

Some Robust Speech Enhancement Techniques using Higher Order AR Estimation

Josep M.SALAVEDRA*, Enrique MASGRAU**, Asunción MORENO*, Joan ESTARELLAS*

* Department of Signal Theory and Communications. Universitat Politècnica de Catalunya. Apartat. 30002. 08080- BARCELONA. SPAIN. Tel/Fax: +34-3- 4017404 / 4016447 . E-mail: mia@tsc.upc.es

** Department of Electrical Engineering and Computers. Universidad de Zaragoza. María de Luna, 3. 50015-ZARAGOZA. SPAIN

Abstract. We study some speech enhancement algorithms based on the iterative Wiener filtering method due to Lim-Oppenheim [2], where the AR spectral estimation of the speech is carried out using a 2nd-order analysis. But in our algorithms we consider an AR estimation by means of cumulant analysis. This work extends some preceding papers due to the authors, where information of previous speech frames is taken to initiate speech AR modelling of the current frame. Two parameters are introduced to design Wiener filter at first iteration of this iterative algorithm. These parameters are the Interframe Factor IF and the Previous Frame Iteration PFI. A detailed study of them shows they allow a very important noise suppression after processing only first iteration of this algorithm, without any appreciable increase of distortion. Two different ways to combine current and previous frame AR modelling are evaluated.

1. Introduction

It is well known, that many applications of speech processing that show very high performance in laboratory conditions degrade dramatically when working in real environments because of low robustness. The solution we propose here concerns to a preprocessing front-end in order to enhance the speech quality by means of a speech parametric modelling insensitive to the noise. The use of HO cumulants for speech AR modelling calculation provides the desirable uncoupling between noise and speech. It is based on the property that for Gaussian processes only, all cumulants of order greater than two are identically zero [1]. Moreover, the non-Gaussian processes presenting a symmetric p.d.f. have null odd-order cumulants. Considering a Gaussian or a symmetric p.d.f. noise (a good approximation of very real environments) and the non-Gaussian characteristic of the speech (principally for the voiced frames) it would be possible to obtain an spectral AR modelling of the speech more independent of the noise by using, e.g., 3rd-order cumulants of noisy speech instead of common 2nd-order cumulant.

2. Iterative Wiener Algorithm

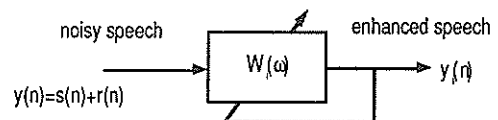
In the original Lim-Oppenheim Method [2], noisy speech is enhanced by means of an iterative Wiener filtering that is defined as:

$$W(\omega) = \frac{P_s}{P_s + P_r} \quad (1)$$

where P_r is the spectrum of the noise signal $r(n)$, estimated in non-speech frames, and P_s is a spectrum estimation of the unavailable clean speech signal. So, both speech and noise spectra estimation must be available to design the Wiener filter at every frame. We talk over signal estimation because both signals are not available and only noisy speech signal can be processed. An iterative Wiener filtering is used to obtain a better estimation of the AR speech modelling (figure 1.). At first sight an improvement of performance can be expected after every iteration since this current AR speech estimation is carried out from a cleaner speech signal than filter estimation

This work was supported by TIC 92-0800-C05-04

of the preceding iteration. But other factors sidetrack this iterative algorithm and a limitation in the number of iterations must be taken in account. Clearly the filtered speech signal contains a smaller residual noise but it presents a larger spectral distortion. Therefore, increasing the number of iterations doesn't always involve a better speech estimation. It is well known that this algorithm leads to a narrowness and a shifting of the speech formants [3], providing an unnatural sounding speech. In [4] a detailed convergence analysis of this algorithm is carried out. It is proved that this estimated Wiener filter tends to cancel all signal frequencies with SNR lower than 4.77dB, and an additional attenuation, proportionally to the noise level, affects signal frequencies with higher SNR, in comparison to the optimum Wiener filter. Only the non-contaminated speech frequencies undergo a null attenuation.



$$W_i(\omega) = \frac{P_{y_i}(\omega)}{P_{y_i}(\omega) + P_r(\omega)} \quad \text{where} \quad P_{y_i}(\omega) = \frac{g^2}{\left| 1 + \sum_{k=1}^p a_k e^{-jk\omega} \right|^2}$$

Figure 1. Scheme of the iterative Wiener algorithm.

3. The Parameterized Algorithm

A parameterized Wiener filtering has been considered to have a better control over noise suppression, intelligibility loss and computational complexity, by adding two parameters ∂ and β in the Wiener filter computation (1). So, we consider now the following equation:

$$W_i(\omega) = \left(\frac{P_y}{P_y + \beta \cdot P_r} \right)^\partial \quad (2)$$

By varying these parameters ∂ and β , filters with different characteristics can be obtained. Thus, if $\partial = \beta = 1.0$ then expression (2) corresponds to the general Wiener filter equation

(1), and if $\delta=0.5$, $\beta=1.0$ it corresponds to power spectrum filtering. In [6], a detailed study of performance was performed. High values of both parameters lead to a more aggressive Wiener filter and so noise suppression is increased but distortion increases too. We found that $\delta=1.0$, $\beta=1.2$ is a good trade-off among noise suppression, distortion and convergence speed of the iterative filtering, when 3rd-order statistics and low SNR are considered.

AR modelling (figure 1.) of the speech spectrum estimation is computed from 3rd-order cumulants that are calculated using the covariance case:

$$C_k(i,j) = \sum_{n=p+1}^N x(n-k).x(n-i).x(n-j) \quad , \quad 0 \leq k,i,j \leq p \quad (3)$$

where $p=10$ is the order of the filter. Then speech AR modelling coefficients a_k are computed by solving the following equations [1]:

$$\sum_{k=0}^p a_k \cdot C_k(i,j) = 0 \quad , \quad 1 \leq i \leq p \quad ; \quad 0 \leq j \leq i \quad (4)$$

As discussed in preceding works due to the authors [5,6], we obtain a twofold benefit by considering this 3rd-order AR modelling: Firstly, an accelerated convergence of the iterative algorithm and so a reduction of both computational complexity and intelligibility loss; Secondly, achievement of a non polluted AR speech parameterization. In comparison to 2nd-order statistics estimation we obtain a good improvement but the price we pay for these advantages is a higher distortion. Thus a higher "peaking" or "narrowness" effect of the speech formants is brought about [4].

When additive noise is AWGN (SNR=0dB) the improvement over second-order algorithm is very appreciated for any number of iterations (see Table 1). While the improvement of second-order approach increases gradually, but slowly, iteration by iteration, 3rd-order one gets a very good improvement, about 3dB, after only two iterations and thus it obtains a faster convergence. Furthermore, in comparison to 4th-order algorithm, third-order one also obtains better results and its computational complexity is much lower. Therefore, 3rd-order cumulants lead to a faster noise reduction because of its higher aggressiveness with respect to both 4th-order cumulants and autocorrelation function [5].

4. The Interframe Factor IF

In table 1, we may appreciate an improvement that increases gradually iteration by iteration. Most part of noise

suppression is obtained after processing two iterations. Third-order cumulants obtain an appreciable noise suppression (about 2dB in Cepstrum distance) after first iteration (see Table 1.b) and then this speech modelling is enhanced a lot in the second iteration because it estimates Wiener filter from cleaner speech signal. At first iteration, speech AR modelling is computed from noisy signal without any initial information about the features of speech signal corresponding to the current frame. However, we know some information of the current speech frame by considering that speech signal features don't vary a lot between two consecutive overlapped frames. Therefore, we propose to obtain the first iteration AR coefficients as a combination between current frame AR estimation and previous frame AR coefficients. Thus, we design the non-causal Wiener filter, corresponding to first iteration, as a linear combination of coefficients a_k , belonging to two consecutive frames, calculated as follows:

$$A_k(n,1) = IF \cdot a_k(n,1) + (1 - IF) \cdot A_k(n-1,PFI) \quad (5)$$

$$0 \leq k \leq P \quad ; \quad 1 \leq PFI \leq MAXITER \quad ; \quad 0 \leq IF \leq 1$$

where n is the current frame, PFI is the Previous Frame Iteration that we consider to help first iteration of the current frame and IF is the Interframe Factor. We write a_k when coefficients are estimated directly from a noisy speech frame and we note capital letter A_k when coefficients are coming from a linear combination of a_k . At the beginning of every speech activity we set parameter IF=1 because information of last speech frame is not related to the current speech frame. Wiener filter designs corresponding to the remaining iterations of the algorithm are estimated over a cleaner speech signal coming from Wiener filtering Output of previous iteration of the same frame:

$$A_k(n,iter) = a_k(n,iter) \quad , \quad 2 \leq iter \leq MAXITER \quad (6)$$

where iter is the iteration number of the current frame. We have two parameters to control this linear combination. First parameter is the Interframe Factor IF that represents the amount of current speech AR estimation $a_k(n,1)$ we put in the AR modelling $A_k(n,1)$ of the filter. The interframe factor is the main parameter to control linear combination (5) because parameter IF=1 represents that only current AR estimation is considered to design Wiener filter at first iteration of current frame and then parameter PFI has no sense to be considered. Thus, parameter IF=1 refers to a situation where no interframe factor is defined. If we decide to consider previous frame information (IF<1) we must consider parameter PFI to answer the following question : Which iteration number (PFI) of preceding frame must we take to obtain a reliable speech AR modelling? Preceding works [5,6] have shown that it has no sense to process more than 5 iterations while third-order statistics are considered. Therefore, parameter MAXITER=5 has been fixed in all our tests.

| a) | SNR | SEGSN | ITAKU | COSH | CEPST |
|---------|------|-------|-------|-------|-------|
| 0 iter. | 0.00 | 0.79 | 9.57 | 11.67 | 12.02 |
| 1 iter. | 7.36 | 4.38 | 9.21 | 10.71 | 11.01 |
| 2 iter. | 8.83 | 5.92 | 8.86 | 10.17 | 9.90 |
| 3 iter. | 9.04 | 6.16 | 7.30 | 9.04 | 9.34 |
| 4 iter. | 9.11 | 6.25 | 6.42 | 8.45 | 9.20 |

| b) | SNR | SEGSN | ITAKU | COSH | CEPST |
|---------|------|-------|-------|-------|-------|
| 0 iter. | 0.00 | 0.79 | 9.57 | 11.67 | 12.02 |
| 1 iter. | 7.92 | 4.86 | 8.18 | 9.78 | 9.82 |
| 2 iter. | 7.60 | 5.31 | 5.94 | 8.16 | 8.47 |
| 3 iter. | 7.59 | 5.59 | 5.11 | 7.55 | 8.15 |
| 4 iter. | 7.36 | 5.79 | 5.15 | 7.64 | 8.30 |

| c) | SNR | SEGSN | ITAKU | COSH | CEPST |
|---------|------|-------|-------|-------|-------|
| 0 iter. | 0.00 | 0.79 | 9.57 | 11.67 | 12.02 |
| 1 iter. | 8.31 | 5.17 | 5.00 | 7.82 | 7.90 |
| 2 iter. | 7.90 | 5.56 | 5.03 | 7.54 | 8.08 |
| 3 iter. | 7.91 | 5.88 | 5.01 | 7.43 | 7.98 |
| 4 iter. | 7.68 | 5.78 | 4.86 | 7.45 | 8.16 |

| d) | SNR | SEGSN | ITAKU | COSH | CEPST |
|---------|------|-------|-------|-------|-------|
| 0 iter. | 0.00 | 0.79 | 9.57 | 11.67 | 12.02 |
| 1 iter. | 7.47 | 4.53 | 8.97 | 10.49 | 10.53 |
| 2 iter. | 7.39 | 4.95 | 7.88 | 9.65 | 9.30 |
| 3 iter. | 7.37 | 5.11 | 6.55 | 8.65 | 8.80 |
| 4 iter. | 7.77 | 5.49 | 5.52 | 7.91 | 8.47 |

Table 1. Distance measures using algorithms based on: a) second order statistics; b) parameterized third order cumulants; c) parameterized third order with interframe factor IF=0.7, considering parameter PFI=3; d) fourth order cumulants at SNR=0dB.

When parameter $PFI=1$ and parameter $IF < 1$, We may evaluate in (5) that current AR coefficients $A_k(n,1)$ to design Wiener Filter are coming from a linear combination of AR estimations belonging to previous frames:

$$A_k(n,1) = IF \sum_{r=2}^n (1-IF)^{n-r} \cdot a_k(r,1) + (1-IF)^{n-1} \cdot a_k(1,1) \quad (7)$$

$0 \leq k \leq P$

where n is the distance in number of frames since last non-speech activity frame appeared. We have a combination of coefficients a_k coming from a lot of different frames and therefore the estimation window may contain some phonemes and it doesn't fulfil our initial assumption. The worst of it is that parameter $IF=0$ means AR estimation of first frame after non-speech activity is assigned to all the frames belonging to the same speech activity. So, distortion effect may increase a lot, specially when parameter IF takes small values.

To avoid this problem, another linear combination to design Wiener Filter is proposed:

$$A_k(n,1) = IF \cdot a_k(n,1) + (1-IF) \cdot a_k(n-1,PFI) \quad (8)$$

$0 \leq k \leq P; 1 \leq PFI \leq MAXITER; 0 \leq IF \leq 1$

This estimation (8),(6) is referred as Method B while previous one (5),(6) is referred as Method A in Figure 2. A comparison of performance of both methods is shown. Cepstrum distance after first iteration of the parameterized third order algorithm is depicted with a value of parameter $PFI=1$. It can be appreciated that Method B give better estimations when parameter IF is lower than 0.5 and similar performance is achieved when the information coming from previous frame is less significant.

On the other hand, parameter $IF=0$ represents that the coming noisy speech frame is filtered by means of a filter

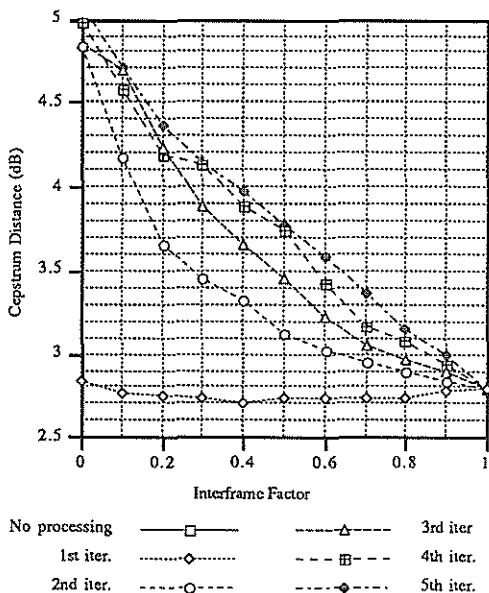


Figure 3. Distortion Effect produced by first iteration processing when some different speech AR estimations belonging to different iterations of previous speech frame are considered.

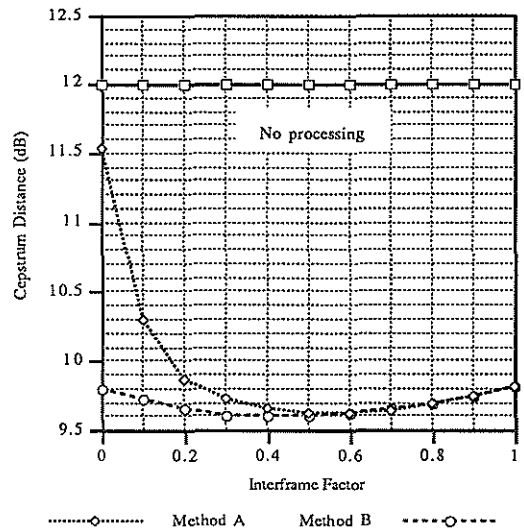


Figure 2. Comparison of two possible methods to combine AR coefficients coming from current and previous speech frames when first iteration is processed and $PFI=1$.

estimation coming from preceding speech frame. Two different situations may be distinguished: $PFI=1$ and $PFI > 1$. When information proceeding from first iteration of previous speech frame ($PFI=1$) is considered, no better results than before ($IF=1$) are expected, because the speech AR estimator is looking at the same noisy speech signal, but in the previous frame, and performance therefore decreases when parameter IF decreases to zero (see 1st iteration line in figure 4.). However, a good improvement (about 1.5dB in Cepstrum distance) is obtained

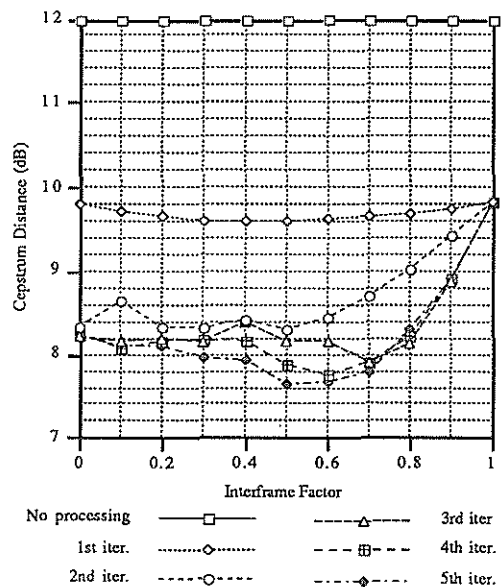


Figure 4. Noise Suppression after processing first iteration of current frame when some different speech AR estimations belonging to different iterations of previous speech frame are considered.

when parameter $PFI > 1$ but distortion effect increases more than 2dB in Cepstrum distance (see figure 3.) because current Wiener filter is designed with speech AR estimation proceeding from the preceding frame over a cleaner speech signal.

In figure 3., Cepstrum distance corresponding to first iteration of current frame has been represented and some different iteration numbers of preceding frame have been evaluated. Clean speech has been processed by this system and so distortion effect corresponding to the iterative algorithm has been depicted. To avoid an appreciable increase of distortion effect all values of parameter IF lower than 0.6 must be discarded. In figure 4., first iteration of current frame corresponding to speech signal disturbed by AWGN at SNR=0dB has been processed and some different speech AR estimations of previous frame have been evaluated (ranging PFI from 1 to 5). We may come to the conclusion that values of parameter IF ranging from 0.6 to 0.8 represent a good tradeoff between distortion and noise suppression. Therefore, we may achieve an improvement of 2dB in Cepstrum distance by introducing parameter IF ($PFI=3$ and $IF=0.7$) to estimate current speech AR modelling without any noticeable increase of distortion (0.25 dB). Thus, we may obtain an improvement higher than 4 dB in Cepstrum distance after processing only first iteration of the iterative Wiener filtering.

As it is depicted in Figure 5. values of parameter $PFI > 2$ give a similar performance. After second iteration most part of linear combinations leads to similar levels of Cepstrum Distance but, in listening tests, it may be appreciated a less distortion effect when parameter $PFI > 2$ and furthermore the best performance is achieved after 3 iterations of Lim-Oppenheim algorithm while 4 iterations are necessary when parameter $IF=1$, to obtain a similar quality. Therefore value $PFI=3$ may be considered as a good trade-off among computational complexity, distortion effect and noise suppression. This fact

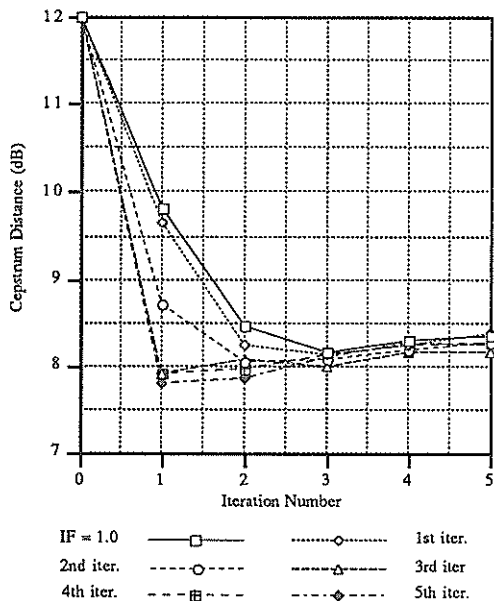


Figure 5: Noise Suppression by using parameter $IF=0.7$ and some different speech AR estimations corresponding to different iterations of previous speech frame.

may be justified looking at Itakura Distance where a very important reduction (about 4.5dB) with only first iteration is achieved and therefore formants estimation in voiced sounds is clearly improved by introducing these two parameters at first iteration of every frame. Obviously this linear combination of coefficients a_k tends to improve quality inside of voiced sounds. Some constraints have been added to the algorithm to protect unvoiced frames against this linear combination and so parameter $IF=1$ is set inside of unvoiced frames and first voiced frame. Similar performance is achieved when diesel engine noise is considered, although differences between AR weighting method ($IF < 1$) and no interframe weighting method ($IF=1$) are smaller.

5. Conclusions

A speech enhancement method based on an iterative Wiener filtering have been proposed. Spectral estimation of speech is obtained by means of an AR modelling based on third-order cumulant analysis to provide the desirable noise-speech uncoupling. Two parameters, IF (Interframe Factor) and PFI (Previous Frame Iteration), have been considered to take advantage of previous speech spectrum estimations to initiate AR modelling corresponding to first iteration of the current speech frame. This approach achieves a noise suppression about 4dB (Cepstrum Distance) after processing only first iteration of the algorithm. This fact represents an improvement about 2dB (Cepstrum Distance) in relation to parameterized third-order algorithm ($IF=1$). Finally, the convergence of the iterative algorithms based on cumulant AR estimation is strongly accelerated. Therefore, a good reduction of computational complexity and processing delay are achieved, while no appreciable increase of distortion effect is generated. All these features are specially esteemed when low and medium SNR are considered.

References

- [1] C.L.Nikias, M.R.Raghuveer, "Bispectrum Estimation: A Digital Signal Processing Framework". Proc. of IEEE, pp 869-891. July 1987.
- [2] J.S.Lim and A.V.Oppenheim, "All-Pole Modeling of Degraded Speech". IEEE Trans ASSP, pp197-210. June 1978.
- [3] J.H.L.Hansen, M.A.Clements, "Constrained Iterative Speech Enhancement with Applications to Speech Recognition". IEEE Trans on Sig. Proc., pp 795-805. April 1991.
- [4] E.Masgrau, J.M.Salavedra, A.Moreno, A.Ardanuy, "Speech Enhancement by Adaptive Wiener Filtering based on Cumulant AR Modelling". Proc. ESCA Workshop on Speech Processing in Adverse Conditions, pp 143-146. Cannes, France. November 92.
- [5] J.M.Salavedra, E.Masgrau, A.Moreno, X.Jové, "A Speech Enhancement System using Higher-order AR estimation in real environments". Proc. EUROSPEECH'93, pp. 223-226. Berlin, Germany. September 21-23, 1993.
- [6] J.M.Salavedra, E.Masgrau, A.Moreno, X.Jové and J.Estarellas, "Robust Coefficients of a Higher-order AR Modelling in a Speech Enhancement System using parameterized Wiener Filtering". Proc. MELECON'94, pp. 69-72. Antalya, Turkey. April 12-14, 1994.