

ANALYSIS OF THE FREQUENCY-DOMAIN WIENER FILTER WITH THE PREDICTION GAIN

Jingdong Chen¹, Jacob Benesty², and Yiteng (Arden) Huang¹

¹: WeVoice, Inc.
9 Sylvan Dr.
Bridgewater, NJ 08807, USA

²: INRS-EMT, University of Quebec
800 de la Gauchetiere Ouest, Suite 6900
Montreal, QC H5A 1K6, Canada

ABSTRACT

This paper presents a theoretical analysis on the performance of the optimal noise-reduction filter in the frequency domain. Using the autoregressive (AR) model to model both the clean speech and noise, we build the relationship between the Wiener filter and the AR parameters of the clean speech and noise signals. We show that if noise is not predictable, the Wiener filter is mostly related to the AR parameters of the desired speech signal. On the contrary, if the desired signal is not predictable, the Wiener filter is then mostly related to the AR parameters of the noise signal. More importantly, we provide the bounds for noise reduction, speech distortion, and SNR improvement, and show that the performance of the Wiener filter in terms of SNR improvement and degree of noise reduction and speech distortion is closely related to the prediction gain of the desired speech and noise signals.

Index Terms— Noise reduction, Wiener filter, autoregressive model, prediction gain.

1. PROBLEM FORMULATION, SIGNAL MODEL, AND PREDICTION GAIN

The noise reduction problem considered in this paper is one of recovering the desired signal (clean speech) $x(k)$, k being the discrete-time index, of zero mean from the noisy observation (microphone signal)

$$y(k) = x(k) + v(k), \quad (1)$$

where $v(k)$ is the unwanted additive noise, which is assumed to be a zero-mean random process (white or colored) and uncorrelated with $x(k)$. Using the z -transform, (1) can be rewritten as

$$Y(z) = X(z) + V(z), \quad (2)$$

where $Y(z)$, $X(z)$, and $V(z)$ are the z -transforms of $y(k)$, $x(k)$, and $v(k)$, respectively. In the rest, we will always take $z = e^{j\omega}$, where j is the imaginary unit ($j^2 = -1$) and ω ($-\pi < \omega \leq \pi$) is the angular frequency.

Since $x(k)$ and $v(k)$ are assumed to be uncorrelated, we have

$$\phi_y(\omega) = \phi_x(\omega) + \phi_v(\omega), \quad (3)$$

where

$$\phi_a(\omega) = E \left[|A(e^{j\omega})|^2 \right] \quad (4)$$

is the power spectral density (PSD) of the signal $a(k)$, $a \in \{x, v, y\}$, and $E[\cdot]$ denotes the mathematical expectation.

An estimate of $X(e^{j\omega})$ can be obtained by multiplying $Y(e^{j\omega})$ with a complex gain, $H(e^{j\omega})$, i.e.,

$$Z(e^{j\omega}) = H(e^{j\omega})Y(e^{j\omega}) = X_f(e^{j\omega}) + V_{rn}(e^{j\omega}), \quad (5)$$

where $Z(e^{j\omega})$ is the frequency-domain representation of the signal $z(k)$, $X_f(e^{j\omega}) \triangleq H(e^{j\omega})X(e^{j\omega})$ is the filtered clean speech, and $V_{rn}(e^{j\omega}) \triangleq H(e^{j\omega})V(e^{j\omega})$ is the residual noise. From (5), we deduce the PSD of $z(k)$:

$$\phi_z(\omega) = |H(e^{j\omega})|^2 \phi_y(\omega) = |H(e^{j\omega})|^2 [\phi_x(\omega) + \phi_v(\omega)]. \quad (6)$$

The objective of noise reduction in the frequency domain is then to find an optimal gain $H(e^{j\omega})$ at each frequency ω that would attenuate the noise as much as possible with as little distortion as possible to the desired signal (speech).

To better understand and analyze the optimal gains in the context of noise reduction, we propose to use the autoregressive (AR) model [5] for both the clean speech and noise signals. With this popular linear stochastic model, these two signals can be written as linear combination of their past values, i.e.,

$$x(k) = \sum_{l=1}^{L_x} a_{x,l} x(k-l) + u_x(k), \quad (7)$$

$$v(k) = \sum_{l=1}^{L_v} a_{v,l} v(k-l) + u_v(k), \quad (8)$$

where $a_{x,l}$ and $a_{v,l}$ are the AR parameters of $x(k)$ and $v(k)$, respectively, and $u_x(k)$ and $u_v(k)$ are two zero-mean random white noise signals. In the frequency domain, (7) and (8) are

$$X(e^{j\omega}) = \frac{U_x(e^{j\omega})}{A_x(e^{j\omega})}, \quad (9)$$

$$V(e^{j\omega}) = \frac{U_v(e^{j\omega})}{A_v(e^{j\omega})}, \quad (10)$$

where $U_x(e^{j\omega})$ and $U_v(e^{j\omega})$ are the frequency-domain representations of $u_x(k)$ and $u_v(k)$, respectively, and

$$A_x(e^{j\omega}) = 1 - \sum_{l=1}^{L_x} a_{x,l} e^{-j\omega l}, \quad (11)$$

$$A_v(e^{j\omega}) = 1 - \sum_{l=1}^{L_v} a_{v,l} e^{-j\omega l}, \quad (12)$$

are two minimum-phase polynomials. With this AR modelling, the PSDs of the speech and noise signals become

$$\phi_x(\omega) = \frac{\sigma_{u_x}^2}{|A_x(e^{j\omega})|^2}, \quad (13)$$

$$\phi_v(\omega) = \frac{\sigma_{u_v}^2}{|A_v(e^{j\omega})|^2}, \quad (14)$$

where $\sigma_{u_x}^2 = E[u_x^2(k)]$ and $\sigma_{u_v}^2 = E[u_v^2(k)]$ are the variances of $u_x(k)$ and $u_v(k)$, respectively.

We are now ready to define a measure that naturally follows from the AR modelling. We define the subband prediction gain of the desired signal $x(k)$ as

$$\mathcal{P}_x(\omega) \triangleq \frac{\phi_x(\omega)}{\sigma_{u_x}^2} = \frac{1}{|A_x(e^{j\omega})|^2}. \quad (15)$$

If $\mathcal{P}_x(\omega) = 1$ then the signal $x(k)$ is completely unpredictable at frequency ω . A larger value of $\mathcal{P}_x(\omega)$ means that the signal $x(k)$ is more predictable at frequency ω .

In a similar manner, we define the fullband prediction gain [6] of the desired signal $x(k)$ as

$$\mathcal{P}_x \triangleq \frac{\sigma_x^2}{\sigma_{u_x}^2}, \quad (16)$$

where $\sigma_x^2 = E[x^2(k)]$ is the variance of $x(k)$. Using the relation between the variance and the PSD of $x(k)$:

$$\sigma_x^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} \phi_x(\omega) d\omega = \frac{\sigma_{u_x}^2}{2\pi} \int_{-\pi}^{\pi} \frac{d\omega}{|A_x(e^{j\omega})|^2}, \quad (17)$$

(16) can be rewritten as

$$\mathcal{P}_x = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{d\omega}{|A_x(e^{j\omega})|^2} = \frac{1}{2\pi} \int_{-\pi}^{\pi} \mathcal{P}_x(\omega) d\omega. \quad (18)$$

It can be checked that $\mathcal{P}_x \geq 1$. As a matter of fact, using (7) we find that the variance of $x(k)$ is

$$\sigma_x^2 = \mathbf{a}_x^T \mathbf{R}_x \mathbf{a}_x + \sigma_{u_x}^2, \quad (19)$$

where

$$\mathbf{a}_x = [a_{x,1} \quad a_{x,2} \quad \cdots \quad a_{x,L_x}]^T,$$

and $\mathbf{R}_x = E[\mathbf{x}(k)\mathbf{x}^T(k)]$ is the correlation matrix of

$$\mathbf{x}(k) = [x(k) \quad x(k-1) \quad \cdots \quad x(k-L_x+1)]^T.$$

From (19), we get another form of the fullband prediction gain:

$$\mathcal{P}_x = 1 + \frac{\mathbf{a}_x^T \mathbf{R}_x \mathbf{a}_x}{\sigma_{u_x}^2}. \quad (20)$$

Since $\mathbf{a}_x^T \mathbf{R}_x \mathbf{a}_x \geq 0$, we deduce from (20) that $\mathcal{P}_x \geq 1$.

If $\mathcal{P}_x = 1$, this means that the signal $x(k)$ is completely unpredictable (white noise). The larger is \mathcal{P}_x , the more predictable is the signal $x(k)$.

Obviously, the same definitions of subband and fullband prediction gains apply to the noise signal $v(k)$ by just replacing in the previous derivations $x(k)$ by $v(k)$.

By analogy to the fullband prediction gain defined in (18), we define the fullband prediction gains of the filtered desired and residual noise signals as, respectively,

$$\mathcal{P}_x(H) = \frac{1}{2\pi} \int_{-\pi}^{\pi} |H(e^{j\omega})|^2 \mathcal{P}_x(\omega) d\omega \quad (21)$$

and

$$\mathcal{P}_v(H) = \frac{1}{2\pi} \int_{-\pi}^{\pi} |H(e^{j\omega})|^2 \mathcal{P}_v(\omega) d\omega. \quad (22)$$

We can verify that if $|H(e^{j\omega})|^2 \leq 1, \forall \omega$, then $\mathcal{P}_x(H) \leq \mathcal{P}_x$ and $\mathcal{P}_v(H) \leq \mathcal{P}_v$ with equalities if and only if $|H(e^{j\omega})|^2 = 1, \forall \omega$. Contrary to \mathcal{P}_x and \mathcal{P}_v , $\mathcal{P}_x(H)$ and $\mathcal{P}_v(H)$ can be smaller than 1.

We will see that the noise reduction (or speech enhancement) problem can be exclusively reformulated as a function of the prediction gains.

2. PERFORMANCE MEASURES

To facilitate the analysis and interpretation of the noise-reduction performance, some performance measures are presented in this section.

The most important and reliable measure in noise reduction is the signal-to-noise ratio (SNR) [3]. The fullband input SNR is defined as the ratio of the intensity of the signal of interest over the intensity of the additive noise, i.e.,

$$\text{iSNR} \triangleq \frac{\sigma_x^2}{\sigma_v^2} = \text{iSNR}_0 \cdot \frac{\mathcal{P}_x}{\mathcal{P}_v}, \quad (23)$$

where $\text{iSNR}_0 \triangleq \sigma_{u_x}^2 / \sigma_{u_v}^2$. It is worth noticing how the fullband prediction gains affect the SNR. Indeed, if the desired signal is more predictable than the noise signal then $\text{iSNR} > \text{iSNR}_0$. But if the noise signal is more predictable than the desired signal then $\text{iSNR} < \text{iSNR}_0$. The subband input SNR is

$$\text{iSNR}(\omega) \triangleq \frac{\phi_x(\omega)}{\phi_v(\omega)} = \text{iSNR}_0 \cdot \frac{\mathcal{P}_x(\omega)}{\mathcal{P}_v(\omega)}. \quad (24)$$

After noise reduction with the frequency-domain model given in (5), the subband output SNR is

$$\text{oSNR} [H(e^{j\omega})] \triangleq \frac{|H(e^{j\omega})|^2 \phi_x(\omega)}{|H(e^{j\omega})|^2 \phi_v(\omega)} = \text{iSNR}(\omega). \quad (25)$$

So, the subband SNR is not affected by the filtering process. But the fullband output SNR is

$$\text{oSNR}(H) \triangleq \frac{\int_{-\pi}^{\pi} |H(e^{j\omega})|^2 \phi_x(\omega) d\omega}{\int_{-\pi}^{\pi} |H(e^{j\omega})|^2 \phi_v(\omega) d\omega} = \text{iSNR}_0 \cdot \frac{\mathcal{P}_x(H)}{\mathcal{P}_v(H)}. \quad (26)$$

Apparently, the prediction gains will affect the fullband output SNR. We also define the fullband SNR gain as

$$\mathcal{G}(H) \triangleq \frac{\text{oSNR}(H)}{\text{iSNR}} = \frac{\mathcal{P}_x(H)}{\mathcal{P}_x} \cdot \frac{\mathcal{P}_v}{\mathcal{P}_v(H)}. \quad (27)$$

For a constant $\alpha \neq 0$, $\mathcal{G}(\alpha H) = \mathcal{G}(H)$. Therefore, changing the gains $H(e^{j\omega})$ by a scaling factor (same over all frequencies) will not affect the fullband SNR gain. We also observe in (27) that making the filtered desired signal more predictable than the desired signal or making the residual noise signal less predictable than the original noise signal, will increase the fullband SNR gain.

It is of great importance to design the gains $H(e^{j\omega})$ in such a way that $\mathcal{G}(H) > 1$; this would mean that the output SNR would be improved and the estimated signal, $z(k)$, would be less noisy than the microphone signal, $y(k)$. We see from (27) that $\mathcal{G}(H)$ depends exclusively on the prediction gains, \mathcal{P}_x , $\mathcal{P}_x(H)$, \mathcal{P}_v , and $\mathcal{P}_v(H)$. Clearly, the prediction gains of the desired and noise signals play a fundamental role in noise reduction and will affect the design and performance of the gains, $H(e^{j\omega})$.

Another important measure in noise reduction is the noise-reduction factor [1], [2], which quantifies the amount of noise being attenuated by the gains. With the frequency-domain formulation, these subband and fullband factors are defined as, respectively,

$$\xi_{\text{nr}} [H(e^{j\omega})] \triangleq \frac{\phi_v(\omega)}{|H(e^{j\omega})|^2 \phi_v(\omega)} = \frac{1}{|H(e^{j\omega})|^2}, \quad (28)$$

$$\xi_{\text{nr}}(H) \triangleq \frac{\int_{-\pi}^{\pi} \phi_v(\omega) d\omega}{\int_{-\pi}^{\pi} |H(e^{j\omega})|^2 \phi_v(\omega) d\omega} = \frac{\mathcal{P}_v}{\mathcal{P}_v(H)}. \quad (29)$$

The larger the value of the noise-reduction factor, the more the noise is reduced. This factor should be lower bounded by 1.

The gains add distortion to the desired signal. In order to evaluate the amount of distortion to the signal of interest, we define the subband and fullband speech-distortion indices [1], [2] as, respectively,

$$v_{\text{sd}}[H(e^{j\omega})] \triangleq \frac{E[|H(e^{j\omega})X(e^{j\omega}) - X(e^{j\omega})|^2]}{\phi_x(\omega)} = |1 - H(e^{j\omega})|^2, \quad (30)$$

$$v_{\text{sd}}(H) \triangleq \frac{\int_{-\pi}^{\pi} E[|H(e^{j\omega})X(e^{j\omega}) - X(e^{j\omega})|^2] d\omega}{\int_{-\pi}^{\pi} \phi_x(\omega) d\omega} = \frac{\int_{-\pi}^{\pi} \phi_x(\omega) |1 - H(e^{j\omega})|^2 d\omega}{\int_{-\pi}^{\pi} \phi_x(\omega) d\omega} = \frac{\int_{-\pi}^{\pi} \mathcal{P}_x(\omega) |1 - H(e^{j\omega})|^2 d\omega}{2\pi \mathcal{P}_x}. \quad (31)$$

The speech-distortion index is lower bounded by 0 and expected to be upper bounded by 1 for optimal gains. The higher the value of this index, the more the desired signal is distorted.

3. WIENER FILTER

In this section, we derive the Wiener filter and explain its relationship with AR processes.

We define the error signal between the estimated and desired signals at frequency ω as

$$\mathcal{E}(e^{j\omega}) \triangleq Z(e^{j\omega}) - X(e^{j\omega}) = H(e^{j\omega})Y(e^{j\omega}) - X(e^{j\omega}). \quad (32)$$

This error can also be put into the form:

$$\mathcal{E}(e^{j\omega}) = \mathcal{E}_x(e^{j\omega}) + \mathcal{E}_v(e^{j\omega}), \quad (33)$$

where

$$\mathcal{E}_x(e^{j\omega}) \triangleq [H(e^{j\omega}) - 1] X(e^{j\omega}) \quad (34)$$

is the speech distortion due to the complex gain, and

$$\mathcal{E}_v(e^{j\omega}) \triangleq H(e^{j\omega})V(e^{j\omega}) \quad (35)$$

represents the residual noise.

The frequency-domain (or subband) mean-squared error (MSE) is then

$$J[H(e^{j\omega})] = E[|\mathcal{E}(e^{j\omega})|^2]. \quad (36)$$

Minimizing the subband MSE with respect to $H(e^{j\omega})$, we easily find the Wiener gain:

$$H_{\text{W}}(e^{j\omega}) = \frac{\phi_x(\omega)}{\phi_y(\omega)} = 1 - \frac{\phi_v(\omega)}{\phi_y(\omega)} = \frac{i\text{SNR}(\omega)}{1 + i\text{SNR}(\omega)}. \quad (37)$$

With this filter, we have the following property.

Property: With the optimal noncausal Wiener filter given in (37), the fullband SNR gain is always greater than or equal to 1, i.e., $\mathcal{G}(H_{\text{W}}) \geq 1$.

The proof of this can be found in [4]. This fundamental property shows that the optimal gain can never amplify the noise. Nevertheless, it is essential to understand when the Wiener filter improves the SNR when it can.

It is more informative to rewrite the optimal gain (37) as a function of the prediction gains:

$$H_{\text{W}}(e^{j\omega}) = \frac{i\text{SNR}_0 \cdot \mathcal{P}_x(\omega)}{\mathcal{P}_v(\omega) + i\text{SNR}_0 \cdot \mathcal{P}_x(\omega)}. \quad (38)$$

From the definition of the fullband SNR gain, we observe that the worst-case scenario for the Wiener filter is when $\mathcal{P}_x(\omega) = \mathcal{P}_v(\omega)$, $\forall \omega$. Indeed, in this situation

$$H_{\text{W}}(e^{j\omega}) = \frac{i\text{SNR}_0}{1 + i\text{SNR}_0} = \text{Constant}, \quad \forall \omega \quad (39)$$

and it is easy to verify that $\mathcal{G}(H_{\text{W}}) = 1$, so the SNR cannot be improved. The case $\mathcal{P}_x(\omega) = \mathcal{P}_v(\omega)$, $\forall \omega$, happens when either the speech and noise signals are completely unpredictable (white random signals) so that $\mathcal{P}_x = \mathcal{P}_v = 1$ [and $\mathcal{P}_x(\omega) = \mathcal{P}_v(\omega) = 1$, $\forall \omega$] or the speech and noise signals are basically the same signals (having the same AR parameters) but they may have different volumes. It is now easy to understand why with babbling noise, for example, the Wiener filter may not perform at its best, since the speech and noise signals may have similar spectra [3]. This simple analysis of the Wiener filter is very insightful; it explains some practical situations where this filter may have limited performances.

To further analyze (38), we assume that $i\text{SNR}_0 = 1$ (for simplicity) and the desired and noise signals are neither white nor they have the same AR parameters. In this case

$$H_{\text{W}}(e^{j\omega}) = \frac{\mathcal{P}_x(\omega)}{\mathcal{P}_v(\omega) + \mathcal{P}_x(\omega)}. \quad (40)$$

If the desired signal is more predictable than the noise signal at all frequencies [i.e., $\mathcal{P}_x(\omega) > \mathcal{P}_v(\omega)$, $\forall \omega$], we have $H_{\text{W}}(e^{j\omega}) > 0.5$, $\forall \omega$. But if the noise signal is more predictable than the desired signal at all frequencies [i.e., $\mathcal{P}_v(\omega) > \mathcal{P}_x(\omega)$, $\forall \omega$] then $H_{\text{W}}(e^{j\omega}) < 0.5$, $\forall \omega$. We deduce, by looking at (27), that the fullband SNR gain is greater in the second case than in the first one but at a price of more speech distortion by inspection of (31). These two cases have opposite behaviors. These results may come as a small surprise. The fact that the noise is predictable may not be, after all, a bad thing as long as its spectrum is different from the speech spectrum. It is often intuitively thought that the non-whiteness of the noise may limit the performance of the Wiener filter but we understand now that this limitation that may occur in practical situations may mainly be due to the fact that noise and speech may have similar spectra.

Let us further exploit (38) by rewriting it as follows:

$$H_{\text{W}}^2(e^{j\omega}) [\mathcal{P}_v(\omega) + i\text{SNR}_0 \cdot \mathcal{P}_x(\omega)] = i\text{SNR}_0 \cdot \mathcal{P}_x(\omega) H_{\text{W}}(e^{j\omega}). \quad (41)$$

Integrating both sides of (41) and multiplying by $1/(2\pi)$, we obtain

$$\mathcal{P}_v(H_{\text{W}}) + i\text{SNR}_0 \cdot \mathcal{P}_x(H_{\text{W}}) = \frac{i\text{SNR}_0}{2\pi} \int_{-\pi}^{\pi} H_{\text{W}}(e^{j\omega}) \mathcal{P}_x(\omega) d\omega \leq i\text{SNR}_0 \cdot \mathcal{P}_x. \quad (42)$$

Therefore, we deduce the bounds for the fullband prediction gain of the filtered desired signal with the Wiener filter:

$$0 \leq \mathcal{P}_x(H_{\text{W}}) \leq \mathcal{P}_x - \frac{\mathcal{P}_v(H_{\text{W}})}{i\text{SNR}_0} \leq \mathcal{P}_x, \quad (43)$$

or

$$0 \leq \mathcal{P}_x(H_{\text{W}}) \leq \mathcal{P}_x \frac{o\text{SNR}(H_{\text{W}})}{1 + o\text{SNR}(H_{\text{W}})} \leq \mathcal{P}_x, \quad (44)$$

