rightarrow \begin{bmatrix} & \\ & \\ & \end{bmatrix}_{Q \times 1}$

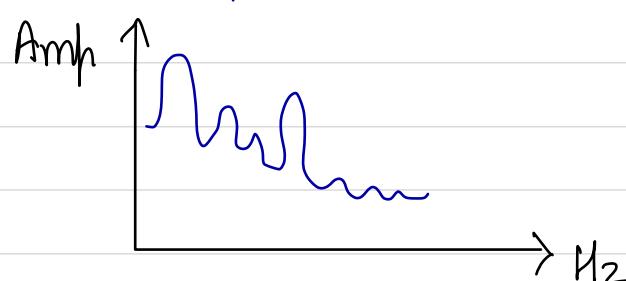
Uniform filter bank



Mel frequency cepstrum



$$s[n] \xrightarrow{\text{DFT}} |s[k]|$$



Hz \rightarrow Physical frequency

Mel \rightarrow Perceived frequency

$$f_{\text{mel}} = 1125 \log_e \left(1 + \frac{f_{\text{Hz}}}{700} \right)$$

$$f_{\text{mel}} = 2595 \log_{10} \left(1 + \frac{f_{\text{Hz}}}{700} \right)$$

e.g.: If $f_{\text{Hz}} = 100 \text{ Hz}$, find f_{mel}

$$\begin{aligned} f_{\text{mel}} &= 1125 \log_e \left(1 + \frac{100}{700} \right) \\ &= 150.22 \end{aligned}$$

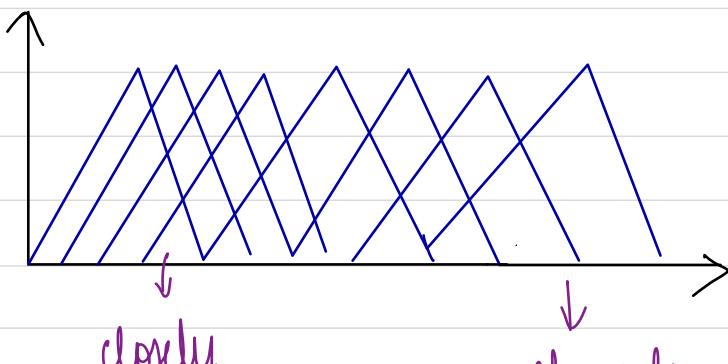
NOTE.

f_{Hz}	f_{mel}
100	150.2
150	218.42
200	282.72
1500	1288.27
1550	1313.55
1600	1338.28

f_{mel}	f_{Hz}
100	64.95
150	99.65
200	135.92
1500	1949.31
1550	2069.49
1600	2195.13

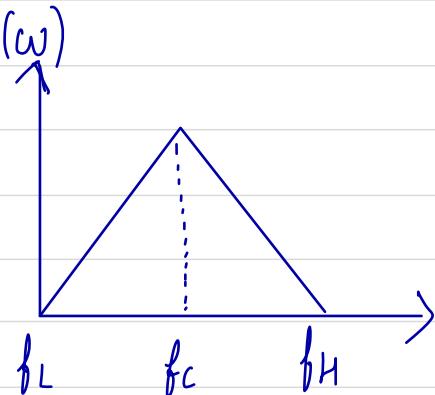
Mel filter bank:

No. of filters = 26



Lower sensitivity for higher frequency.

Consider one filter: $H(\omega)$



$$H(\omega) = \begin{cases} \frac{f - f_L}{f_C - f_L}, & f_L \leq f \leq f_C \\ \frac{f_H - f}{f_H - f_C}, & f_C < f \leq f_H \\ 0, & \text{otherwise} \end{cases}$$

Lipstrum analysis

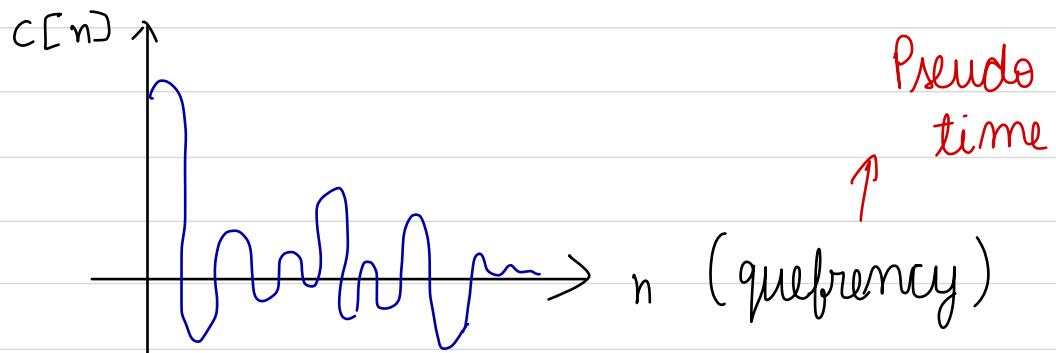
$$S[n] = h[n] * e[n]$$

↓ Homomorphic processing

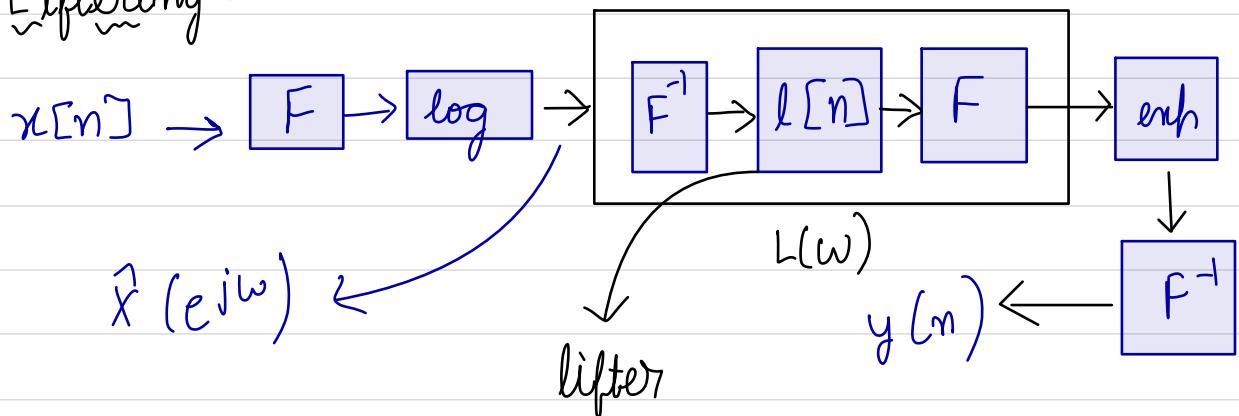
$$\hat{h}[n] + \hat{e}[n]$$

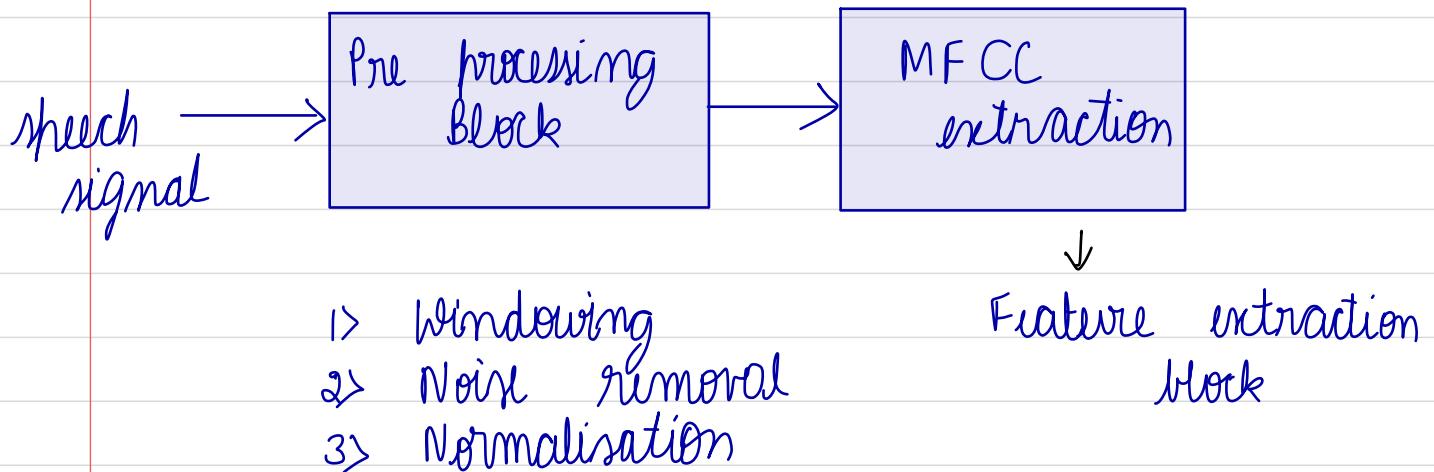
$$\hat{h}[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log \hat{x}(e^{j\omega}) e^{j\omega n} d\omega \rightarrow \text{complex}$$

$$c[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log |\hat{x}(e^{j\omega})| e^{j\omega n} d\omega \rightarrow \text{real}$$



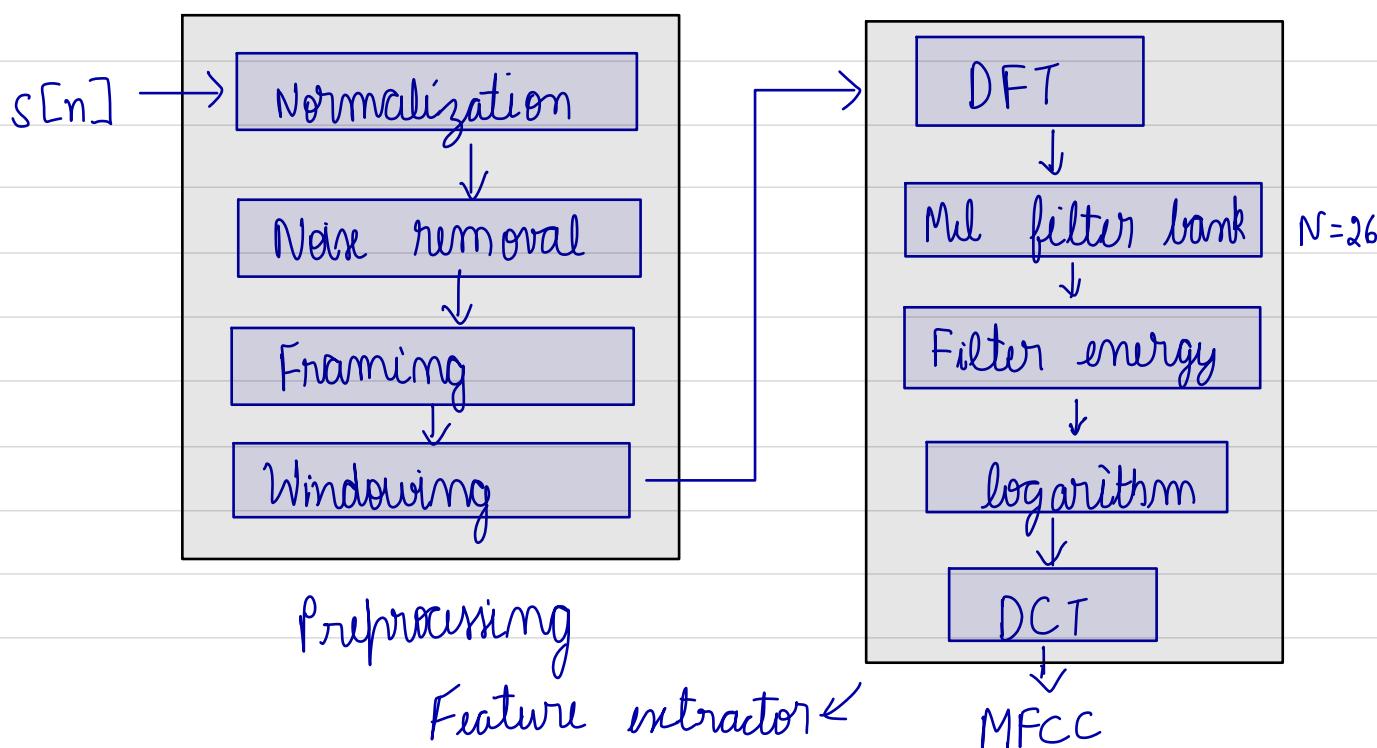
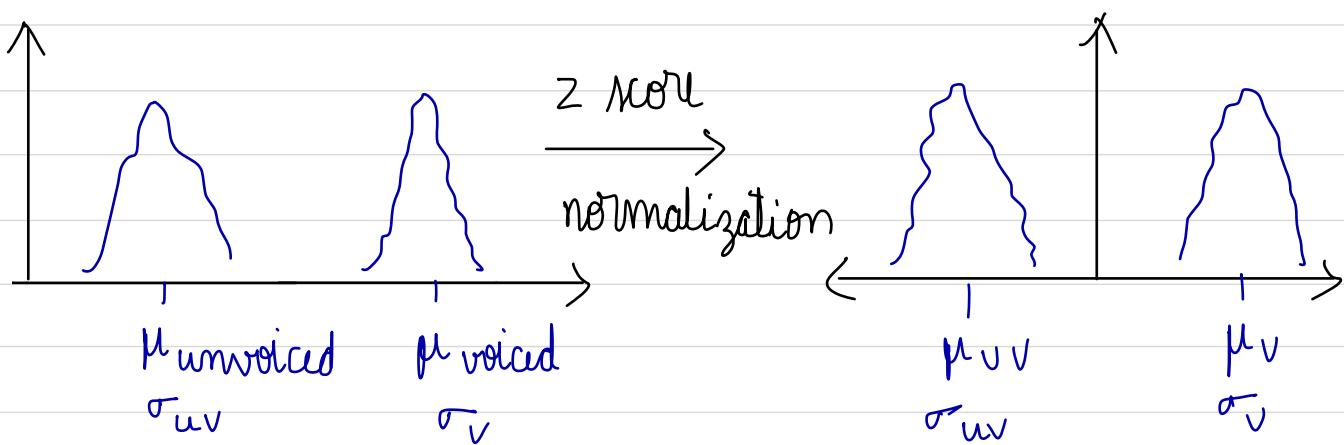
Lifting:





Z score normalization:

$$z = \frac{(x - \mu)}{\sigma}$$



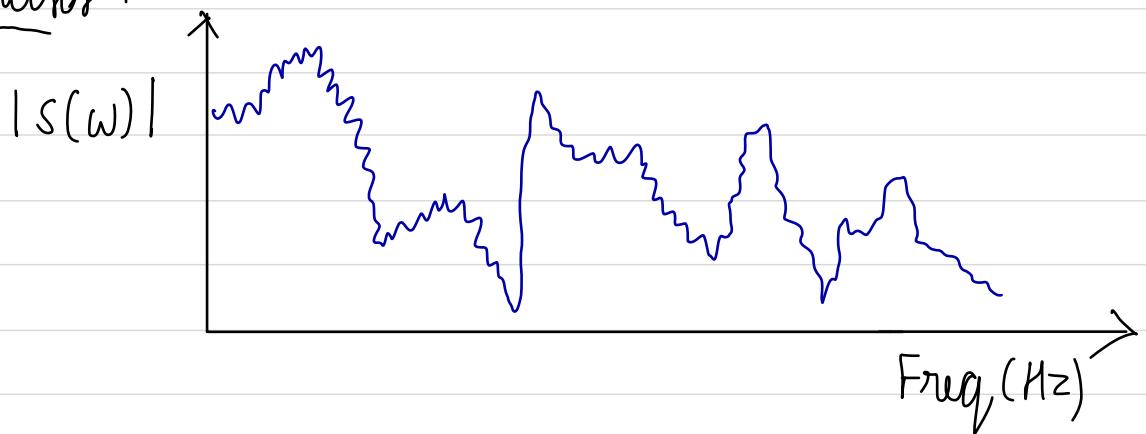
$$\text{Windows: Hamming: } 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right)$$

$$\text{Hamming: } 0.5 \left[1 - \cos\left(\frac{2\pi n}{N-1}\right) \right]$$

$$\text{MFCC delta: } mfcc_{i+1} - mfcc_i$$

$$\text{MFCC delta delta: } (mfcc_{i+1} - mfcc_i) - (mfcc_i - mfcc_{i-1})$$

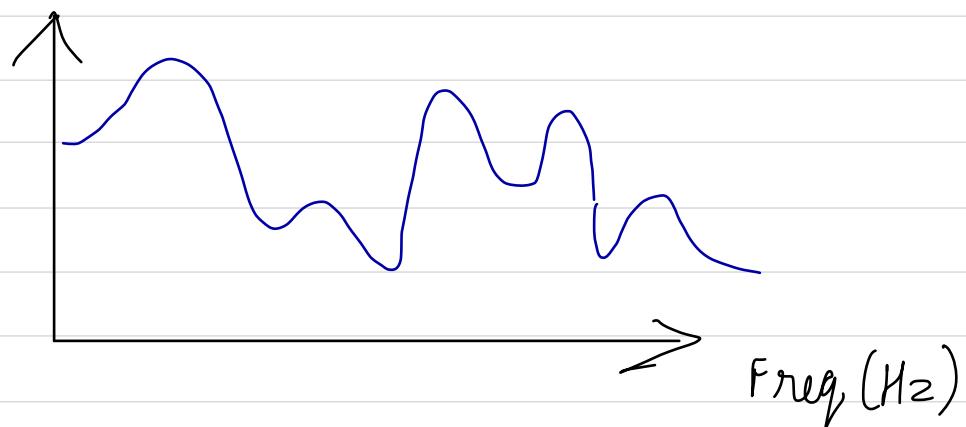
Preemphasis:



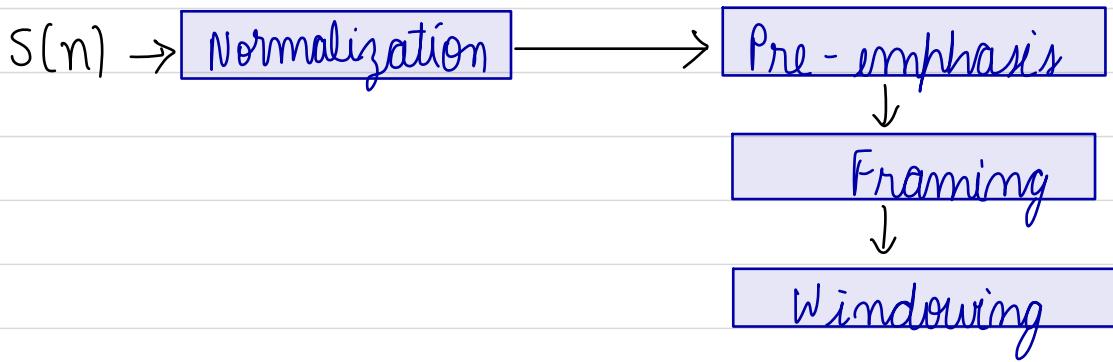
$$P(z) = (1 - az^{-1}) \quad (\text{acts as HPF})$$

$a = 0.68$ is mostly used for speech.

→ Passing above signal through HPF



'a' decides amount of pre-emphasis



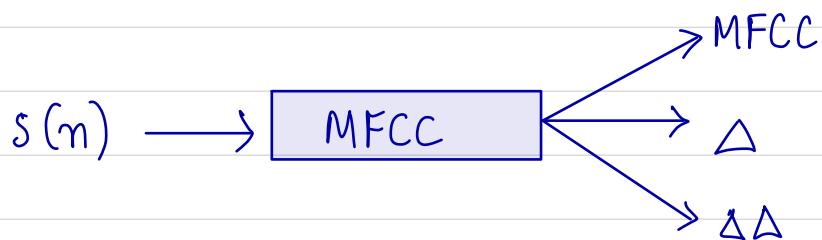
MFCC_n

$$\text{MFCC} = [c_1, c_2, c_3, \dots, c_N] \quad N \leq 26$$

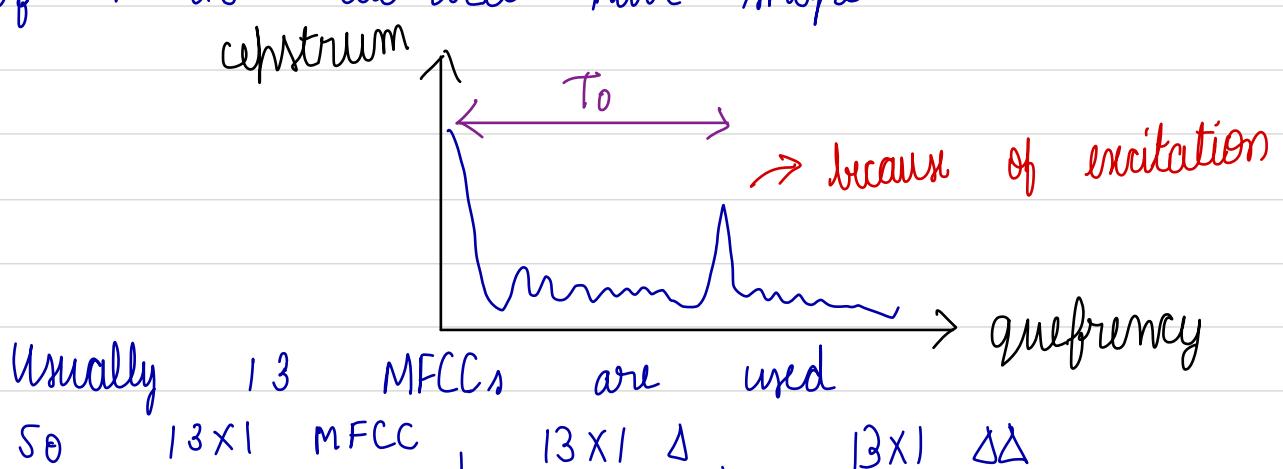
$$\Delta_i = c_{i+1} - c_{i-1}$$

$$\Delta = \frac{\sum_{n=1}^N n (c_{n+i} - c_{n-i})}{2 \sum_{n=1}^N n^2}$$

$$\Delta\Delta_i = \Delta_{k+i} - \Delta_{k-i}$$



If $N = 26$ all will have shape 26×1



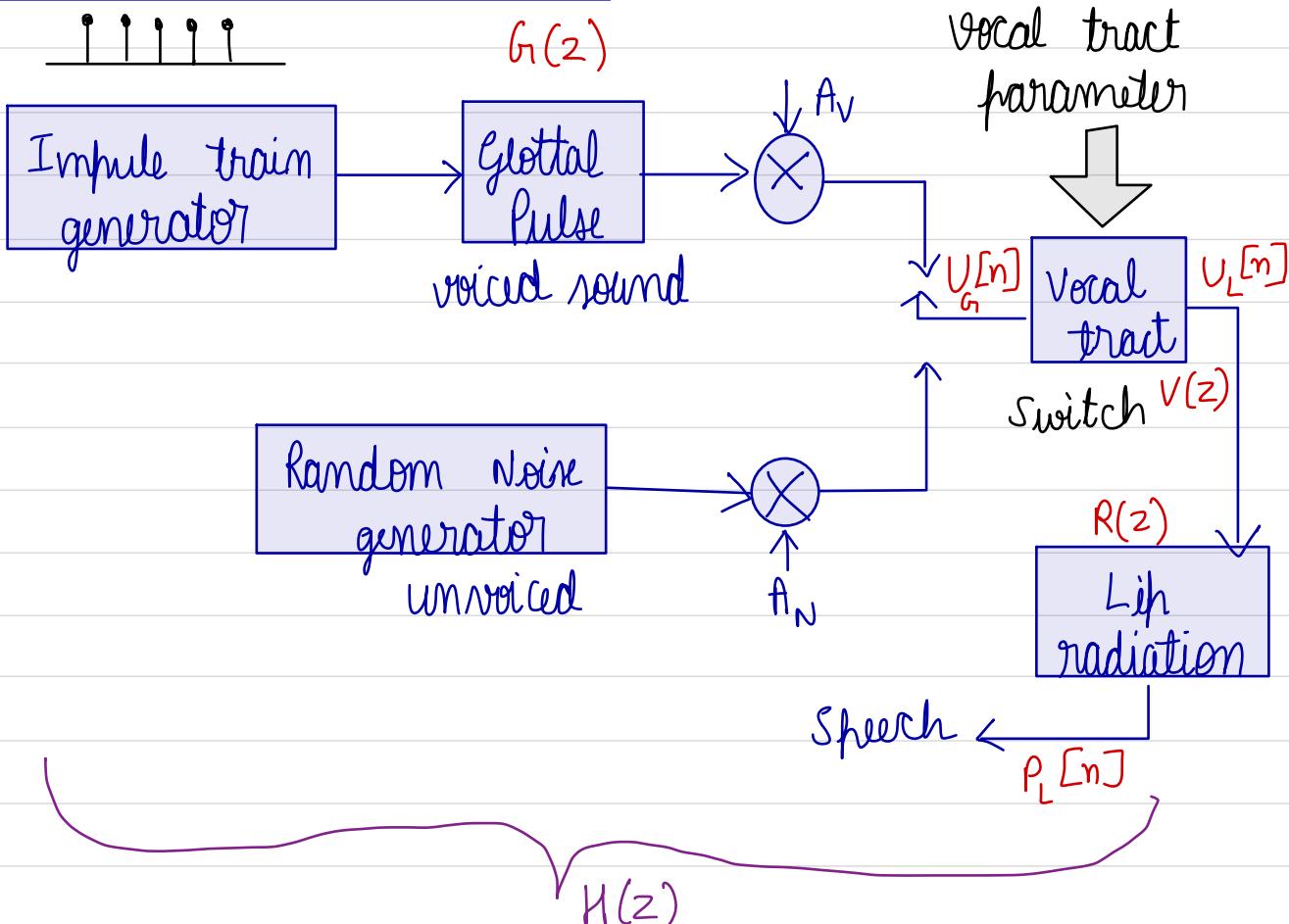
$$\text{MFCC}^T_{\text{full}} = [\text{MFCC}, \Delta, \Delta\Delta]_{1 \times 39}$$

$$\text{Fundamental frequency} = \frac{1}{T_0} \text{ (from MFCCs)}$$

→ This method of finding fundamental frequency is more robust to noise than auto correlation method.

Linear Prediction coefficient: (LPC)

Speech Production Model:



$$H(z) = G(z) \cdot V(z) \cdot R(z)$$

where $G(z)$ = glottal transfer function

$V(z)$ = Vocal Tract transfer function

$R(z)$ = Lip radiation transfer function

$$G(z) = \frac{1}{(1 - e^{-CT} z^{-1})^2} \quad e^{-CT} \approx 1$$

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

$$V(z) = \frac{G}{\prod_{k=1}^{N/2} (1 - 2\gamma_k \cos \theta_k z^{-1} + \gamma_k^2 z^{-2})^2}$$

$$R(z) = R_0 (1 - z^{-1})$$

$$H(z) = \frac{G R_0}{1 - \sum_{k=1}^P a_k z^{-k}} \quad \text{let } G R_0 = r$$

$$= \frac{r}{1 - \sum_{k=1}^P a_k z^{-k}} = \frac{S(z)}{U(z)}$$

Now $S(z) = H(z) \cdot U(z)$

$$S(z) \left(1 - \sum_{k=1}^P a_k z^{-k} \right) = r U(z)$$

$$S(z) = \sum_{k=1}^P a_k z^{-k} S(z) + r U(z)$$

Taking inverse z transform

$$s[n] = \sum_{k=1}^P a_k s[n-k] + r u(n)$$

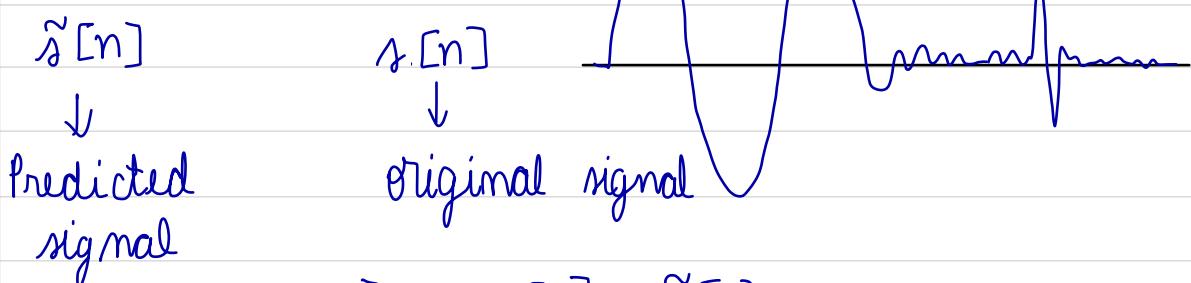
when $u[n] < 0$

$$s[n] = \sum_{k=1}^p a_k s[n-k] \Rightarrow a_1 s[n-1] + a_2 s[n-2] + \dots + a_p s[n-p]$$

↑
current sample

Here $\{a_1, a_2, a_3, \dots, a_p\}$ are called LPC.

Prediction error:



$$e[n] = s[n] - \tilde{s}[n]$$

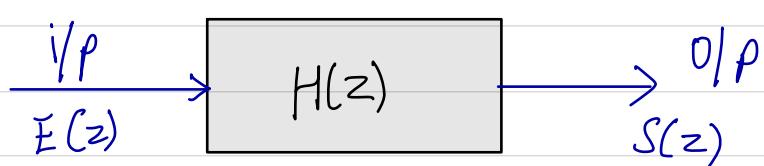
error

$e[n]$ is higher during transition from one region to another region.

$$e[n] = s[n] - \sum_{k=1}^p a_k s[n-k]$$

$$\begin{aligned} E(z) &= S(z) - \sum_{k=1}^p a_k z^{-k} S(z) \\ &= S(z) \left\{ 1 - \sum_{k=1}^p a_k z^{-k} \right\} \\ &= S(z) A(z) \end{aligned}$$

$$\text{where } A(z) = 1 - \sum_{k=1}^p a_k z^{-k}$$



$$H(z) = \frac{S(z)}{E(z)} = \frac{\sum_{k=0}^M b_k z^{-k}}{\sum_{k=0}^N a_k z^{-k}}$$

$$\sum_{k=0}^N a_k z^{-k} s(z) = \sum_{k=0}^M b_k E(z) z^{-k}$$

Taking inverse Z transform:

$$\sum_{k=0}^N a_k s[n-k] = \sum_{k=0}^M b_k e[n-k]$$

$$a_0 s[n] + a_1 s[n-1] + \dots + a_N s[n-N] = \\ b_0 e[n] + b_1 e[n-1] + \dots + b_M e[n-M]$$

Prediction : $\tilde{s}[n] = - \sum_{k=1}^P a_k s[n-k]$

Mean square error:

Instantaneous error:

$$e[n] = s[n] - \tilde{s}[n] \\ = s[n] + \sum_{k=1}^P a_k s[n-k]$$

$$MSE = \sum_n e^2[n] \\ = \sum_n \left(s[n] + \sum_{k=1}^P a_k s[n-k] \right)^2$$

For $k=1$, $\frac{\partial E}{\partial a_1} = 0$

$$\frac{\partial E}{\partial a_1} = \sum_n (s[n] s[n-1] + \sum_{k=1}^P a_k s[n-1] s[n-k]) = 0$$

for $k=2$, $\frac{\partial E}{\partial a_2} = \sum_n (s[n] s[n-2] + \sum_{k=1}^P a_k s[n-2] s[n-k]) = 0$

$$\frac{\partial E}{\partial a_k} = \sum_n \left(a[n] a[n-p] + \sum_{k=1}^p a_k a[n-p+k] \right) = 0$$

$$\Rightarrow \sum_n s[n] s[n-p] = - \sum_{k=1}^p a_k \sum_n s[n-k] s[n-p]$$

Autocorrelation method:

$$R(i) = \sum_{n=-\infty}^{\infty} s[n] s[n+i]$$

$$\text{Also } R(i) = R(-i)$$

$$R(i) = \sum_{n=-\infty}^{\infty} s[n] s[n-i]$$

$$R(i) = - \sum_{k=1}^p a_k \sum_{n=-\infty}^{\infty} s[n-k] s[n-i]$$

$$R(i) = -\sum_{k=1}^{r-1} a_k R(i-k)$$

where $i = 1, 2, \dots, p$

Total minimum error (E_p):

$$E = \sum_n (s[n] + \sum_{k=1}^p a_k s[n-k])^2$$

$$= \sum_{n=-\infty}^{\infty} \left\{ s^2[n] + \left(\sum_{k=1}^p a_k s[n-k] \right)^2 + 2s[n] \sum_{k=1}^p a_k s[n-k] \right\}$$

$$E = \underbrace{\sum_{n=-\infty}^{\infty} \delta^2[n]}_{\text{Term 1}} + \sum_{n=-\infty}^{\infty} \left(\sum_{k=1}^p a_k s[n-k] \right) \left(\sum_{i=1}^p a_i s[n-i] \right)$$

$$+ \underbrace{s[n] \sum_{k=1}^p a_k s[n-k]}_{\text{Term 2}}$$

$$\begin{aligned}
 E &= \text{Term 1} + \sum_{i=1}^p a_i \sum_{n=-\infty}^{\infty} \sum_{k=1}^p a_k s[n-k] s[n-i] + \\
 &\quad + 2 \sum_{n=-\infty}^{\infty} \sum_{k=1}^p a_k s[n] s[n-k] \\
 &= \sum_{n=-\infty}^{\infty} s^2[n] - \sum_{k=1}^p a_k \sum_{n=-\infty}^{\infty} s[n] s[n-k] + \\
 &\quad 2 \sum_{n=-\infty}^{\infty} \sum_{k=1}^p a_k s[n] s[n-k] \\
 &= \sum_{n=-\infty}^{\infty} s^2[n] + \sum_{n=-\infty}^{\infty} \sum_{k=1}^p a_k s[n] s[n-k] \\
 &\quad \downarrow \qquad \qquad \qquad \downarrow \\
 \therefore E_p &= R(0) + \sum_{k=1}^p a_k R(k)
 \end{aligned}$$

Levinson Durbin (LD) Recursive approach

$$R(0) + \sum_{k=1}^p a_k R(k) = E_p$$

$$R(i) + \sum_{k=1}^p a_k R(|i-k|) = 0 \quad ; \quad i=1, 2, \dots, p$$

$$\begin{array}{ccc}
 R^{(i-1)} & a^{(i-1)} & = e^{(i-1)} \\
 \downarrow & \downarrow & \downarrow \\
 \text{matrix} & \text{vector} & \text{vector}
 \end{array}$$

$$\begin{bmatrix} R(0) & R(1) & \dots & R(i-2) & R(i-1) \\ \vdots & \vdots & & \vdots & \vdots \\ R(i-1) & R(i-2) & & R(1) & R(0) \end{bmatrix} \begin{bmatrix} a_1^{(i-1)} \\ \vdots \\ a_{i-1}^{(i-1)} \end{bmatrix} = \begin{bmatrix} E_{i-1} \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

i^{th} iteration autocorrelation matrix:

$$\begin{array}{c} \left[\begin{array}{cccccc} R(0) & R(1) & \cdots & R(i-2) & R(i-1) & R(i) \\ R(1) & R(0) & \cdots & \cdots & \vdots & \vdots \\ \vdots & \vdots & & & \vdots & \vdots \\ R(i-1) & R(i-2) & \cdots & \cdots & R(0) & R(1) \\ R(i) & R(i-1) & \cdots & \cdots & R(1) & R(0) \end{array} \right] \begin{bmatrix} 1 \\ a_1^{(i-1)} \\ \vdots \\ a_{i-1}^{(i-1)} \\ 0 \end{bmatrix} = \begin{bmatrix} E_{i-1} \\ 0 \\ \vdots \\ 0 \\ r^{i-1} \end{bmatrix} \\ 1 \leftarrow \end{array}$$

$$\gamma^{(i-1)} = R(i) - \sum_{j=1}^{i-1} a_{i-j}^{(i-1)} R(j)$$

$$2 \leftarrow \begin{bmatrix} R(0) & R(1) & \cdots & R(i) \\ \vdots & \vdots & & \vdots \\ R(i) & R(i-1) & \cdots & R(0) \end{bmatrix} \begin{bmatrix} 0 \\ a_{i-1}^{(i-1)} \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} \gamma^{i-1} \\ 0 \\ \vdots \\ E_{i-1} \end{bmatrix}$$

$$R[i] \left\{ \begin{bmatrix} 1 \\ a_1^{(i-1)} \\ \vdots \\ a_{i-1}^{(i-1)} \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ a_{i-1}^{(i-1)} \\ \vdots \\ a_1^{(i-1)} \\ 1 \end{bmatrix} \right\} = \left\{ \begin{bmatrix} E_{i-1} \\ 0 \\ \vdots \\ \gamma^{i-1} \\ E_{i-1} \end{bmatrix} - k_i \begin{bmatrix} \gamma^{i-1} \\ 0 \\ \vdots \\ E_{i-1} \end{bmatrix} \right\}$$

$$E_i = E_{i-1} - k_i \gamma^{i-1}$$

$$a_i^{(i)} = a_1^{i-1} - k_i a_{i-1}^{(i-1)}$$

$$a_{i-1}^{(i)} = a_{i-1}^{(i-1)} - k_i a_1^{(i-1)}$$

$$a_i^{(i)} = 0 - k_i$$

Superscript = iteration

$$\text{Also } \gamma^{i-1} - k_i E_{i-1} = 0 \Rightarrow k_i = \frac{\gamma^{i-1}}{E_{i-1}}$$

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

LD Algorithm

$a_k^{(i)}$ \rightarrow kth LPC coefficient value at ith iteration

$E^{(i)}$ \rightarrow Residual error after ith iteration.

Steps: 1> Initially set $E^{(0)} = R(0)$; $i=1$

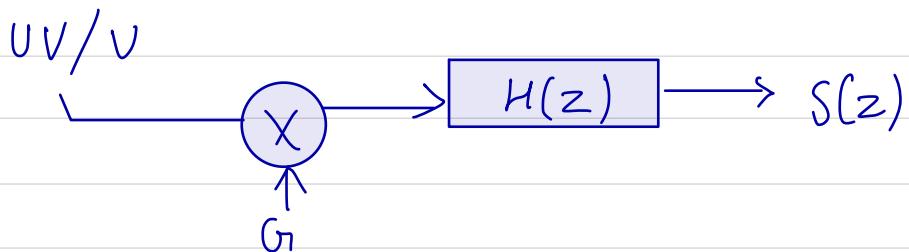
2> Calculate $k_i = \frac{[R(i) - \sum_{j=1}^{i-1} a_j^{(i-1)} R(i-j)]}{E^{(i-1)}}$

3> Set $a_i^{(i)} = k_i$ & $a_j^{(i)} = a_j^{(i-1)} - k_i a_{i-j}^{(i-1)}$; $1 \leq j \leq i$

4> Calculate

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

5> Repeat steps ②, ③, ④, until $i = p$
where p is linear prediction order.



$$H(z) = \frac{1}{1 + \sum_{k=1}^p a_k z^{-k}} \quad \hookrightarrow \text{LP coefficient}$$

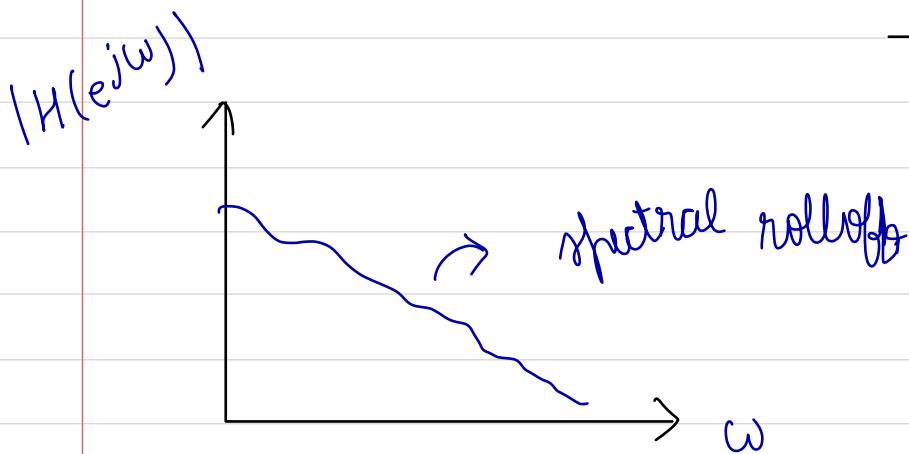
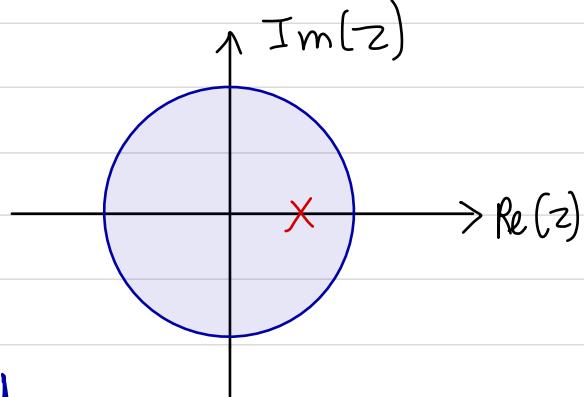
$$\text{Freq. domain: } z = e^{j\omega}$$

$$|H(e^{j\omega})| = \left| \frac{1}{1 + \sum_{k=1}^p a_k e^{-j\omega k}} \right|$$

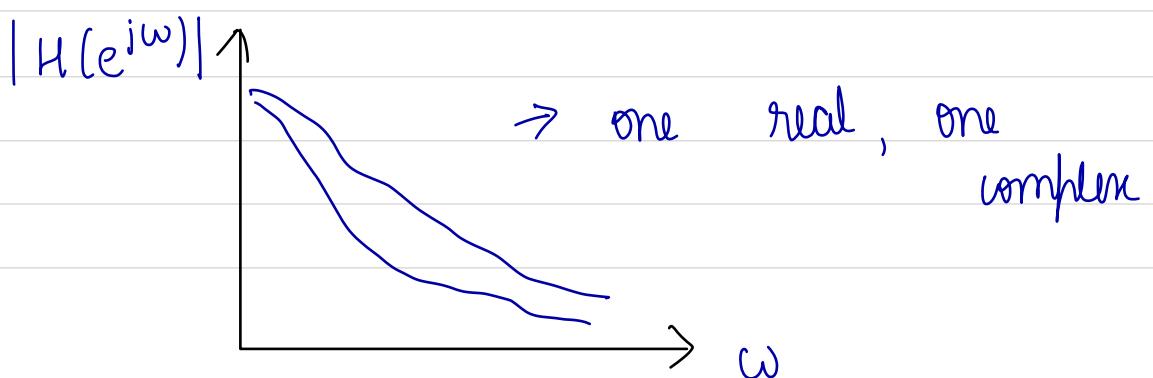
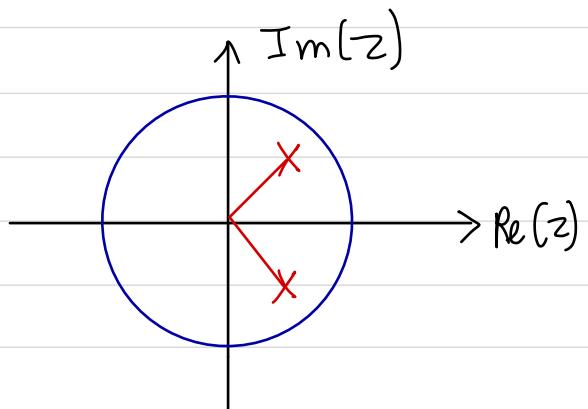
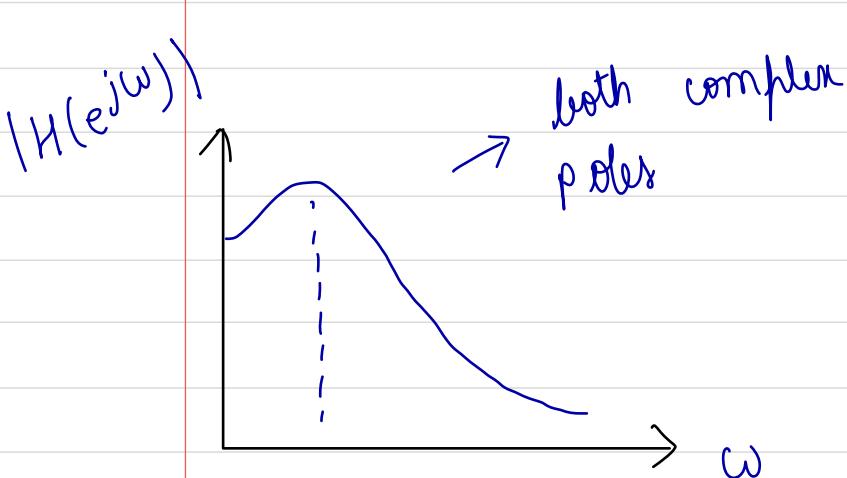
Denominator : $1 + a_1 z^{-1} + a_2 z^{-2} + \dots + a_p z^{-p}$

1) Let $p = 1$

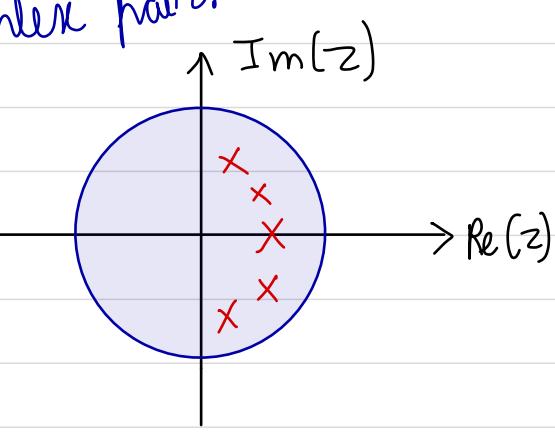
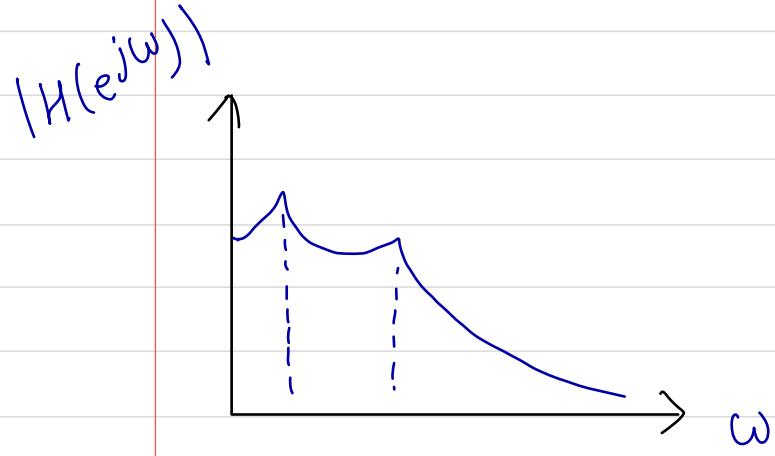
$$H(z) = \left| \frac{1}{1 + a_1 z^{-1}} \right|$$



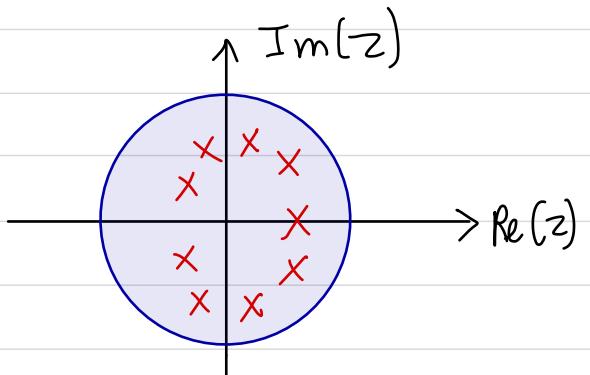
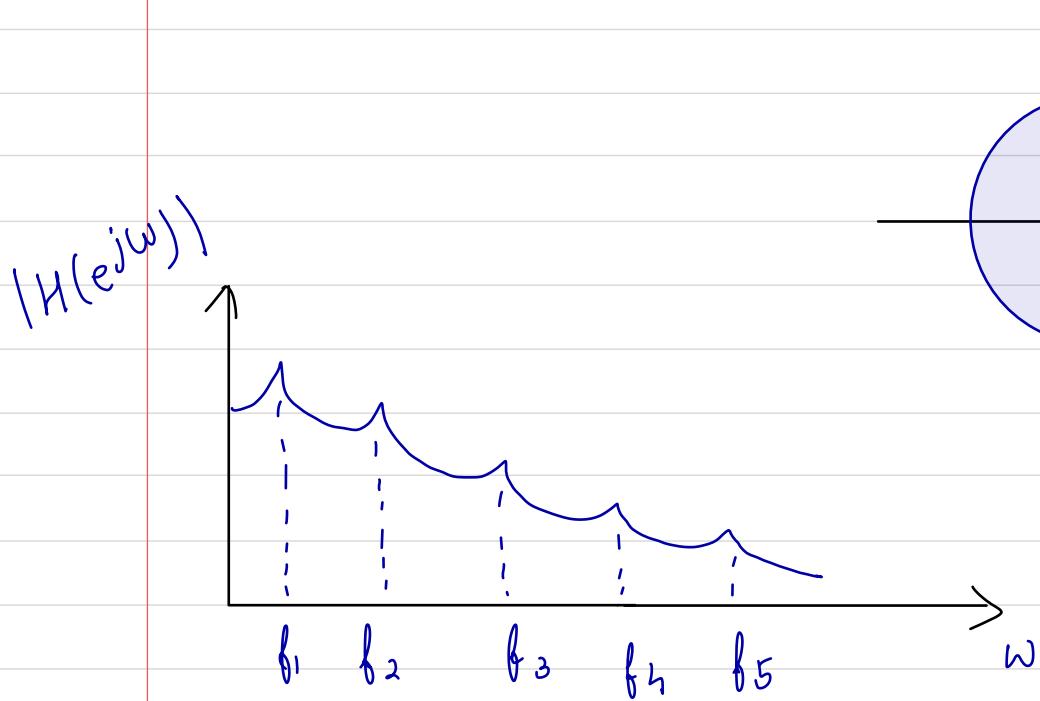
2) Let $p = 2 \rightarrow$ real & complex

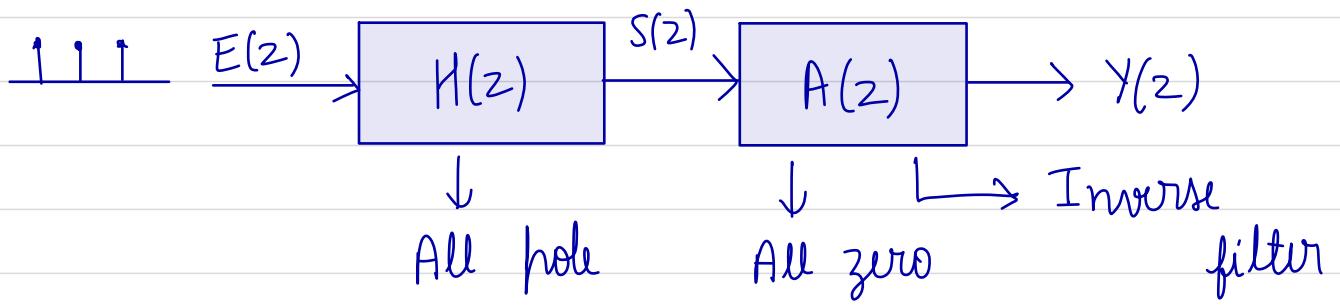


3) let $p = 5 \rightarrow 2$ complex pairs



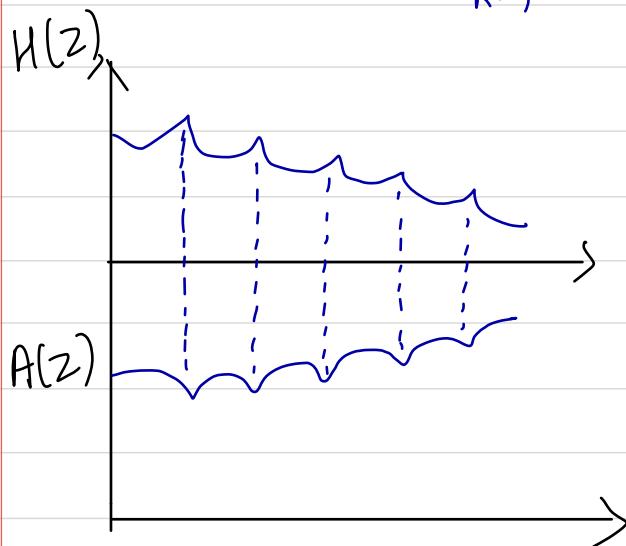
4) let $p = 11 \rightarrow 5$ complex pairs





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

$$A(z) = 1 + \sum_{k=1}^P a_k z^{-k}$$



LP residual

$$A(z) = \frac{Y(z)}{S(z)}$$

$$Y(z) = \left[1 + \sum_{k=1}^P a_k z^{-k} \right] s(z)$$

$$\begin{aligned} \text{MMSSE} &= s(n) - \tilde{s}(n) \\ &= s(n) + \sum_{k=1}^P a_k s(n-k) \end{aligned}$$

$$y(n) = s(n) + \sum_{k=1}^P a_k s(n-k)$$

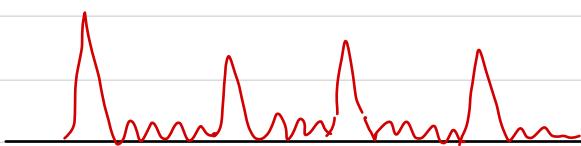


\bar{DEGGr} = Differenced electrogastrogram

Fricative: Excitation

LP residual

$$\Lambda(n) \rightarrow (N \times 1)$$



Assignment:
 1> Compute $LPC_n \Rightarrow LD$ method
 i/p: speech, LP order (p)
 o/p: LPC (2D matrix)

- 2> Plot LP spectrum
- 3> Compute LP residual & find out F_0

Gaussian mixture models

- It estimates the data distribution into multiple gaussian models.
- It is a generative model.

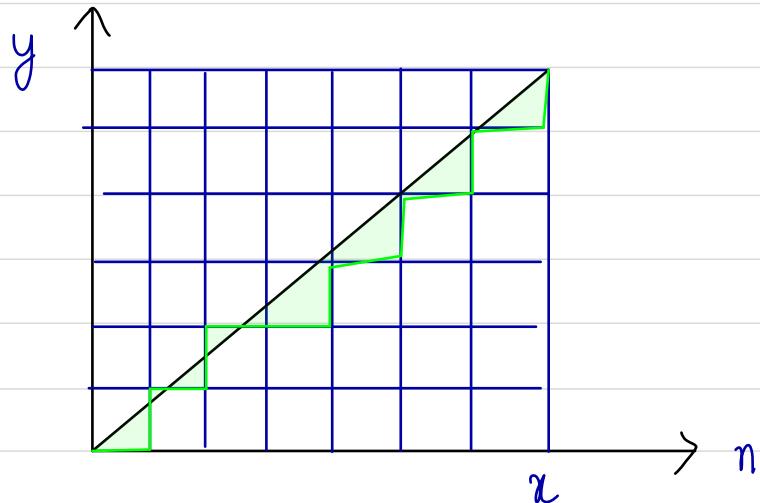
Different types of models:

- 1> K - Means
- 2> GMM - Gaussian mixture models
- 3> VQ - Vector Quantisation

Use cases:

- 1> Vowel v/s consonant classification
- 2> Keyword spotting

Dynamic time warping:



$$w_1 \rightarrow 1, 2, \dots T_x = 90 \text{ features}$$

$$w_2 \rightarrow 1, 2, \dots T_y = 130 \text{ features}$$

Need to map them properly. { since they have different number of features }

$\phi = \{\phi_x(k), \phi_y(k)\}$
 ↳ Non linear warping function

$$\phi(x, y) = \sum_{k=1}^T d(\phi_x(k), \phi_y(k)) \cdot \frac{m(k)}{m(\phi)}$$

$m(k)$ = path weighting coefficients
 $m(\phi)$ = path normalization coefficients

$$d(x, y) = \min_{\phi} \{ d(\phi(x, y)) \}$$

Asynchronous Mover approach:

Bellmann optimality:

$\phi(i, j)$ - cost involved from city i to city j

$\phi(i, l)$ - cost involved from city i to intermediate city l

$\varepsilon_i(i, j)$ - cost of direct travel from i to j

$$\phi(i, j) = \min_l [\phi(i, l) + \varepsilon_i(l, j)]$$

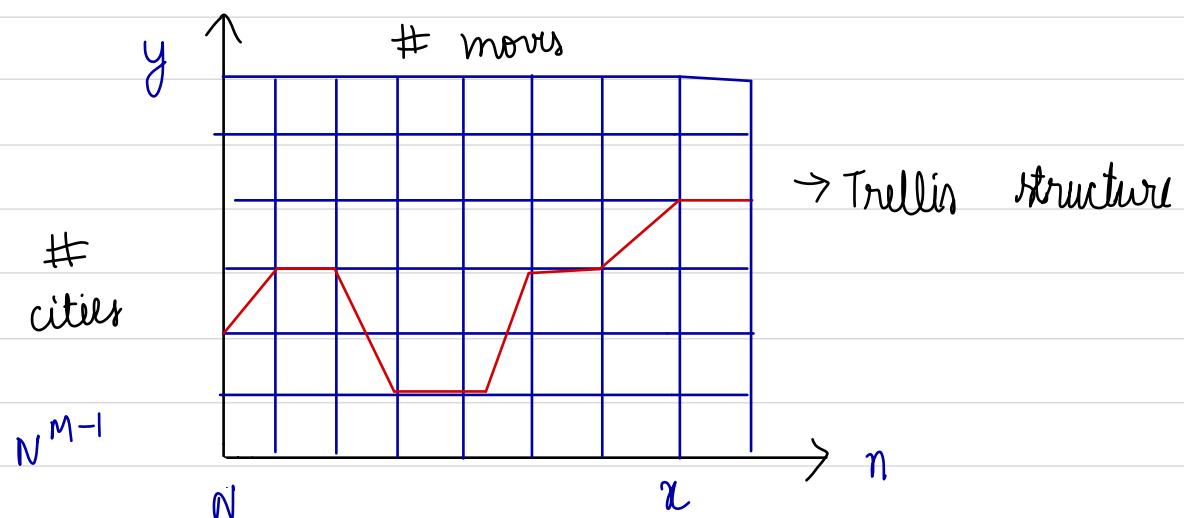
For direct travel $\phi(i, j) = \varepsilon_i(i, j)$
(No halt) where $i, j = 1, 2, \dots, N$

For one halt:

$$\phi_2(i, j) = \min_{1 \leq l \leq N} [\phi_1(i, l) + \varepsilon_i(l, j)]$$

$$\phi_3(i, j) = \min_{1 \leq l \leq N} [\phi_2(i, l) + \varepsilon_i(l, j)]$$

Synchronous move approach



Viterbi Algorithm:

Recursive initialisation:

$$\phi(i, n) = \xi(i, n)$$

For path tracing: $\xi'(n) = i$, $n = 1, 2, \dots, N$

Recursive: $\phi_{m+1}(i, n) = \min_{1 \leq l \leq N} [\phi_m(i, l) + \xi(l, n)]$

$$\xi'(n) = \arg [\phi_{m+1}(i, n)]$$

To do path tracing from $\xi'_{m+1}(n)$:

Constraints of template matching for Viterbi

1> Endpoint constraint:

- [ra] [re] [ga] [ma] -
- [ma] [ga] [re] [ra] -

In endpoint constraint: $\phi_x(1) = 1 = \phi_y(1)$
for starting point

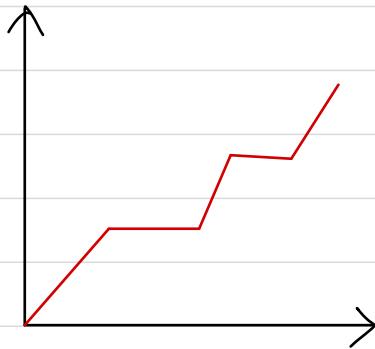
For endpoint:

$$\phi_x(T) = T_x$$

$$\phi_y(T) = T_y$$

where T_x = number of time frames in X-direction
 T_y = number of time frames in Y-direction

2> Monotonicity constraint



monotonic



non monotonic

$$\text{So } \phi_x(k+1) \geq \phi_x(k)$$

$$\phi_y(k+1) \geq \phi_y(k)$$

3> Local continuity constraint: (LCC)

$$\phi_x(k+1) - \phi_x(k) \leq 1$$

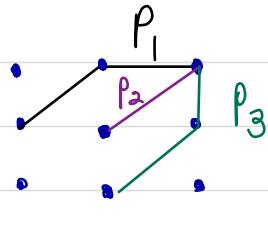
$$\phi_y(k+1) - \phi_y(k) \leq 1$$

$$d_\phi(x, y) = \sum_{k=1}^T d(\phi_x(k), \phi_y(k)) \frac{m(k)}{M_\phi}$$

Normalization

Path weighing
constant

Type a: $m(k) = \min [\phi_x(k) - \phi_x(k-1), \phi_y(k) - \phi_y(k-1)]$



$$P_1 : 2+1 = 3$$

$$P_2 : 1+0 = 1$$

$$P_3 : 1+2 = 3$$

P_2 is chosen as it is minimum.

Path normalizing factor :

$$m_k = \phi_x(k) - \phi_x(k-1)$$

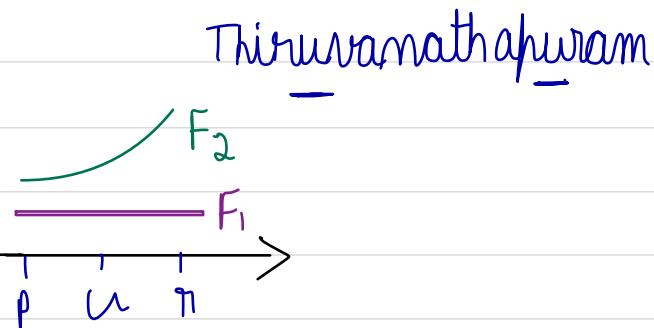
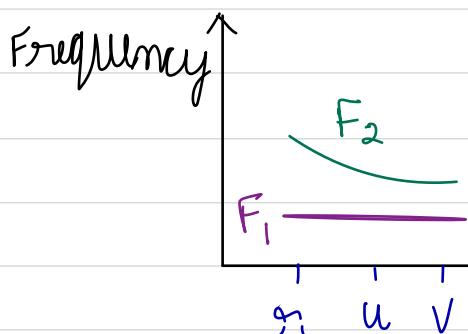
$$M_\phi = \sum_{k=1}^T m(k)$$

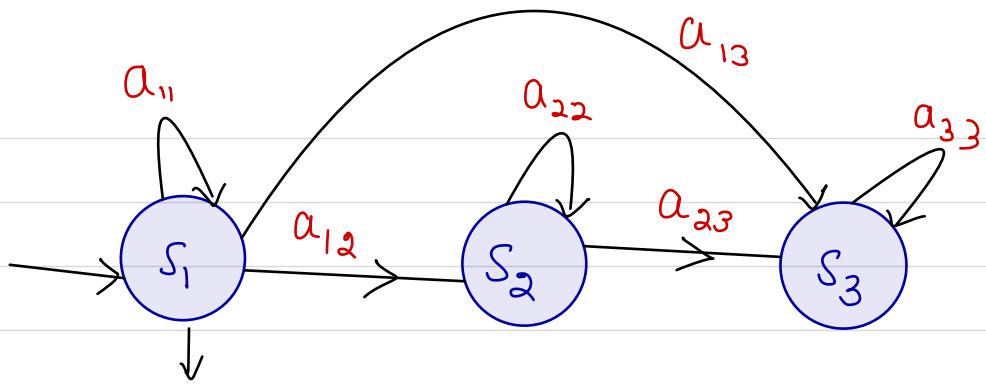
Hidden Markov models

- In DTW $N = \text{large}$
- In HMM $M \ll N$

- In DTW parameter is mean
- In HMM parameters are mean & covariance

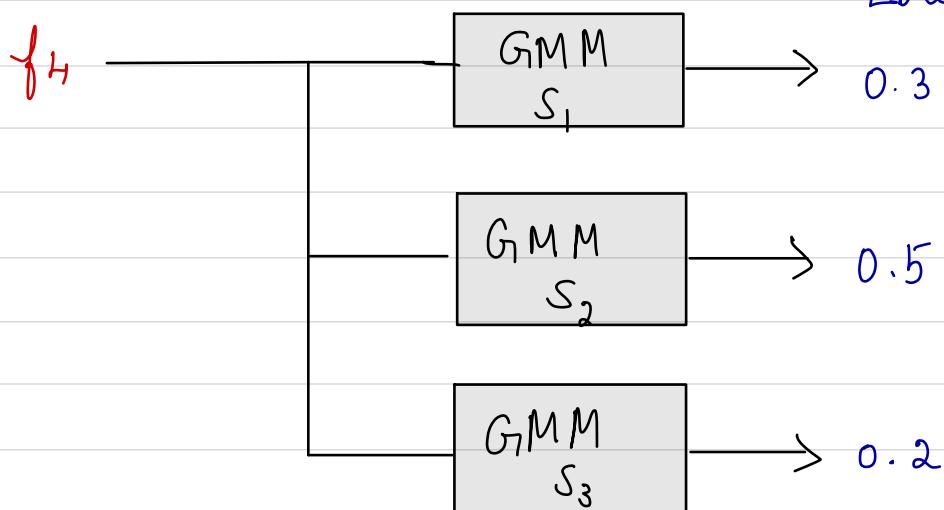
- Goal of DTW is to minimize distance.
- Goal of HMM is to minimize likelihood as it is a probabilistic model.



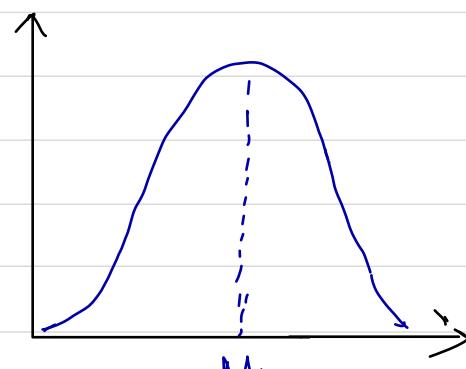


State distribution learning

Likelihood



Gaussian mixture models



1D Gaussian

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

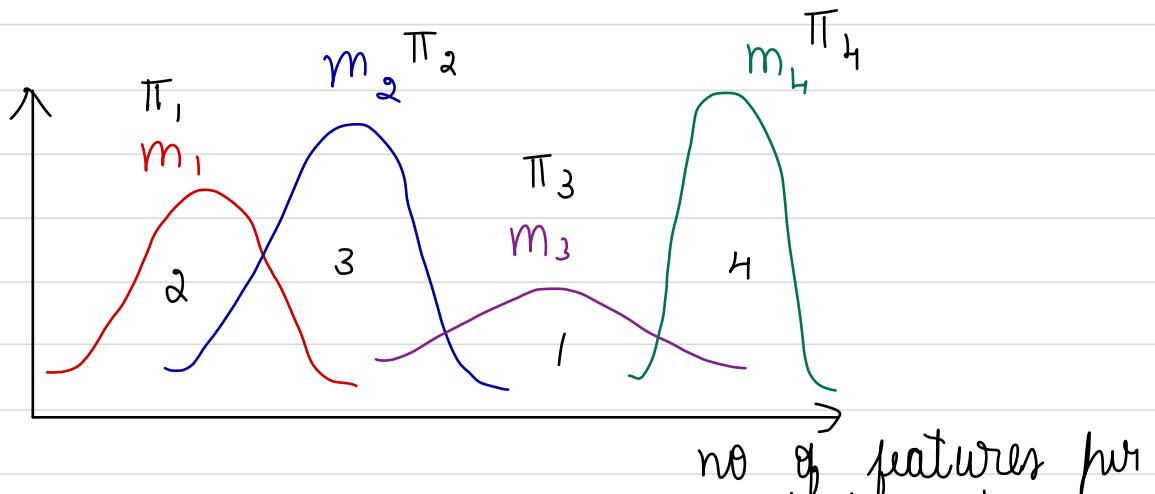
Multivariate gaussian distribution:

Eg: 39 D MFCC:

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

$$\vec{\mu} = \begin{bmatrix} \mu_0 \\ \mu_1 \\ \vdots \\ \mu_{38} \end{bmatrix}_{39 \times 1}$$

$$\Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \dots & \sigma_{1N}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 & \dots & \sigma_{2N}^2 \\ \vdots & \vdots & & \vdots \\ \sigma_{N1}^2 & \sigma_{N2}^2 & \dots & \sigma_{NN}^2 \end{bmatrix}$$



$$\pi_1 = 0.2 \quad \pi_2 = 0.3$$

π_i = weights

$$P(x | A) = \sum_{i=1}^m \pi_i P_i(x | \mu_i, \Sigma_i)$$

Probability distribution of i^{th} mixture

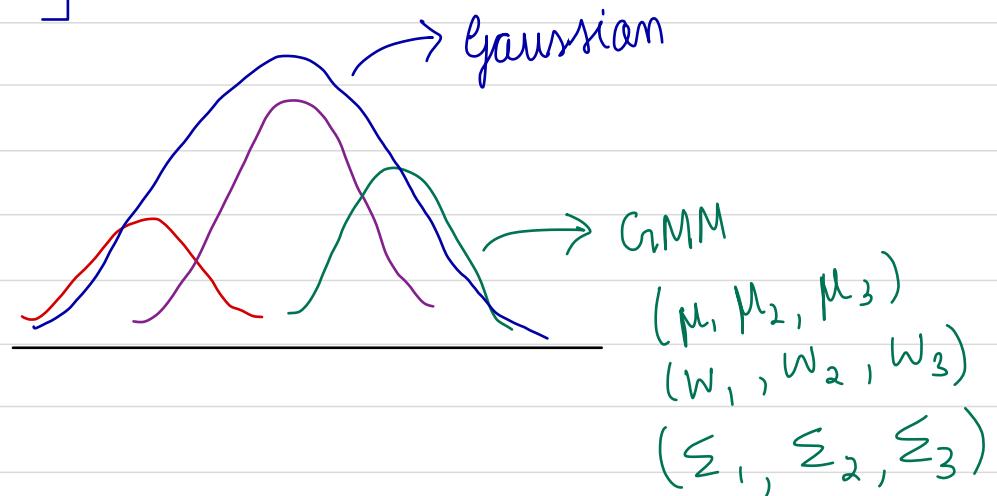
weight of i^{th} component or mixture

Expectation maximization (EM) Algorithm :

→ Maximize the probability of generating feature vector from the model.

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_T \end{bmatrix}$$

$$E(x) = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_T \end{bmatrix}$$



Procedure:

$$X = \{x_1, x_2, \dots, x_T\} \quad \text{let } T=1000$$

$$(39 \times 1000) \quad p(x|\bar{\lambda}) \geq p(x|\lambda)$$

↓ ↗

model in the model in previous iteration
current situation

$$p(i | x_t, \lambda) = \frac{w_i p(x_t | \mu_i, \Sigma_i)}{\sum_{k=1}^M w_k p(x_t | \mu_k, \Sigma_k)}$$

$i = 1, \dots, M$

i^{th} Gaussian mixture
Current feature vector
Model GMM

Also $\sum_{i=1}^M p(i | x_t, \lambda) = 1$

Procedure:

$$\bar{w}_i = \frac{1}{T} \sum_{t=1}^T p(i | x_t, \lambda)$$

$$\bar{\mu}_i = E(x) = \frac{\sum_{t=1}^T x_t p(i | x_t, \lambda)}{\sum_{t=1}^T p(i | x_t, \lambda)}$$

w_i = weights

$$\begin{aligned} \sigma_i^2 &= E(x^2) - (E(x))^2 \\ &= \frac{\sum_{t=1}^T x_t^2 p(i | x_t, \lambda)}{\sum_{t=1}^T p(i | x_t, \lambda)} - (\bar{\mu}_i)^2 \end{aligned}$$

$$\bar{\Sigma}_i = \begin{bmatrix} \sigma_{11}^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma_{22}^2 & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \sigma_{NN}^2 \\ 0 & \dots & \dots & \ddots & \sigma_{NN}^2 \end{bmatrix}$$

$$\lambda = \{w_i, \mu_i, \Sigma_i\}$$

$$\bar{\lambda} = \{\bar{w}_i, \bar{\mu}_i, \bar{\Sigma}_i\}$$

$$p(x|\bar{\lambda}) - p(x|\lambda) \leq \text{Threshold} \quad (\text{eg: } 0.0001)$$

↓
Termination condition

Digit recognition task:

Vocabulary : [0, 1, ..., 9]

↓
500 examples / class

→ Speaker variability / gender

→ speaking mode

→ channel variability

Training : $500 \times 10 = 5000$ examples

MFCC : $[39 \times 1]$

let $FV = 7000$

$$X = \{x_1, x_2, \dots, x_T\}$$

$\downarrow \quad \downarrow \quad \downarrow$

$$c_0 \quad c_0 \quad c_9$$

$[39 \times 7000]$

GMM : $\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_g$

\downarrow

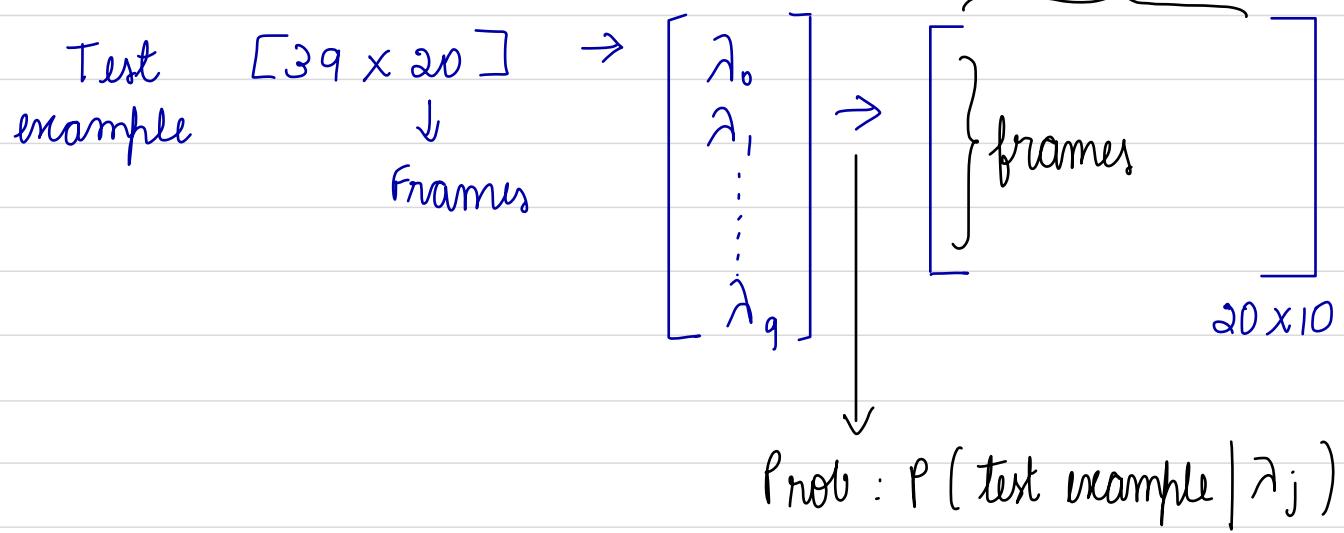
$$(w_0, \mu_0, \Sigma_0)$$

\downarrow

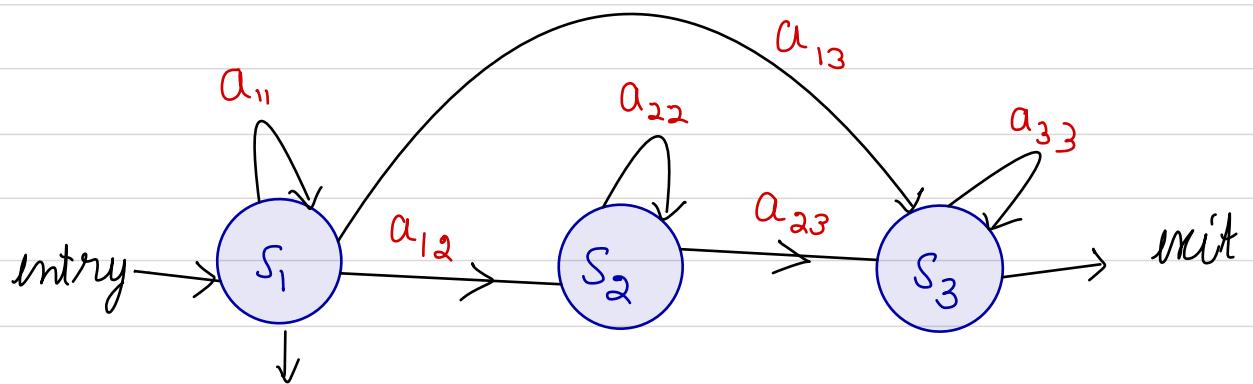
$$(w_1, \mu_1, \Sigma_1)$$

$fV \in C_0 \Rightarrow K_{\text{matrix}} \Rightarrow \begin{matrix} \mu_0^i \\ \Sigma_0^i \\ w_0^i \end{matrix} \Rightarrow EM \Rightarrow \begin{matrix} \mu_0 \\ \Sigma_0 \\ w_0 \end{matrix}$

Testing $\rightarrow 100 \times 10 = 1000$ examples classes



Transition probability



State distribution learning

$$P(q_{t+j} = j \mid q_{t+1}, q_{t+2}, \dots)$$

First order hidden markov model

$P(q_{t+1} = j \mid q_{t-1} = i)$ is the objective.

Observable markov model (OMM);

Here we can see the past sequence clearly, like in the case of tossing a coin continuously.

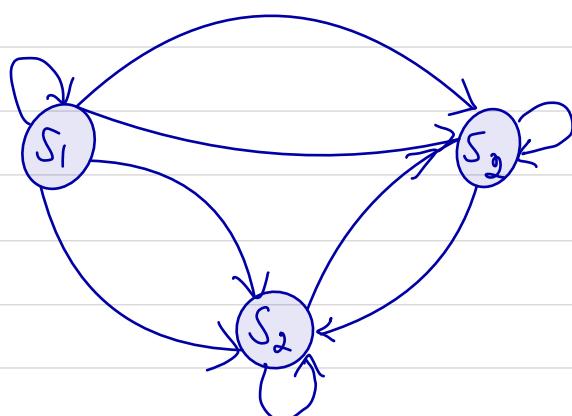
Eg: Take a biased coin with $P(H) = 0.7$ & $P(T) = 0.3$:

observation: T, T, H, H, H, T, H, ... ↓
↓
| 1000

Now consider 3 biased coins.

C_0	C_1	C_2
$P(H) = 0.3$	$P(H) = 0.2$	$P(H) = 0.6$
$P(T) = 0.7$	$P(T) = 0.8$	$P(T) = 0.4$

Observation : T, T, H, T .. T
 Hidden : C₂, C₁, C₁, C₂ .. C₀



1> a_{ij} = state transition probability

final
initial

$$h(q_{t+1} = j \mid q_t = i)$$

2> $N \neq M$
 $N = \text{states}$
 $M = \text{observation}$

3> State Transition Matrix:

$$A = [a_{ij}]_{N \times N}$$

$$= \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1N} \\ a_{21} & a_{22} & \ddots & \dots & a_{2N} \\ \vdots & \vdots & & & \\ a_{N1} & a_{N2} & \ddots & \dots & a_{NN} \end{bmatrix}$$

$$\sum_{j=1}^N a_{ij} = 1$$

(row sum = 1)

4> States: modelled as Histogram
 GMM

M distant Observation $\leftarrow \frac{1}{T}$

$$V = \{v_1, v_2, \dots, v_M\} \Rightarrow \text{observation symbol}$$

$$\text{observation seq, } O = \{O_1, O_2, \dots, O_T\}$$

$$T = 1000$$

observation symbol probability:

$$b_j(k) = p(\theta_t = v_k \mid q_t = j)$$

probability of producing observation v_k at time t given that j th state at t .

$$\sum_{k=1}^M b_j(k) = 1$$

$$B = [b_j(k)]_{N \times M}$$

$$\begin{array}{l} j: I \rightarrow N \\ k: I \rightarrow M \end{array}$$

$$\sum_{k=1}^M b_j(k) = 1$$

→ stochastic constraint

Initial State Probability:

$$\pi_i = p(q_0 = i) \quad i: I \rightarrow N$$

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

$$\lambda = \left\{ \frac{A}{N}, \frac{B}{M}, \pi \right\} \quad (\text{parameters})$$

Testing Procedure
Optimal State Seq finding
Training Procedure

Problem 1:

Testing Procedure : The process of efficiently computing $p(O|\lambda)$, given the observation seq 'o' & model ' λ '.

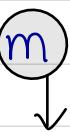
Problem 2:

Optimal state seq finding: This involves finding the optimal state seq, given observation 'o' & model λ

Problem 3: Training Procedure :

Given observation seq 'o' & initial model ' λ '
How to re-estimate model parameters.
 $p(O|\bar{\lambda}) \geq p(O|\lambda)$

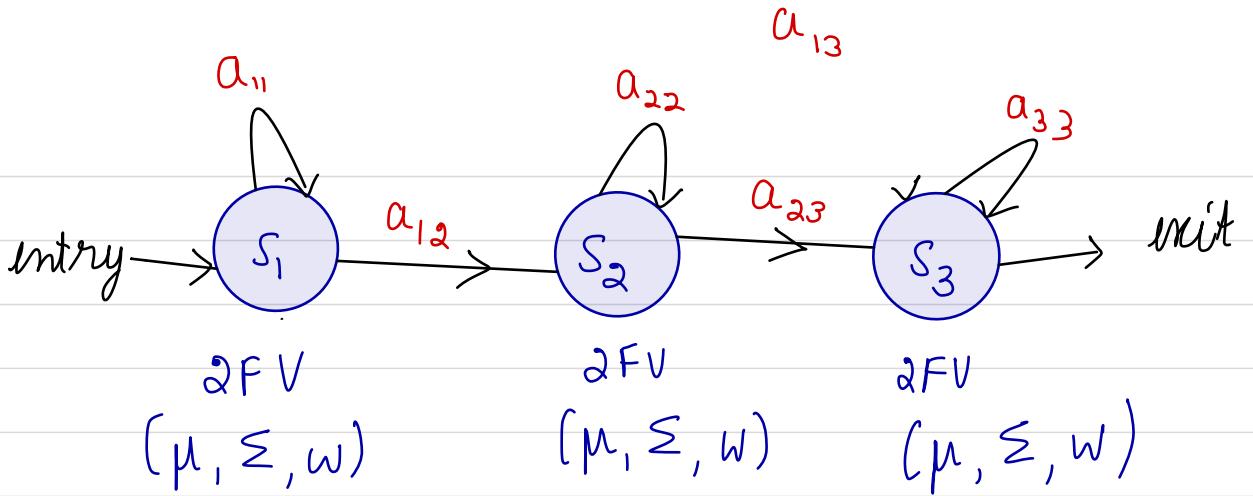
Consider "Mera Bhaarat Mahaan"

silence -  e n a b h a a n a t m a h a a n - sil

3 state HMM

Let 1 sec = 98 Frames (25 msec, 10 msec)

\Rightarrow This will be $[39 \times 98] \rightarrow$ 6 Feature vector/
Phoneme
 $\rightarrow 2(FV / \text{state}) / \text{Phoneme}$



$$a_{ij} = \begin{cases} 0.5 & \text{if } j=1 \text{ or } j=i+1 \\ 0 & \text{otherwise} \end{cases}$$

$$O = \{o_1, \dots, o_T\}$$

$$q = \{q_1, \dots, q_T\}$$

$P(O|\lambda) \rightarrow$ Observation -

Solution to problem: Efficiently compute $P(O|\lambda)$

Direct Method / Brute force:

$P(O|\lambda) \rightarrow$ Probability of ' O ' being generated by λ .

$$O = \{o_1, \dots, o_T\}$$

$$q = \{q_1, \dots, q_T\}$$

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

$$= b_{q_1}(o_1) \cdot b_{q_2}(o_2) \cdot \dots \cdot b_{q_T}(o_T)$$

$$= \prod_{t=1}^T b_{q_t}(o_t)$$

$$P(O|\lambda) = \sum_{\text{all } q \rightarrow N^T} P(O|q, \lambda) \cdot P(q|\lambda)$$

$$P(q|\lambda) = P(q_1, q_2, \dots, q_T | \lambda)$$

$$= [\pi_{q_1} \ a_{q_1, q_2} \ \dots \ a_{q_{T-1}, q_T} \ \dots \ a_{q_{T-1}, q_T}]$$

$\downarrow \quad \downarrow$
 $s_1 \quad s_2$

$$P(O|\lambda) = \sum_{\forall q \rightarrow N^T} \pi_{q_1} b_{q_1}(O_1) a_{q_1, q_2} \cdot b_{q_2}(O_2) \cdot \dots \cdot a_{q_{T-1}, q_T} b_{q_T}(O_T)$$

Forward procedure for $P(O|\lambda)$

Forward variable:

$$\alpha_t(i) = p(O_1, O_2, \dots, O_t, q_t = i | \lambda)$$

Initialization at $t=1$

$$\alpha_1(i) = p(O, q_1 = i | \lambda) = \pi_i b_i(O_i)$$

\hookrightarrow joint Probability

Recursion:

$$\alpha_t(i) = p(O_1, O_2, \dots, O_t, q_t = i | \lambda)$$

$$\alpha_{t+1}(j) = p(O_1, O_2, \dots, O_t, O_{t+1}, q_{t+1} = j | \lambda)$$

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

$1 \leq t \leq T-1$
 $1 \leq j \leq N$

Termination:

$$\alpha_T(j) = \sum_{i=1}^N \alpha_{T-1}(i) a_{ij} b_j(O_T)$$

$$p(O|\lambda) = \sum_{j=1}^n \alpha_T(j) p(O_1, \dots, O_T, q_T | \lambda)$$

Method 3: Backward procedure:

Backward Variable:

$$B_t(i) = P(O_{t+1}, O_{t+2}, \dots, O_T | q_{t+1} = i, \lambda)$$

Initialization: $B_T(i) = P(O_T | q_T = i, \lambda)$

Recursion: $B_t(i)$

$$B_{t+1}(j) = P(O_{t+2}, \dots, O_T | q_{t+1} = j, \lambda)$$

$$B_{t+1}(j) = \sum_{i=1}^N a_{ij} B_t(i) b_j(O_{t+1}) \quad 1 \leq t \leq T-1$$

$$1 \leq i \leq N$$

Termination: $B_1(i) = \sum_{j=1}^N a_{ij} B_2(j) b_j(O_2)$

$$p(O|\lambda) = \sum_{i=1}^N B_1(i) \pi_i b_i(O_i)$$

Problem 2: Optimal state sequence estimation
given $O, \lambda \rightarrow$ find out q_1, q_2, \dots, q_T

$$\delta_T(i) = \max_{q_1, q_2, \dots, q_T} (q_1, q_2, \dots, q_T = i, O_1, O_2, \dots, O_T | \lambda)$$

Path tracing : $\Psi_t(j) = \arg \delta_{t-1}(j)$

Initialization :

$$\begin{aligned} s_i(i) &= p(q_1 = i, o_1 | \lambda) = \pi_i b_i(o_1) \\ \Psi_0(i) &= 0 \end{aligned}$$

Recursion:

$$s_{t-1}(i) = \max_p p(q_1, q_2, \dots, q_{t-1} = i, o_1, o_2, \dots, o_{t-1} | \lambda)$$

$$\delta_t(j) = \max_p p(q_1, q_2, \dots, q_t = j, o_1, \dots, o_t | \lambda)$$

$$s_t(j) = \max_i \delta_{t-1}(i) a_{ij} b_j(o_t)$$

$$\Psi_t(j) = \arg s_t(i)$$

Termination: $s_T(j) = \max_i s_{T-1}(i) a_{ij} b_j(o_T)$
 $\Psi_T(j) = \arg s_T(j)$

$$i = \max_j s_T(j) \quad q^* = \Psi_{T+1}(q_{T+1}^*)$$

Problem 3: Training of HMM

$$\lambda = \{\mathbf{A}, \mathbf{B}, \boldsymbol{\pi}\}$$

$$\begin{aligned} Y_t(i) &= p(q_t = i | \lambda) \\ &= p(q_t = i | \lambda) - \frac{p(o | \lambda)}{p(o | \lambda)} \\ &= \frac{p(q_t = i, o | \lambda)}{\sum_{j=1}^N p(q_t = j, o | \lambda)} \end{aligned}$$

$$\gamma_t(i) = \frac{P(O_0, \dots O_t, q_t = i | \lambda) P(O_{t+1}, \dots O_T, q_t = i | \lambda)}{\sum_{j=1}^N P(q_t = j, o | \lambda)}$$

$$\gamma_t(\lambda) = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)}$$

$$\xi_t(i, j) = P(q_t = i, q_{t+1} = j | \lambda) \left(\frac{P(O | \lambda)}{P(O' | \lambda)} \right)$$

Vowel identification
Speaker recognition \rightarrow HMM

$$\begin{aligned} \xi_t(i, j) &= \frac{P(q_t = i, q_{t+1} = j, o | \lambda)}{\sum_{i=1}^N \sum_{j=1}^N P(q_t = i, q_{t+1} = j, o | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)} \end{aligned}$$

$$\bar{\pi}_i = P(q_1 = i) = \gamma_i(i)$$

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

$$P(O | \bar{\lambda}) \geq P(O | \lambda)$$

$$\bar{b}_j(k) = \frac{\sum_{t=1}^T \gamma_t(j) \delta(O_t, v_k)}{\sum_{t=1}^T \gamma_t(j)}$$