

Speech Enhancement

Single and Dual Channel Noise Reduction

- Noisy speech is everywhere. For example, in a bus/train.
- Objective: ① Remove the noise
or ② Estimate the good speech from noisy speech.
- We consider additive model.

$$y(k) = x(k) + n(k)$$



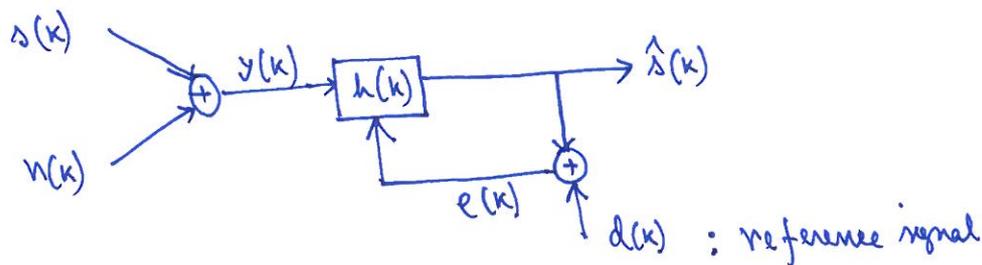
- Assumption: $s(k)$ and $n(k)$ are statistically independent.
So, $E\{s(k)n(i)\} = 0, \forall k, i.$

Linear MMSE Estimators

- We do some linear operation on $y(k)$ and get an $\hat{s}(k)$.
~~However~~ MMSE means minimum $E\{|s(k) - \hat{s}(k)|^2\}$.

Non-causal IIR Wiener Filter

- Linear Filter



$$\hat{s}(k) = \sum_{l=-\infty}^{+\infty} h(l) y(k-l) = y(k) * h(k)$$

- Assume target signal $d(k) = s(k)$. and error $e(k) = \hat{s}(k) - s(k)$.

$$\begin{aligned} \bullet \text{ minimize } E\{e^2(k)\} &= E\{(\hat{s}(k) - d(k))^2\} \\ &= E\{(\hat{s}(k) - s(k))^2\} \\ &= E\left\{\left(\sum_{l=-\infty}^{+\infty} h(l) y(k-l)\right) - s(k)\right\}^2 \end{aligned}$$

$$\bullet \text{ for } i, \frac{\partial E\{e^2(k)\}}{\partial h(i)} = 0$$

$$\text{or, } E\left\{2e(k) \cdot \frac{\partial e(k)}{\partial h(i)}\right\} = 0$$

$$\text{or, } E\left\{e(k) \cdot \frac{\partial \hat{s}(k)}{\partial h(i)}\right\} = 0$$

$$\text{or, } E\{e(k) \cdot y(k-i)\} = 0$$

$$\text{or, } E\{(\hat{s}(k) - s(k)) y(k-i)\} = 0$$

$$\text{or, } E\{\hat{s}(k) y(k-i)\} = E\{s(k) y(k-i)\}$$

$$\text{or, } E\left\{\left(\sum_{l=-\infty}^{+\infty} h(l) y(k-l)\right) y(k-i)\right\} = E\{s(k) y(k-i)\}$$

$$\text{or, } \boxed{\sum_{l=-\infty}^{+\infty} h(l) \psi_{yy}(i-l) = \psi_{ys}(i)}$$

[Assumption: all signals are wide-sense stationary (second-order stationary)]

$$\sum_{l=-\infty}^{+\infty} h(l) \psi_{yy}(i-l) = \psi_{ys}(i)$$

$$\text{or, } h(i) * \psi_{yy}(i) = \psi_{ys}(i) \quad [\text{convolution relation}]$$

• Note: Convolution of two signals can be represented by multiplication in the power spectral domain.

$$H(e^{j\Omega}) \cdot \phi_{yy}(e^{j\Omega}) = \phi_{ys}(e^{j\Omega})$$

$$\text{or, } H(e^{j\Omega}) = \frac{\phi_{ys}(e^{j\Omega})}{\phi_{yy}(e^{j\Omega})} \quad \left[\begin{array}{l} h(k) = \frac{1}{2\pi} \int_{-\pi}^{\pi} H(e^{j\Omega}) e^{j\Omega k} d\Omega \\ \text{and also, assuming } \phi_{yy}(e^{j\Omega}) \neq 0 \end{array} \right]$$

$$\text{Now, } y(k) = s(k) + n(k)$$

$$\begin{aligned} \therefore E\{y(k) y(k-i)\} &= E\{(s(k) + n(k)) (s(k-i) + n(k-i))\} \\ &= E\{s(k) s(k-i)\} + E\{n(k) s(k-i)\} \\ &\quad + E\{s(k) n(k-i)\} + E\{n(k) n(k-i)\} \end{aligned}$$

$$\text{or, } \psi_{yy}(i) = \psi_{ss}(i) + \psi_{nn}(i)$$

[as $s(k)$ and $n(k)$ are statistically independent]

$$\text{or, } \sum_{\forall i} \psi_{yy}(i) e^{-j\Omega i} = \sum_{\forall i} \psi_{ss}(i) e^{-j\Omega i} + \sum_{\forall i} \psi_{nn}(i) e^{-j\Omega i}$$

$$\text{or, } \phi_{yy}(e^{j\Omega}) = \phi_{ss}(e^{j\Omega}) + \phi_{nn}(e^{j\Omega})$$

~~$$H(e^{j\Omega}) = \frac{\phi_{ys}(e^{j\Omega})}{\phi_{yy}(e^{j\Omega})}$$~~

$$\begin{aligned}
 Y_{ys}(i) &= E \{ s(k) y(k-i) \} \\
 &= E \left\{ s(k) \left(s(k-i) + n(k-i) \right) \right\} \\
 &= E \{ s(k) s(k-i) \} + E \{ s(k) n(k-i) \} \\
 &= Y_{ss}(i) \quad \left[\text{statistical independence} \right]
 \end{aligned}$$

$$\begin{aligned}
 \therefore \phi_{yy}(e^{j\Omega}) &= \sum_{\forall i} Y_{ys}(i) e^{-j\Omega i} \\
 &= \sum_{\forall i} Y_{ss}(i) e^{-j\Omega i} = \phi_{ss}(e^{j\Omega}).
 \end{aligned}$$

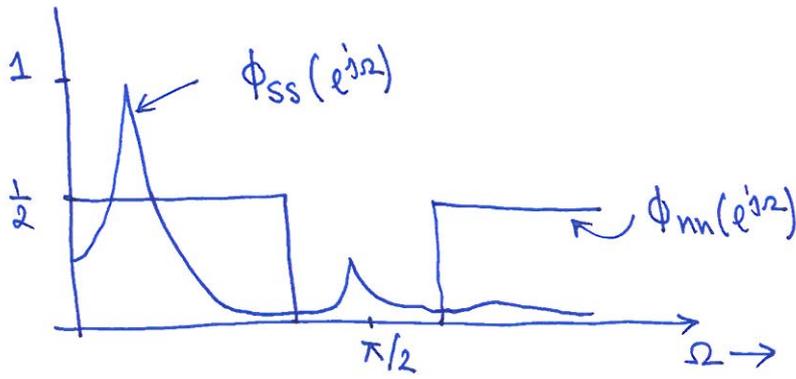
$$\boxed{H(e^{j\Omega}) = \frac{\phi_{ys}(e^{j\Omega})}{\phi_{yy}(e^{j\Omega})} = \frac{\phi_{ss}(e^{j\Omega})}{\phi_{ss}(e^{j\Omega}) + \phi_{nn}(e^{j\Omega})}} = \frac{\frac{\phi_{ss}(e^{j\Omega})}{\phi_{nn}(e^{j\Omega})}}{1 + \frac{\phi_{ss}(e^{j\Omega})}{\phi_{nn}(e^{j\Omega})}}.$$

Filter PSD.

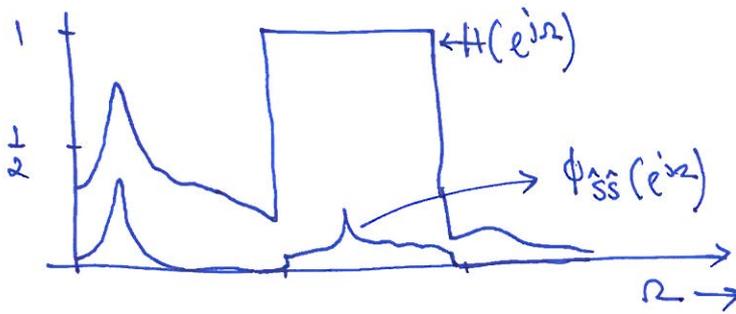
$$\hat{S}(e^{j\Omega}) = H(e^{j\Omega}) \cdot Y(e^{j\Omega}) \quad \text{Output PSD.}$$

The non-causal IIR Wiener Filter evaluates signal-to-noise (SNR) ratio $\frac{\phi_{ss}(e^{j\Omega})}{\phi_{nn}(e^{j\Omega})}$ at a given frequency Ω .

(5)



• When SNR is large,
 $H(e^{j\Omega}) \rightarrow 1$.
The corresponding frequency
component comes out
without attenuation.



• When SNR is low,
 $\phi_{ss}(e^{j\Omega}) \ll \phi_{nn}(e^{j\Omega})$
and $H(e^{j\Omega}) \approx 0$.
So, corresponding frequency
component is attenuated.

- Important Note : The noise reduction task will be efficient if the speech signal and noise do not occupy the same frequency bands.

FIR WIENER FILTER

- We now restrict the impulse response $h(k)$ of the optimal filter to be of finite and even order N , and to be causal, i.e.,

$$h(k) = \begin{cases} \text{arbitrary} & 0 \leq k \leq N \\ 0 & \text{otherwise} \end{cases}$$

- So, $\hat{s}(k) = h(k) * y(k) = \sum_{l=0}^N h(l) y(k-l)$

- We again try to minimize

$$E\{e^2(k)\} = E\{(\hat{s}(k) - d(k))^2\}$$

- Assuming $d(k) = s(k)$, we will again can get

$$\sum_{l=0}^N h(l) \gamma_{yy}(i-l) = \gamma_{ys}(i), \forall i$$

also, due to statistical independence $\gamma_{ys}(i) = \gamma_{sy}(i)$.

$$\therefore \sum_{l=0}^N h(l) \gamma_{yy}(i-l) = \gamma_{sy}(i), \forall i$$

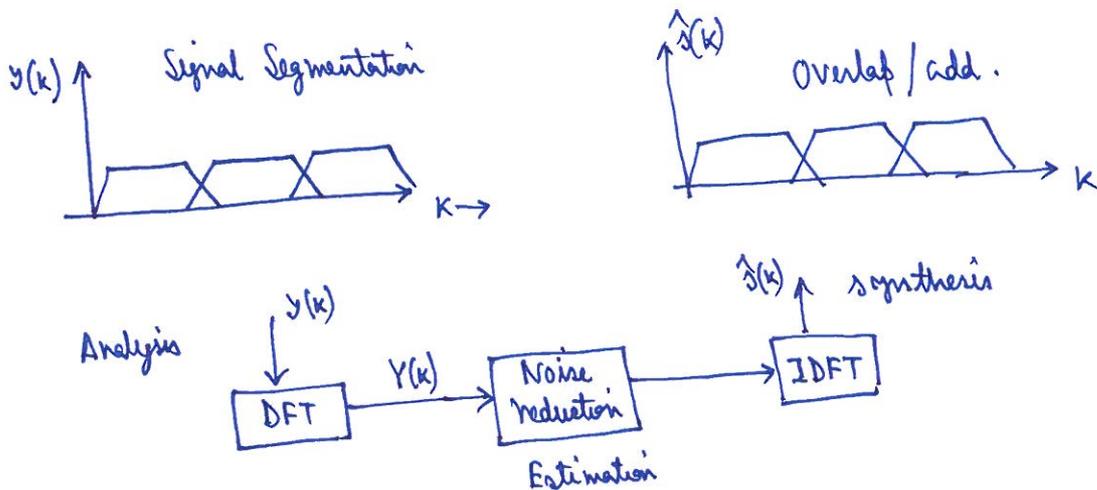
$$\underbrace{\begin{bmatrix} \gamma_{yy}(0) & \gamma_{yy}(1) & \dots & \gamma_{yy}(N) \\ \gamma_{yy}(1) & \gamma_{yy}(0) & \dots & \gamma_{yy}(-N+1) \\ \vdots & \vdots & \ddots & \vdots \\ \gamma_{yy}(N) & \dots & \dots & \gamma_{yy}(0) \end{bmatrix}}_{R_{yy}} \underbrace{\begin{bmatrix} h(0) \\ h(1) \\ \vdots \\ h(N) \end{bmatrix}}_{\underline{h}} = \underbrace{\begin{bmatrix} \gamma_{sy}(0) \\ \gamma_{sy}(1) \\ \vdots \\ \gamma_{sy}(N) \end{bmatrix}}_{\underline{\gamma}_{ss}}$$

$$\therefore R_{yy} \underline{h} = \underline{\gamma}_{ss}$$

$$\underline{h} = R_{yy}^{-1} \underline{\gamma}_{ss}$$

Speech Enhancement in the DFT domain

- IIR Wiener Filter is a time domain filter, but it gives us clear understanding in the frequency domain.
- For a time domain signal, the frequency domain behavior and techniques can be realized by using either Fourier Transform or filter banks.
- Let us treat the estimation in the DFT domain.



- M-size segment (or frame)

$$\underline{y}(k) = (y(k-M+1) \quad y(k-M+2) \quad \dots \quad y(k))^T$$

$$\underline{s}(k) = (s(k-M+1) \quad s(k-M+2) \quad \dots \quad s(k))^T$$

$$\underline{n}(k) = (n(k-M+1) \quad n(k-M+2) \quad \dots \quad n(k))^T$$

- An analysis window $\underline{w} = (w(0) \quad w(1) \quad \dots \quad w(M-1))^T$

$$\underline{Y}(k) = \text{DFT} \left\{ \underline{w} \otimes \underline{y}(k) \right\}$$

$$\underline{S}(k) = \text{DFT} \left\{ \underline{w} \otimes \underline{s}(k) \right\}$$

$$\underline{N}(k) = \text{DFT} \left\{ \underline{w} \otimes \underline{n}(k) \right\}$$

[Here \otimes denotes multiplication element-by-element]

- $\underline{Y}(k) = (Y_0(k) \dots Y_\mu(k) \dots, Y_{M-1}(k))^T$
same for $\underline{S}(k)$ and $\underline{N}(k)$.

Now, we come back to Wiener Filter

- Time domain convolution \leftrightarrow Frequency domain multiplication (Fourier)

- Enhanced signal $\hat{\underline{s}}(k)$ has the DFT $\hat{\underline{S}}(k)$

- $\hat{\underline{S}}(k) = \underline{H}(k) \otimes \underline{Y}(k)$
where $\underline{H}(k) = \text{DFT} \{ \underline{h}(k) \}$

- We want to minimize

$$E \{ \| \underline{s}(k) - \hat{\underline{s}}(k) \|^2 \}$$

- But for windowed signal, we have to minimize

$$E \{ \| (\underline{w} \otimes \underline{s}(k)) - (\underline{w} \otimes \hat{\underline{s}}(k)) \|^2 \}$$

$$= E \{ \| \underline{s}(k) - \hat{\underline{s}}(k) \|^2 \} \quad \left[\begin{array}{l} \text{As } \underline{s}(k) = \text{DFT} \{ \underline{w} \otimes \underline{s}(k) \} \\ \text{and DFT is an orthonormal} \\ \text{transform} \end{array} \right]$$

- assuming statistical independence between the DFT coefficients, we can minimize

$$E \{ | S_\mu(k) - \hat{S}_\mu(k) |^2 \}, \forall \mu.$$

$$\bullet \quad E \left\{ |S_{\mu}(k) - \hat{S}_{\mu}(k)|^2 \right\} = E \left\{ (S_{\mu}(k) - H_{\mu}(k) Y_{\mu}(k)) (S_{\mu}(k) - H_{\mu}(k) Y_{\mu}(k))^* \right\}$$

$$\bullet \quad \textcircled{a} \quad \frac{\partial E \left\{ |S_{\mu}(k) - \hat{S}_{\mu}(k)|^2 \right\}}{\partial \operatorname{Re} \left\{ H_{\mu}(k) \right\}} = 0$$

$$\Rightarrow \operatorname{Re} \left\{ H_{\mu}(k) \right\} = \frac{E \left\{ |S_{\mu}(k)|^2 \right\}}{E \left\{ |Y_{\mu}(k)|^2 \right\}} = \frac{E \left\{ |S_{\mu}(k)|^2 \right\}}{E \left\{ |S_{\mu}(k)|^2 \right\} + E \left\{ |N_{\mu}(k)|^2 \right\}}$$

$$\bullet \quad \textcircled{b} \quad \frac{\partial E \left\{ |S_{\mu}(k) - \hat{S}_{\mu}(k)|^2 \right\}}{\partial \operatorname{Im} \left\{ H_{\mu}(k) \right\}} = 0$$

$$\Rightarrow \operatorname{Im} \left\{ H_{\mu}(k) \right\} = 0.$$

$$\text{Therefore,} \quad H_{\mu}(k) = \frac{E \left\{ |S_{\mu}(k)|^2 \right\}}{E \left\{ |S_{\mu}(k)|^2 \right\} + E \left\{ |N_{\mu}(k)|^2 \right\}}.$$

- The DFT based solution is straight-forward. But, suffers in a very difficult practical problem.

- The clean speech frequency domain information is not available.

- To get a practical solution, we discuss "SPECTRAL SUBTRACTION".

Spectral Subtraction

- The basic idea is to subtract an estimate of the noise floor from an estimate of the spectrum of the noisy signal.

- $\Omega_\mu = \frac{2\pi}{M} \cdot \mu, \mu = 0, 1, \dots$

- $$\begin{aligned} E \{ |S_\mu(k)|^2 \} &= E \{ |Y_\mu(k)|^2 \} - E \{ |N_\mu(k)|^2 \} \\ &= E \{ |Y_\mu(k)|^2 \} \left[1 - \frac{E \{ |N_\mu(k)|^2 \}}{E \{ |Y_\mu(k)|^2 \}} \right] \\ &= E \{ |Y_\mu(k)|^2 \} \cdot |\tilde{H}_\mu(k)|^2 \end{aligned}$$

- A DFT may be interpreted as an analysis filter bank, $E \{ |S_\mu(k)|^2 \}$ represents the power of a complex-valued subband signal $S_\mu(k)$ for any fixed μ .

- The spectral subtraction method may be interpreted as a filter with frequency response

$$\tilde{H}_\mu(k) = \sqrt{1 - \frac{E \{ |N_\mu(k)|^2 \}}{E \{ |Y_\mu(k)|^2 \}}} = \sqrt{\frac{E \{ |S_\mu(k)|^2 \}}{E \{ |Y_\mu(k)|^2 \}}} = \sqrt{H_\mu(k)}$$

- ~~This~~ Since we are subtracting in the PSD domain, this approach is called "power subtraction".

- We can use other variations of the spectral gain function, such as "magnitude subtraction"

$$\sqrt{E\{|Y_{\mu}(k)|^2\}} - \sqrt{E\{|N_{\mu}(k)|^2\}} = \sqrt{E\{|Y_{\mu}(k)|^2\}} \left[1 - \frac{\sqrt{E\{|N_{\mu}(k)|^2\}}}{\sqrt{E\{|Y_{\mu}(k)|^2\}}} \right]$$

These variations can be written in a generalized version

$$|\tilde{H}_{\mu}(k)|^2 = \left[1 - \left(\frac{E\{|N_{\mu}(k)|^2\}}{E\{|Y_{\mu}(k)|^2\}} \right)^{\beta} \right]^{\alpha}$$

- Using magnitude squared DFT spectra and short-time estimates $|\hat{N}_{\mu}(k)|^2$ and $|\hat{Y}_{\mu}(k)|^2$, a practical final approach can be

$$|\hat{S}_{\mu}(k)|^2 = |Y_{\mu}(k)|^2 \left[1 - \left(\frac{|\hat{N}_{\mu}(k)|^2}{|\hat{Y}_{\mu}(k)|^2} \right)^{\beta} \right]^{\alpha}$$

- One point: In practice we have to make sure $|\hat{S}_{\mu}(k)|^2$ is non-negative and real valued.