

## 1 Introduction

- Linear Prediction suffers from well known problems when the 2-norm error minimization criterion is employed in the analysis and coding of voiced speech.
- The usual approach is to find coefficients for the short-term and long-term signal correlation in two different steps, leading to inherently suboptimal solutions.
- In this work we define a joint estimation approach based on the observation of the behavior of the short-term and long-term cascade polynomial.
- Imposing sparsity on a high order predictor, we obtain a polynomial that can be easily factorized into long-term and short-term predictors.
- This method incorporated into an ACELP scheme shows to have better performance than traditional cascade methods and other joint estimation methods.

## 2 Joint Estimator

- In order to remove near-sample redundancies and distant-sample redundancies, a cascade of a short-term linear predictor  $F(z)$  and a long-term linear predictor  $P(z)$  is employed.
- The cascade of the two predictors corresponds the multiplication in the  $z$ -domain of the two transfer functions:

$$\begin{aligned} A(z) &= F(z)P(z) = 1 - \sum_{k=1}^K a_k z^{-k} \\ &= \left(1 - \sum_{k=1}^{N_f} f_k z^{-k}\right) \left(1 - \sum_{k=1}^{N_p} g_k z^{-(T_p+k)}\right). \end{aligned}$$

- $A(z)$  will therefore be highly sparse. Sparsity is then taken into account in new error minimization criterion:

$$\min_{\mathbf{a}} \|\mathbf{x} - \mathbf{X}\mathbf{a}\|_2^2 + \gamma \|\mathbf{a}\|_1,$$

where the 1-norm is employed as a relaxation of the non-convex 0-norm and:

$$\mathbf{x} = \begin{bmatrix} x(N_1) \\ \vdots \\ x(N_2) \end{bmatrix}, \mathbf{X} = \begin{bmatrix} x(N_1-1) \cdots x(N_1-N_f) \\ \vdots \\ x(N_2-1) \cdots x(N_2-N_f) \end{bmatrix}$$

- $A(z)$  is now similar to the multiplication between a short-term and a long-term predictor:

$$A(z) \approx A_{LTP}(z)A_{stp}(z).$$

- The first  $N_{stp}$  coefficients are used as the estimated coefficients of the short-term predictor  $A_{stp}(z)$ .
- $A_{LTP}(z)$  is created by taking the quotient of the division between  $A(z)$  by  $A_{stp}(z)$ . The minimum value and its position will correspond to our estimate of the pitch gain and delay (parameters of the predictor  $P(z)$ ):

$$g_p = \min\{a_{LTP}\},$$

$$T_p = \arg \min\{a_{LTP}\}.$$

where  $\{a_{LTP}\}$  are the coefficients of  $A_{LTP}(z)$ . An example is shown Figure 1.

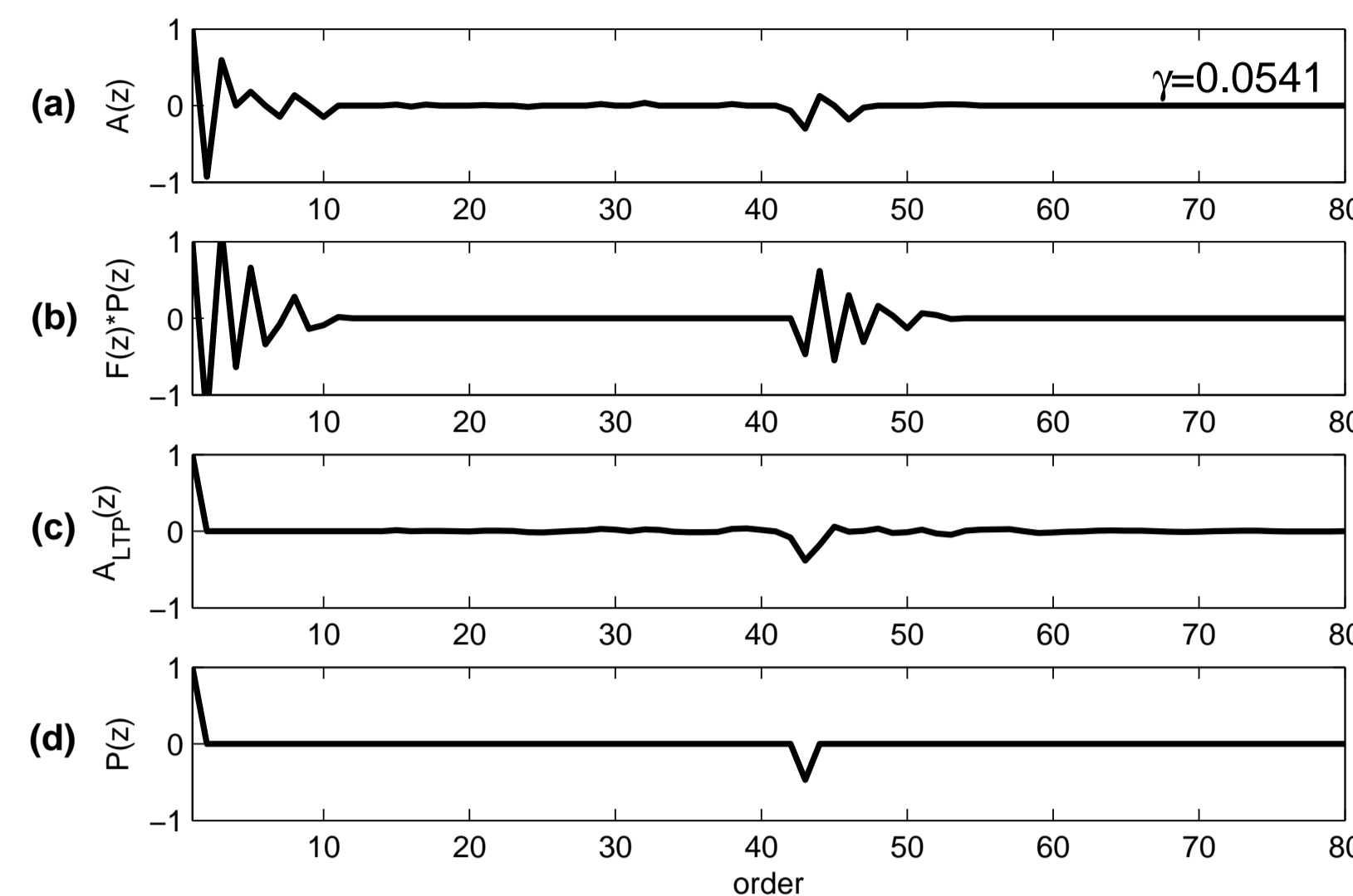


Figure 1: (a) and (b) show a comparison between the polynomial obtained with regularized minimization  $A(z)$  and multiplication of the two predictors  $F(z)P(z)$  obtained in cascade; (c) and (d) a comparison of the two long-term predictors  $A_{LTP}(z)$  and  $P(z)$ .

## 3 Regularization Parameter

- The regularization parameter  $\gamma$  is intimately related to the *a priori* knowledge that we have on the coefficients vector  $\{a_k\}$  (how sparse  $\{a_k\}$  is) considering our minimization criterion from a Bayesian point of view.
- The best trade-off between the 2-norm of the residual and the 1-norm of the solution vector is found finding the point of maximum curvature of the curve ( $\|\mathbf{x} - \mathbf{X}\mathbf{a}_\gamma\|_2, \|\mathbf{a}_\gamma\|_1$ ) (*L*-curve).
- $\gamma$  is bounded ( $0 < \gamma < \|\mathbf{X}^T \mathbf{x}\|_\infty$ ).
- We investigate three approaches for the selection of  $\gamma$  based on the magnitude of the difference between the encoded-decoded signal and the original signal:
  - constant ( $\mathbf{R}_c$ )**. The value that on average gave the best result.
  - adaptive ( $\mathbf{R}_a$ )**. The value of  $\gamma$  intimately related to the pitch gain  $g_p$ . We update  $\gamma$  using the following approximate relation:
 
$$\gamma(n+1) = -0.18g_p^2(n) + 0.2.$$
  - optimal ( $\mathbf{R}_o$ )**.  $\gamma$  is tuned for every frame analyzed in order to obtain the best result.
- Selection of  $\gamma$  is based on the magnitude of the difference between the encoded-decoded signal and the original signal.

## 4 Validation

- The joint method is implemented in an ACELP scheme.
- The order of the optimization problem is  $K = 110$  and the frame length is  $N = 160$  (20 ms). The order of the short-term and long-term predictors are respectively  $N_{stp} = 12$  and  $N_{LTP} = 1$ .
- Using  $K = 110$  we can cover pitch frequencies in the interval [82 Hz, 571 Hz].
- Residual vector is encoded using 40 non-zero samples constrained with  $\pm 1$  values and a gain (Algebraic Codebook).
- For each method ( $\mathbf{R}_c$ ,  $\mathbf{R}_a$ ,  $\mathbf{R}_o$ ,  $\mathbf{A}_j$ ), the signals coming out of the encoding-decoding scheme are compared to the original speech and the traditional ACELP  $\mathbf{A}_c$ , PESQ evaluation is then performed.

METHOD	$\Delta$ DIST	$\Delta$ MOS
$\mathbf{R}_o$	$2.05 \pm 0.06$ dB	$0.11 \pm 0.00$
$\mathbf{R}_a$	$1.65 \pm 0.11$ dB	$0.07 \pm 0.00$
$\mathbf{R}_c$	$1.04 \pm 0.27$ dB	$0.03 \pm 0.03$
$\mathbf{A}_j$	$0.32 \pm 0.13$ dB	$0.00 \pm 0.02$

Improvements over conventional ACELP  $\mathbf{A}_c$  in the decoded speech signal in terms of reduction of log magnitude distortion ( $\Delta$ DIST) and Mean Opinion Score ( $\Delta$ MOS). A 95% confidence intervals is given for each value.

## 5 Discussion

- The increase in accuracy is given by the more precise search of the algebraic codeword (spectrally white residual) and improved pitch tracking.
- Number of taps is highly customizable and can be chosen using an Analysis-by-Synthesis scheme or a Model Order Selection criterion.
- Lower emphasis on peaks is achieved by intrinsically taking into consideration the possible outliers due to the pitch excitation in the minimization process. This reflects in a lower sensitivity of the short-term predictor to quantization than traditional LP.
- The cascade  $A_{stp}(z)P(z)$  has a very low instability rate (less than 0.01%).
- The optimization problem can be posed as a quadratic programming problem and solved efficiently using an interior-point algorithm.

## 6 Conclusion

- A new formulation for the minimization process involved in the linear prediction has been presented.
- We have obtained a better statistical fitting for the model of speech that makes analysis and coding more straightforward and accurate.
- Higher accuracy than with traditional LP have been obtained due to whiter residual, improved pitch tracking and predictors that are less sensitive to quantization.

## References

- [1] J. Makhoul, "Linear Prediction: A Tutorial Review", *Proc. IEEE*, vol. 63(4), pp. 561–580, April 1975.
- [2] P. Kabal and R. P. Ramachandran, "Joint Optimization of Linear Predictors in Speech Coders", *Proc. IEEE*, vol. 37(5), pp. 642–650, May 1989.
- [3] H. Zarrinkoub and P. Mermelstein, "Joint Optimization of Short-Term and Long-Term Predictors in CELP Speech Coders", *Proc. IEEE Int. Conf. Acoust., Speech, Signal Processing*, vvol. 2, pp. 157–160, 2003.
- [4] D. Giacobello, M. G. Christensen, J. Dahl, S. H. Jensen and M. Moonen, "Sparse Linear Predictors for Speech Processing", *Proc. INTERSPEECH*, September 2008.
- [5] P. C. Hansen, "Analysis of Discrete Ill-Posed Problems by Means of the L-Curve", *Proc. SIAM Review*, vol. 34, no. 4, pp. 561–580, December 1992.
- [6] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [7] S. J. Wright, *Primal-Dual Interior-Point Methods*, SIAM, 1997.