② Zero-crossing Rate: (ZCR)

    Unvoiced has high
    Voiced has low

③ Presence of pitch:

    Voiced has high pitch
    Unvoiced has low pitch. (constriction).

④ Auto-correlation.

    Voiced has high auto-correlation coefficient
    Unvoiced has low auto-correlation coefficient

⑤ Bond Ratio: $\dfrac{\text{Energy upto 1KHz}}{\text{Total energy}}$

    High for voiced
    Low for unvoiced.

⑥ Linear prediction error:

$$\hat{S}n[k] = \Sigma\, a_k s[n{+}k]$$

$$\hat{s}[n] = \Sigma\, a_k s[n{+}k]$$

$$\boxed{error = s[n] - \Sigma\, a_k s[n{+}k]}$$

→ more for unvoiced  }→ Explanation:
   less for voiced        autocorrelation
                           coefficient.

Normalized LP error analysis:

⑦ First LP coefficient:

    Consider $a_1$,

    $a_1 >$ in voiced than unvoiced because voiced has high autocorrelation
    coefficient than in unvoiced.

Pitch modification
Time-scale modification $\}$ Trade off

modify time-scale, pitch changes
modify pitch, time scale changes.

Voiced and unvoiced speech detection:

three types of regions:

Voiced — all vowels
unvoiced
silence. — noise

Voiced region has better features.

Speech coding (AMR) — adaptive multi-rate coder.                    Pitch detection

64 Kbps : 4.4 Kbps – 12.2 Kbps
(silence)   (Voiced)

frontend block to speaker recognition.

Voiced vs unvoiced detection algorithms:

① Energy Threshold:

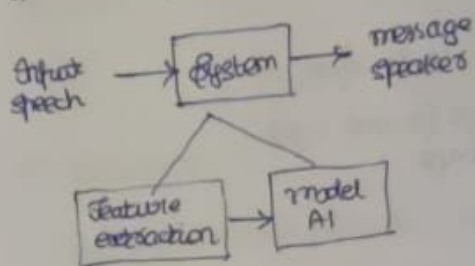Voiced ↑          Drawback: Noise also have high energy and will be detected
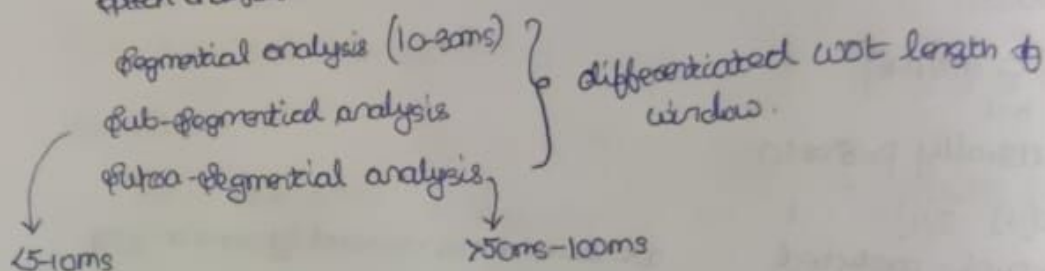Unvoiced ↓                    as unvoiced/voiced signal.
                              Not suitable method.

LP spectrum — Envelope of the DFT
↓
outputs envelope.

Speech Analysis:

Input speech → System → message speaker

Feature extraction → model AI

Speech analysis:

Segmental analysis (10-30ms)
Sub-segmental analysis      } differentiated wrt length of window.
Supra-segmental analysis

<5-10ms                                    >50ms-100ms

Segmental is most popular because the duration can cover around 2 to 3 pitch cycles

20ms = fixed frame size
2 to 3 pitch cycles.

Segmental analysis — VTS, Excitations
Sub Segmental analysis — Excitations
Supra Segmental analysis — Prosody

Pitch will not be there when the speaker is unvoiced (not talking
Pitch can vary within a person.

Short-time fourier transform works good if signal in my window is stationary.

Frame size??

FT works if whole signal is stationary. But speech is dynamic.

Spectrogram —— Implements STFT

STFT! for a signal STFT is not unique,
depends on window size (Frame size).
window shape

Usually, good frame size $\sim 20ms$.

LTI model of speech:

Linear prediction:

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

usually $p = 8$ or $10$.

$$e[n] = s[n] - \hat{s}[n]$$
actual - predicted.

$a_k$ values can be obtained by minimizing error $e[n]$.

$$e[n] \quad \frac{d}{da_k}(e^2[n]) = 0$$

Z Transform
$$\Rightarrow s[n] - \sum_{k=1}^{P} a_k s[n+k]$$

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

$$\frac{E(z)}{S[z]} = 1 - \sum_{k=1}^{P} a_k z^{-k}$$

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

→ Transfer Function of vocal track system.

Speech of/
Vocal
excitations
$e[n]/e[n]$

unvoiced
excitations

→ | Vocal Track System | → $S(t)/s[n]$

↑
$h[n]$ (LTI)
$H(z)$ (LTI)

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

$H(z) = \dfrac{S(z)}{e(z)}$ Synthesis equation

$$\Rightarrow \frac{1}{1 - \sum_{k=1}^{P} a_k z^{-k}} \quad \text{(all pole system)}$$

$$\frac{1}{H(z)} = \frac{e(z)}{S(z)} = 1 - \sum_{k=1}^{P} a_k z^{-k}$$
(analysis equation)
(all zero system)

Semivowels: Vocal track is not completely closed nor completely open.

ಐ ಒ ಎ ಔ
ಯ ರ ಲ v

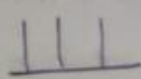Fricatives: Friction in the excitation part

h sh shh h s

Co-articulation:

Before one sound production ends, other sound production starts.

Consider (ಏಳಿ, ಅಳು) —— Co-articulation.

Both are different

Formants:

air → vocal folds → |||| → VTA → Resonate frequency
            vibrations

These frequencies are more or less same for a phone by every person.

Frequencies of a syllable are approximately same produced by any person.
These are called formants.

Lecture: 4

Recap: Vocal track system,
       exhitationes
       most important things of speech production.

Vowels = 150-300ms, consonants duration = <120ms
Consider a signal with vowels a, e, i
and another signal with vowels e, i, a
How to identify the vowels

Fourier transform? Not useful because FT spectrums of both the signals
                   are more or less the same.

Answer: STFT
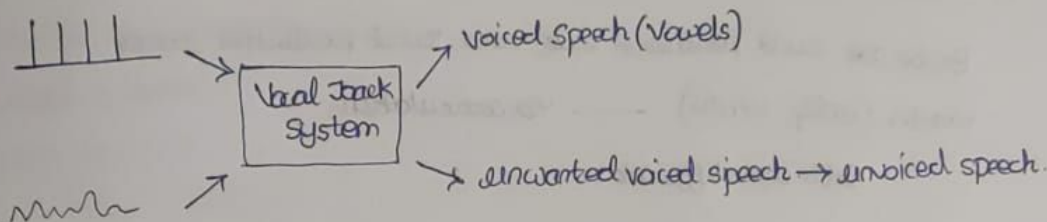        Short time Fourier transform.

manner of articulation ——— excitation characteristics

Place of articulation ——— system (vocal track) characteristics.

Generally,
  vowels have high strengths and long duration
  consonants have low strength and low duration



→ voiced speech (vowels)

Vocal Track System

→ unwanted voiced speech → unvoiced speech.

For all vowels, vocal track system must be opened.
All vowels are voiced and have no constriction.

for consonants:

PoA: Place where you put a constriction in producing a consonant in vocal track system.

| Place of articulation | Unvoiced unaspirated | Unvoiced aspirated | Voiced unaspirated | Voiced aspirated | Nasals | Semivowel |
|---|---|---|---|---|---|---|
| Velar | k | Kh | g | gh | Kn | |
| Palatel | ch | chh | j | jh | chn | y |
| Alvedar | T | Th | D | Dh | Tn | r |
| Dental | t | th | d | dh | n | l |
| Bilabial | P | ph | b | bh | m | v |

* Phone is the basic sound production unit

Syllable is the combination of phones.

90% Indian syllables are consonant vowel combination. (kg tha pa)

## Linear and Time invariant systems:

$x(t)$ —[ ]— $y(t)$

$x_2(t)$ —[ ]— $y_2(t)$

$x(t)$ — $y(t)$

$x(t-t_0)$ — $y(t-t_0)$
time invariant system.

$\alpha x_1(t) + \beta x_2(t)$ —[ ]— $\alpha y_1(t) + \beta y_2(t)$

Linear system

## Z-transform:

$$X(z) = \sum x(n) z^{-n} \qquad \boxed{z = \sigma e^{j\theta}}$$

$\sigma$ is added because sometimes $\sum x(n)$ is not summable but $\sum x[n] \sigma^{-n}$ may be summable.

When $\sigma = 1$, $z = e^{j\theta}$, $X(z)$ becomes discrete time fourier transform (DTFT).

assume if only $X(z)$ is given, it is insignificant without ROC.

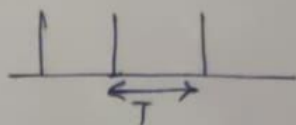## Physiological model of speech production:

message formulation → language coding → Neuromuscular commands → Vocal track System → Speech

acts a resonator (produces frequ..)

↳ Extraction source

laryn: acts as exhicitationer

Rate of vocal fold vibration is called speech (soray pitch)

children > female > male

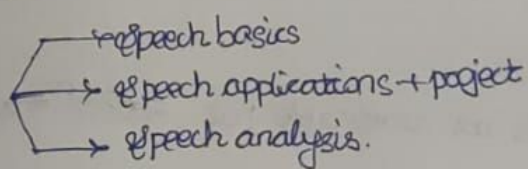Vocal folds when they are open, air passes through vocal folds as an impulse.

$$\frac{1}{T} = pitch = rate\ of\ vocal\ fold\ vibrations.$$

Each impulse = epochs / ISE.

# Speech Signal Processing.

Grading policy:

Quiz- 10%.
midsem - 20%.
assignments - 20%.
Project - 25%. (10%-mid, 15%-end)
Endexam - 25%.

- Speech basics
- Speech applications + project
- Speech analysis.

## why speech processing needed?
Human Computer interaction.
Easy accesable of technology to all classes of poeple.

## How speech is different from other signals?
↳ Involves speech production.

Legal sequence of legal sounds produced by human being.

## what is signal processing:
mathematical approach to extract or manipulate the signals.

## Fourier transform:

$$X(\omega) = \int_{-\infty}^{\infty} x(t) e^{-j\omega t} dt$$

$$x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} d\omega$$

Continuous fourier transform

0-4KHz sounds are understandable by humans.

$x(t)$
↓ sampling
$x[n]$

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

$$x[n] = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(e^{j\omega}) e^{j\omega n} d\omega$$

## Sampling DTFT = DFT

$$X(K) = \sum_{0}^{N-1} x[n] e^{(-j\frac{2\pi}{N})Kn}$$

$$x[n] = \frac{1}{N} \sum X(K) e^{j\frac{2\pi}{N}Kn}$$