

Spoken Lang. Processing

Midterm Exam - Fall 2017

Question 1:

(a)


/k/ plosive, Un-Voiced
/sh/ Fricative, Un-Voiced
/Ah/ Vowel, Voiced
/M/ Nasal, Voiced

(b)

* plosives (stops) e.g. /k/: Very little power (or no power)
Followed by sudden explosion and aspiration when constriction is released.

Waveform e.g. 

* Fricatives (e.g. /sh/): Non-periodic sound (friction) - random energy across wide freq. range.

Waveform e.g. 

* Vowels (e.g. /Ah/): periodic waveform with fundamental frequency corresponding to vibration rate of vocal-folds.
Vocal Tract enhances some frequencies ^(Formants) and suppresses others.
Energy of Vowels sounds is relatively higher than other sounds.

Waveform e.g. 

Q1

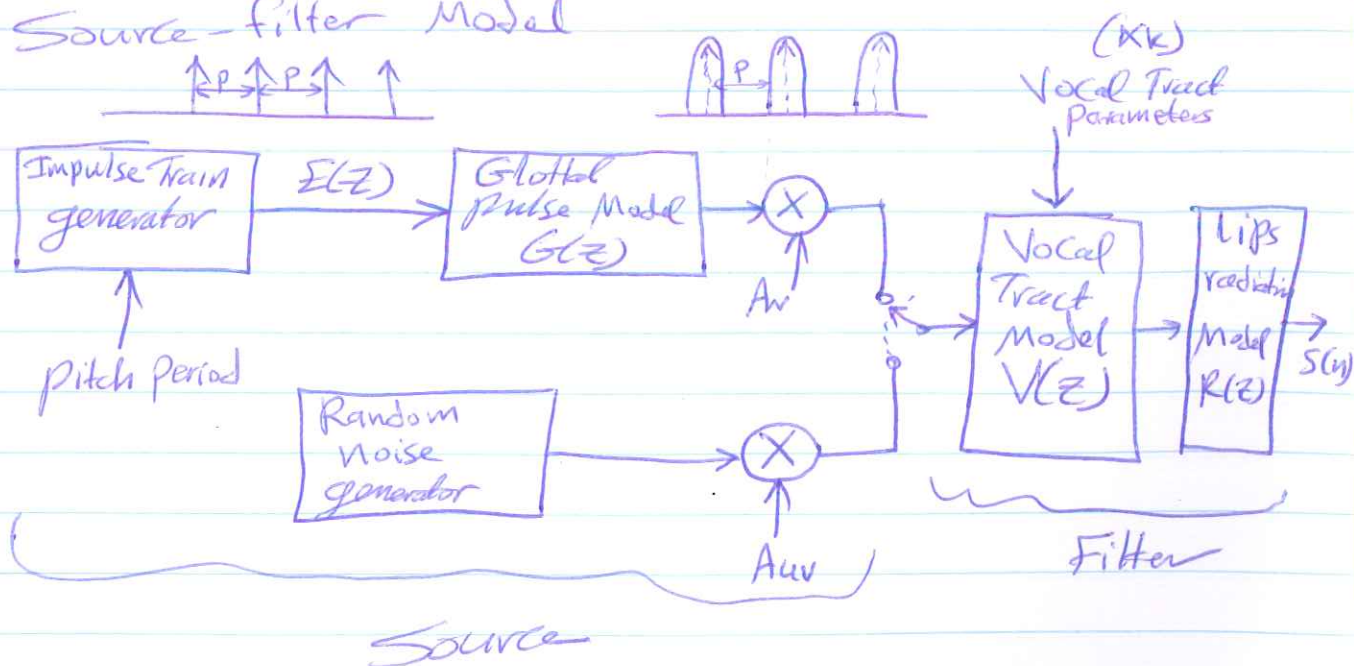
(C) Spectrogram: Time-frequency representation of speech signal. It is a concatenated spectrum of short frames (10-20 ms) and shows the energy distribution over frequency and time scale.



	Window Length	Freq. resolution	Temporal Resolution
Narrow-Band	Wide	Fine	Poor
Wide-Band	Narrow	Poor	Fine

Question 2

(a) Source-filter Model



(2)

- * Impulse Train generator corresponds to Vocal folds Vibration with frequency of F_0 (Fundamental) for Voiced sounds
- * Glottal Model: corresponds to "opening/closing" of the vocal folds in the larynx which shapes the impulses to be like pulses (puffs)
- * For Un-Voiced sounds (eg. /sh/) the output of the larynx is noise-like signal. This is modeled by a random noise generator.
- * In general, Voiced sounds have greater amplitude than Un-Voiced. This is represented by multiplying output of larynx by two constants A_v (Voiced) and A_{uv} (unvoiced)
 $A_v > A_{uv}$.
- * Vocal Tract is modeled by a filter with sufficient number of poles (All-pole filter). Usually 10-14 poles for both Voiced and Un-Voiced sounds.
- * Lip radiation is modeled by a high-pass filter of a transfer function $R(z)$.

(b) Impulse Train: $E(z) = \frac{1}{1 - z^{-P}}$, P : Pitch Period.

Glottal model:

$$G(z) = \begin{cases} \frac{1}{(1 - \frac{1}{P} z^{-1})^2} \approx \frac{1}{(1 - z^{-1})^2}, & \text{for Voiced sounds} \\ 1, & \text{for Unvoiced sounds} \end{cases}$$

Vocal Tract model:

$$V(z) = \begin{cases} \frac{1}{1 + \sum_{k=1}^P a_k z^{-k}}, & \text{for Voiced sound} \\ \frac{1}{1 + \sum_{k=1}^{p+2L} a_k z^{-k}}, & \text{for Un-Voiced sound} \end{cases}$$

$P = \text{no. of poles}$

$L = \text{no. of zeros}$

OR

$$V(z) = \frac{1}{1 + \sum_{k=1}^P a_k' z^{-k}}, \quad P \approx 10-14 \quad \left. \vphantom{\frac{1}{1 + \sum_{k=1}^P a_k' z^{-k}}} \right\}^* \text{For Both Voiced and Un-Voiced Speech.}$$

Lips Radiation:

$$R(z) = 1 - z^{-1}, \quad \text{First-order HPF.}$$

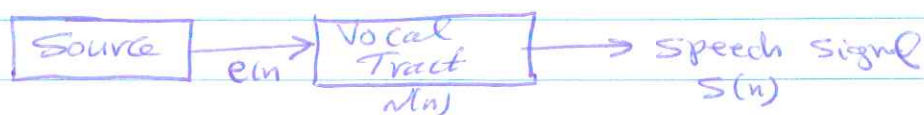
* Overall

$$\frac{S(z)}{E(z)} = \frac{G}{1 + \sum_{k=1}^P a_k' z^{-k}}, \quad 10 \leq P \leq 14$$

Question 2

(c) Right-hand graph is spectrum of speech frame
Left-hand graph is the response of the vocal tract

* Left-hand graph (Vocal tract) is the envelop of the speech spectrum, OR Low Frequency Components of the speech spectrum.



$$s(n) = e(n) * v(n)$$

$$S(z) = E(z) \cdot V(z)$$

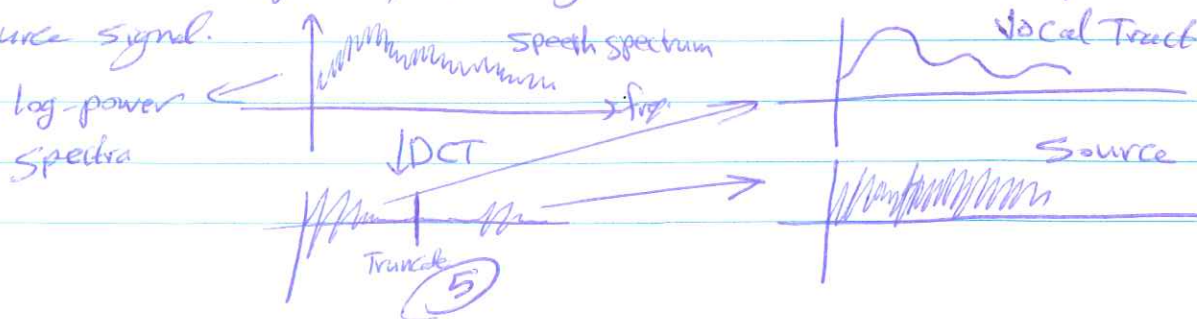
* This looks like a voiced sound because the power is higher in the low-frequency components i.e. power at $f_1 > \text{power at } f_2 > \dots$

(d) Two benefits of using DCT in the MFCC technique.

1. Un-correlated features, so diagonal Covariance Matrix can be used in the modelling.

2. Separate source signal (excitation) from the vocal tract filter response by truncating DCT coefficients.

The low order coefficients (or first 12) corresponds to vocal tract filter response, and high order coefficients corresponds to source signal.



Question 38

(a) $F_s = 600$ Sample/sec.

$$S(n) = \{-0.35, 0.0, 3.15, -2.5, -3.54, 2.8, 0.0, -0.28\}, \quad N=8$$

$$\begin{aligned} \text{(i)} \quad E_{\text{Energy}} &= \frac{1}{N} \sum_{n=0}^{N-1} S(n)^2 = \frac{1}{8} \sum_{n=0}^7 S(n)^2 \\ &= \frac{1}{8} \left[(-0.35)^2 + (0)^2 + (3.15)^2 + (-2.5)^2 + (-3.5)^2 + (2.8)^2 + (0)^2 + \right. \\ &\quad \left. + (-0.28)^2 \right] = 36.4634 \end{aligned}$$

$$E_{\text{dB}} = 10 \log_{10} E = 15.618 \text{ dB}.$$

$$\begin{aligned} \text{(ii)} \quad Z_{\text{CC}} &= \frac{1}{N} \sum_{n=0}^{N-1} \frac{1}{2} \left| \text{sign}(S(n)) - \text{sign}(S(n-1)) \right| \\ &= \frac{1}{8} \cdot (4) = \frac{1}{2} \end{aligned}$$

$$\text{(iii)} \quad R(k) = \frac{1}{N} \sum_{n=0}^{N-1} S(n)S(n-k)$$

$$150 \leq f_0 \leq 300 \Rightarrow \frac{600}{300} \leq k \leq \frac{600}{150} \Rightarrow 2 \leq k \leq 4$$

We compute $R(2)$, $R(3)$ and $R(4) \Rightarrow$ The pitch period corresponds to k which gives $R(k)$ max. value.

Q3. (a) (iii) Continue ---

$S(n)$	-0.35	0	3.15	-2.5	-3.54	2.8	0	-0.28
$S(n-1)$	0	0	-0.35	0	3.15	-2.5	-3.54	2.8
$S(n-2)$	0	0	0	-0.35	0	3.15	-2.5	-3.54
$S(n-3)$	0	0	0	0	-0.35	0	3.15	-2.5

$$R(2) = \frac{1}{8}(3.15)(-0.35) + (3.15)(-3.54) + (2.8)(-2.5) + (2.8)(-0.28) = -2.673$$

$$R(3) = \frac{1}{8}(-2.5)(-0.35) + (2.8)(3.15) + (-0.28)(-3.54) = 1.335$$

$$R(4) = \frac{1}{8}(-3.54)(-0.35) + (-0.28)(-2.5) = 0.2425$$

So, $R(3)$ is the max. value $\Rightarrow k=3$ corresponds to Pitch period.

$$\Rightarrow \text{Pitch period} = 3 * \frac{1}{600} = \frac{1}{200} \text{ sec.} \Rightarrow F_0 = 200 \text{ Hz.}$$

Q3 (b)

High energy + low Z_{cc} + nonzero $F_0 \Rightarrow$ Voiced speech

Low energy + High Z_{cc} + zero $F_0 \Rightarrow$ Unvoiced speech

This speech segment is Voiced because Fundamental frequency (F_0) is not zero ($F_0 = 200 \text{ Hz}$).

Question 3 (c)

- (1) No, because min. sound pressure level for 100 Hz is around 25 dB
- (2) a tone of 1000 Hz at 20 dB because human ear is more sensitive to 1000 Hz than 100 Hz
- (3) B 2 kHz - 5 kHz.