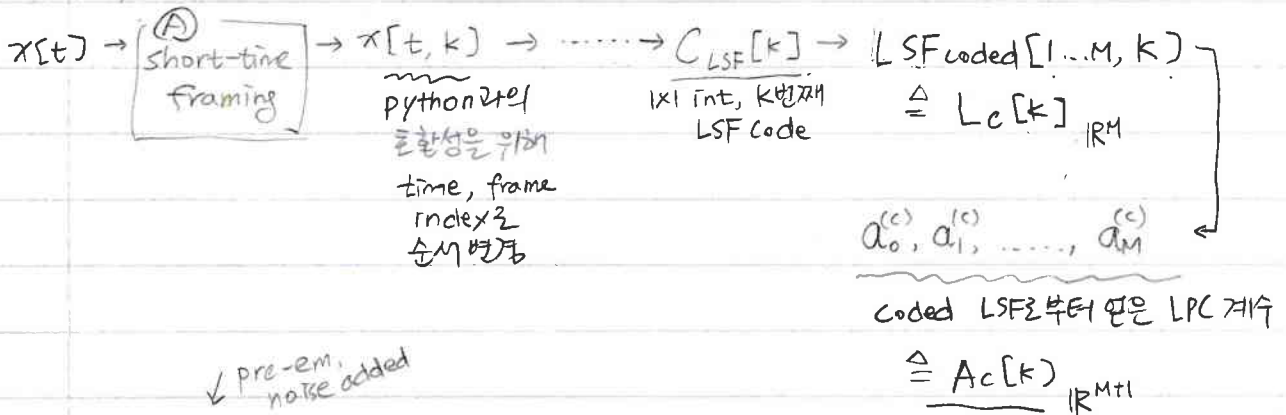


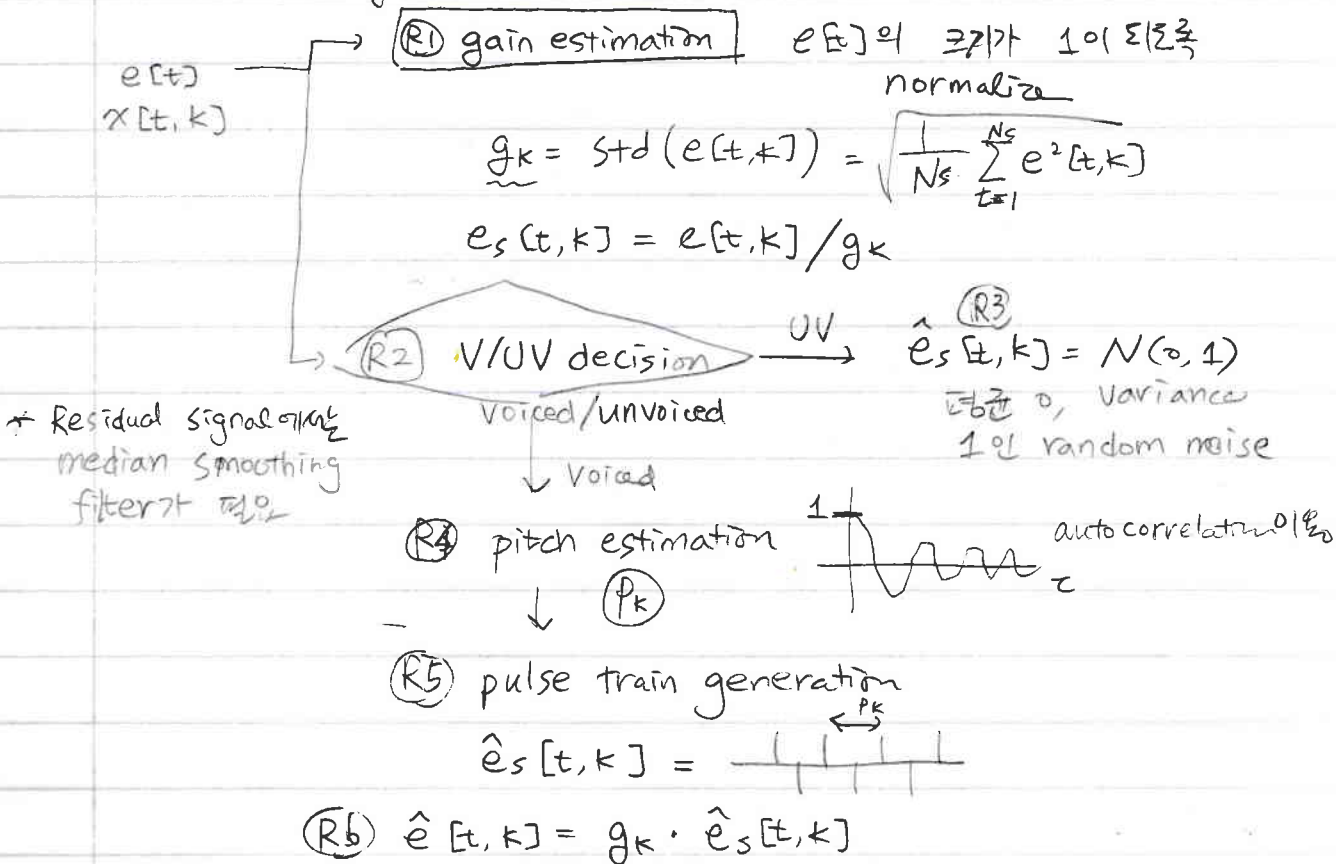
2021. 2. 25 Residual ^{frame} signal encoding

< Encoder >



$\Rightarrow x[t, k] \begin{cases} C_{LSF}[k] : \text{envelope} \\ e[t] = x_{pn} * A_c : \mathbb{R}^{N_s} \end{cases}$

< Residual encoding >



(R1) gain estimation

음성은 보통 다음과 같이 모델링된다

$$X(\omega) = \frac{G \cdot E(\omega)}{A(\omega)}$$

$G \in \mathbb{R}^1$: scalar 값

$E(\omega)$: excitation / residual

$A(\omega)$: all-pole envelope

gain은 의미함.
마이크를 통해 녹음되는

frame gain

short-time framing을 했기 때문에 frame energy의 square-root를 구하면 된다 (즉, standard deviation)

$$g_k = \text{std}(e[t, k]) \approx \sqrt{\frac{1}{N_s} \sum e^2[t, k]}$$

(A) unbiased estimation이 ($N_s - 1$ 로 나누기)

통계적으로 의미있는 frame gain을 구할 수 있지만

N_s 는 고정되어 있고, gain은 reconstruction에서

다시 곱해주기 때문에 굳이 $N_s - 1$ 로 나눌 필요는 없다.

(정확한 frame gain이 필요할 때 사용)

(codec에서는 영향 없음)

(B) 정확한 standard deviation 계산을 위해서는

$$(e[t, k] - \mu_k)^2, \quad \mu_k = \frac{1}{N_s} \sum e[t, k], \quad \text{가 맞지만}$$

dc 성분이라고 함

preemphasis에서 0-frequency (dc) 성분은 제거가 되고

그 이후의 과정에서 dc값이 바뀌는 현상이 있기 때문에

굳이 평균 차감은 필요없다.

* 극단적인 burst, 비선형 noise가 유입된 경우



문제가 될 수도 있기 때문에 평균 차감은 하는 것이 안전하기는 하다.

- R2** Voiced / unvoiced decision (무성음)
- residual (excitation) 신호의 모델링 방법은 voiced (유성음) / unvoiced (무성음)에 따라 크게 달라지기 때문에 매우 중요하다.
 - 정확하게 추정하기는 매우 어렵다 (보조 강의 자료에 model-based 방법 있음)
 - 다음의 방법을 우선 적용해 보고, 잘 안되면 다른 방법들 시도한다.

① Voiced sound의 특징

power가 -6dB/octave (주파수가 2배가 되면 $\frac{1}{2}$ 로 크기감소)

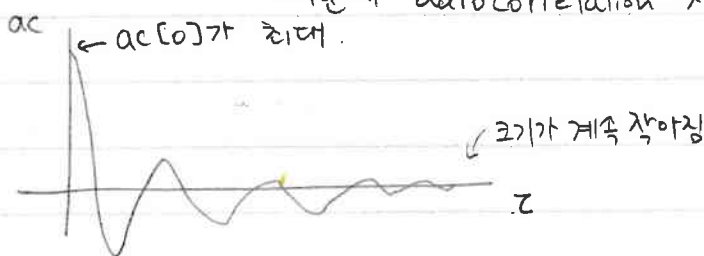
\equiv energy가 -12dB/octave

$ac[\tau] = x[t] \cdot x[t-\tau]$: 간격이 τ 인 autocorrelation coef ($\tau = 0, \dots, N_s$)

$\approx \sum_{t=\tau}^{N_s} x[t] \cdot x[t-\tau]$: short-time frame

등다 OK ($\approx \frac{1}{N_s} \sum_{t=\tau}^{N_s} x[t] \cdot x[t-\tau]$: 범위를 벗어나면 $x[t] = 0$ 이므로 $t = \tau$ 부터 시작

\uparrow normalized denominator (분모)는 N_s 그대로, 즉 τ 가 커질수록 작아짐
원래 infinite series에 대해서는 normalize 불가능하기 때문에 autocorrelation 정의에서는 sample 수로 나누는 것 없음



② unvoiced sound : 주기성이 없으므로 이론적으로 $ac[0] = ac[1] = \dots$

③ decision algorithm

• Compute $ac[0], ac[1], \dots, ac[N_s]$ from $x[t, k]$ framing만 한 신호

• if $ac[0] < \epsilon \rightarrow$ return unvoiced * silence로 따로 분류하는 안고려함도 있지만 어차피 unvoiced와 encoding 방법 같음

elif $ac[1]/ac[0] > \theta_{uv} \rightarrow$ return UV

else return Voiced

* θ_{uv} 는 0.5 권장

unvoiced excitation encoding

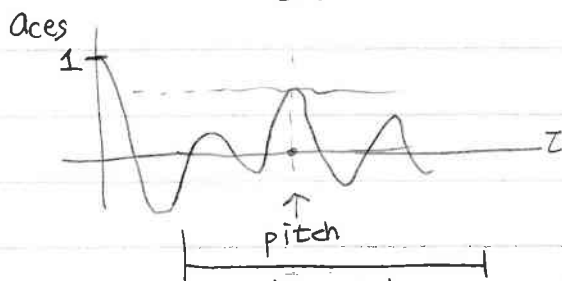
- [R3] Input : g_k , $e_s[t, k]$ ($= e[t, k]/g_k$), $VUV = UV$ (boolean)
 generate Gaussian random noise with $\mu=0$, $\sigma^2=1$ of size N_s
 $\hat{e}_s[t, k] = N(0, 1) [1 \dots N_s]$

- [R4] pitch estimation: 최대 주기 찾기.

* clipping, median filtering 등의 technique은 강의자료를
 읽어보기 바람 여기서 auto correlation을 쓰는 기본적인 방법
 * 주의: VUV decision에서는 raw frame $x[t, k]$ 를 썼지만
 pitch estimation에서는 $x_{pn}[t, k]$ 의 LPC residual (normalized)
 $e_s[t, k]$ 를 쓴다.

① main concept: acco]를 제외한 나머지에서 max 위치가 pitch

$$A_{ces}[z, k] = \sum_{t=z}^{N_s} e_s[t, k] \cdot e_s[t-z, k]$$



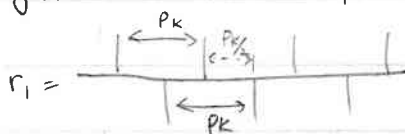
$$p[k] = \arg \max_{p_{min} \leq z \leq p_{max}} A_{ces}[z, k]$$

valid pitch range \Rightarrow 찾아보기 바람

② median filter 등을 통해 smoothing

- [R5] pulse train generator w/ $p[k]$

① generate



* 실제 glottal pulse의 모양은



가 더 적합하지만

② time alignment: $e[t, k]$ 와 correlation이
 가장 높은 alignment 찾는 것

+/- 신호의 간격을 정확히
 추정 못하기 때문에
 $p_k / 3$ 정도로 +/- 생략.

$$z_a = \arg \max_{0 \leq z \leq p_{k-1}} e[t, k] \cdot r_1[t+z]$$

$$r_2[t] = r_1[t+z]$$

(주의) pulse의 개수가 같아야 함

③ normalize $\hat{e}_s[t, k] = r_2[t] / \sqrt{\frac{1}{N_s} \sum_{t=1}^{N_s} r_2^2[t]}$

Encoder output: $C_{LSF}[k] \in \mathbb{Z}$ (integer)

$flag_{vuv} \in \{T, F\}$

$g[k] \in \mathbb{R}$

$p[k] \in \mathbb{Z}$ (int) \leftarrow voiced only

Decoding:

① $C_{LSF} \rightarrow LSF[1..M] \rightarrow (a_0, a_1, \dots, a_M) \in \mathbb{R}^{M+1}$

② if $flag_{vuv}$ is F (unvoiced or silence)

$\hat{e}_s[t, k] = N(0, 1)[1..N_s] \in \mathbb{R}^{N_s} \dots \boxed{R3}$

③ else (voiced)

$\hat{e}_s[t, k] = \text{periodic pulse train} \dots \boxed{R5}$

④ apply gain: $\hat{e}[t, k] = g[k] \cdot \hat{e}_s[t, k] \dots \boxed{R6}$

⑤ $\hat{x}_{pn}[t, k] = \frac{\hat{e}[t, k]}{a[0..M]}$ (IIR filtering)

Programming assignment 3

R1 ~ R6 구현하고 encoder / decoder 구현

- V/UV 결과를 spectrogram 위에 그린다.

- pitch 도 그린다.