

Audio signal segmentation using recursive Bayesian change-point detectors

ROMAN CMEJLA, PAVEL SOVKA
Department Circuit Theory
Czech Technical University in Prague
Technicka 2, 166 27 Prague 6
CZECH REPUBLIC

Abstract: - Two basic types of change-point detectors are used for a localization of abrupt changes in audio signals, especially in speech and music signals. The Bayesian autoregressive change-point detector (BACD) and Bayesian polynomial change-point detector (BPCD) are analyzed and modified to enable sequential signal segmentation. The modification consists in a recursive evaluation of functions used in these detectors. Normalization based on Bayesian evidence is discussed together with a proper choice of parameters ensuring the numerical stability. Suggested detectors seem to be computationally effective and numerical stable as shown by experiments. Some illustrative examples of the segmentation of real music and speech signals are given.

Key-Words: - Bayesian change-point detector, Bayesian evidence, sequential signal segmentation, recursive algorithm

1 Introduction

The signal segmentation has been challenge for many years. A great number of methods and applications can be found, e.g. a speech, music and biological signal segmentation. The segmentation is based on searching change-points detection using suitable signal parameters. Very robust and reliable methods are based the maximum likelihood and Bayesian approach e.g. [1], [2], [3], [4]. Bayesian detectors are very effective because they remove nuisance parameters from the analysis by a marginalization process. Moreover, offer powerful tools for the model order selection. Effective way of the Bayesian detector implementation is based on recursions for a change-point position and a new data. Robust implementation of a recursive growing window algorithm for one change-point detection is suggested in [2], [5]. When sequential signal segmentation is required the problems with the occurrence of more change-points has to be solved. Problems connected with multiple change-points detection can be overcome by e.g. Markov Chain Monte Carlo Methods (MCMC) and their modifications e.g. [6], [7]. In spite of very high robustness of these methods some other solutions can be found (e.g. [1], [8], [9], [10], [11]). These methods offer relatively easy implementation of the sequential detection of abrupt changes, especially for audio signals. This paper makes an effort to use simple Bayesian change-point detectors for the sequential signal segmentation. Two basic types of Bayesian detectors are used and modified using sliding window algorithm to be suitable for the sequential signal segmentation. First detector is the BACD second the BPCD. The BACD is based on autoregressive modelling of signals the BPCD uses the polynomial signal modelling.

These two types of methods are chosen with respect to signal types. In order the segmentation to be successful a proper type of change-point detector should be used. As shown below the BACD seems to be suitable for speech, violin, oboe and clarinet signals while the BPCD seems to be better for trumpet and drum signals. Also parameters of segmentation methods should be optimised for given signal characteristics.

2 Problem Formulation

This section defines two types of Bayesian detectors the BACD and BPCD and describes their recursive implementation.

2.1 BACD

The BACD requires the signal model consisting of two parts: “left” generated by AR model with M_1 parameters a_k and “right” generated by another AR model with M_2 parameters b_k

$$d[n] = \begin{cases} \sum_{k=1}^{M_1} a_k \cdot d[n-k] + e[n], & n \leq m \\ \sum_{k=1}^{M_2} b_k \cdot d[n-k] + e[n], & n > m \end{cases} \quad n = 1, \mathbf{L}, N, \quad (1)$$

or in matrix form

$$\mathbf{d} = \mathbf{G}_A \cdot \mathbf{b}_A + \mathbf{e}. \quad (2)$$

The matrix \mathbf{G}_A has the Jordan form

$$\mathbf{G}_A = \begin{bmatrix} d[0] & d[-1] & \mathbf{L} & 0 & 0 & \mathbf{L} \\ d[1] & d[0] & \mathbf{L} & 0 & 0 & \mathbf{L} \\ d[2] & d[1] & \mathbf{L} & 0 & 0 & \mathbf{L} \\ \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{M} & \mathbf{M} & \mathbf{O} \\ d[m-1] & d[m-2] & \mathbf{L} & 0 & 0 & \mathbf{L} \\ 0 & 0 & \mathbf{L} & d[m] & d[m-1] & \mathbf{L} \\ \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{M} & \mathbf{M} & \mathbf{O} \\ 0 & 0 & \mathbf{L} & d[N-1] & d[N-2] & \mathbf{L} \end{bmatrix} \quad (3)$$

and it depends on the unknown index of change-point $m = 1, \dots, N$, which is determined by the maximum (MAP) of the posterior probability density function (pdf) given by [2]

$$p(m | \mathbf{d}, \mathbf{M}) \propto \frac{\left[(D - \mathbf{g}_A \Phi_A \mathbf{g}_A^T) \right]^{\frac{N-M_1-M_2}{2}}}{\sqrt{\Delta_A}}. \quad (4)$$

Matrix $\Phi_A = (\mathbf{G}_A^T \mathbf{G}_A)^{-1}$ is the inverse correlation matrix, $D = \mathbf{d}^T \mathbf{d}$ is the signal energy, $\mathbf{g}_A = \mathbf{d}^T \mathbf{G}_A$ is correlation vector, and $\Delta_A = \det(\mathbf{G}_A^T \mathbf{G}_A)$.

2.2 BPCD

The signal model for this type of detector is modelled by two different polynomials

$$d[n] = \begin{cases} \sum_{p=0}^{P_1} a_p \cdot t^p[n] + e[n], \\ \sum_{p=0}^{P_2} b_p \cdot t^p[n] + e[n], \end{cases} \quad n = 1, \mathbf{L}, N. \quad (5)$$

This equation can be written in matrix form

$$\mathbf{d} = \mathbf{G}_p \cdot \mathbf{b}_p + \mathbf{e}, \quad (6)$$

$$\mathbf{G}_p = \begin{bmatrix} t^0[1] & t^1[1] & \mathbf{K} & t^R[1] & 0 & 0 & \mathbf{K} & 0 \\ t^0[2] & t^1[2] & \mathbf{K} & t^R[2] & 0 & 0 & \mathbf{K} & 0 \\ t^0[3] & t^1[3] & \mathbf{K} & t^R[3] & 0 & 0 & \mathbf{K} & 0 \\ \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{M} & \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{M} \\ t^0[m] & t^1[m] & \mathbf{K} & t^R[m] & 0 & 0 & \mathbf{K} & 0 \\ 0 & 0 & \mathbf{K} & 0 & t^0[m+1] & t^1[m+1] & \mathbf{K} & t^R[m+1] \\ \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{M} & \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{M} \\ 0 & 0 & \mathbf{K} & 0 & t^0[N] & t^1[N] & \mathbf{K} & t^R[N] \end{bmatrix} \quad (7)$$

The probability density is now given by

$$p(m | \mathbf{d}, \mathbf{M}) \propto \frac{\left[(D - \mathbf{g}_p \Phi_p \mathbf{g}_p^T) \right]^{\frac{N-M_1-M_2}{2}}}{\sqrt{\Delta_p}}. \quad (8)$$

2.3 Recursive algorithm for BACD and BPCD

When both detectors are to be used for the change-point detection of non-stationary signals then the signal segmentation is necessary. The length of the window is dependent on signal parameters because both signal models (AR and polynomial model) presuppose one change in one signal segment. This requirement can lead to an inapplicable short window length. Moreover, results from different windows are not comparable. The first problem can be avoided by using sliding window algorithm for the evaluation of functions (4) and (8). The latter problem can be solved by the normalization of (4) and (8) using the Bayesian evidence. The need of this normalization follows from the fact that data vector is not constant (as supposed in the marginalization process used in [2]) when sliding window is applied. If the marginalization process is used for two parameters (not for three as in [2]) the resulting formula for Bayesian evidence

$$BE_I \propto \frac{\left[(D - \mathbf{g}_I \Phi_I \mathbf{g}_I^T) \right]^{\frac{N-M}{2}}}{\sqrt{\Delta_I}}, \quad M = M_1 \text{ or } M_2 \quad (9)$$

slightly differs from the formula given in [2]. Index I stands for A (in the case of the evidence for an autoregressive model and BACD) and P (in the case of the evidence for polynomial model and BPCD). The inverse of the correlations matrix $\Phi_I = (\mathbf{G}_I^T \mathbf{G}_I)^{-1}$ and cross-correlation vector $\mathbf{g}_I = \mathbf{d}_I^T \mathbf{G}_I$ are now derived from the data matrix \mathbf{G}_I given by

$$\mathbf{G}_I = \begin{bmatrix} d[1] & 0 & \mathbf{L} & \mathbf{L} & 0 & 0 \\ d[2] & d[1] & \mathbf{L} & \mathbf{L} & 0 & 0 \\ \mathbf{M} & \mathbf{M} & \mathbf{O} & \mathbf{L} & \mathbf{M} & \mathbf{M} \\ d[M] & d[M-1] & \mathbf{L} & \mathbf{O} & d[2] & d[1] \\ \mathbf{M} & \mathbf{M} & \mathbf{L} & \mathbf{L} & \mathbf{M} & \mathbf{M} \\ \mathbf{M} & \mathbf{M} & \mathbf{L} & \mathbf{L} & \mathbf{M} & \mathbf{M} \\ d[N-1] & d[N-2] & \mathbf{L} & \mathbf{L} & d[N-M+1] & d[N-M] \end{bmatrix} \quad (10)$$

The final normalized BACD (BACDN) or BPCD (BPCDN) detectors are given by

$$\tilde{p}(m | \mathbf{d}, \mathbf{M}) \propto \frac{\left[(D - \mathbf{g}_I \Phi_I \mathbf{g}_I^T) \right]^{\frac{N-M}{2}}}{\sqrt{\Delta_I}} \cdot BE_I^{-1}, \quad (11)$$

Using this function enables to compare results between different signal segments and thus it enables to form sliding window algorithm.

The recursive evaluation of pdf (4), (8) and (9) requires two types of the update for functions $D, \mathbf{g}_I, \Phi_I, \Delta_I$ (where I is A or P): the update of a change-point position m and the update given a new data [2], [5]. While the

position update is used here without any change, the data update is modified. Instead the growing memory algorithm used in [2] and [5] the sliding window algorithm is used.

This approach can be further simplified. Let $\tilde{p}_l(m|\mathbf{d}, \mathbf{M}), l=1, \dots, L, m=1, \dots, N$ represents the sliding window (length N) evaluation of the normalized posterior density (11) for L segments. Thus the evaluation of this time-dependent function requires updating $N \times L$ values. This number of operations can be further reduced. Instead of using all $N \times L$ samples of function $\tilde{p}_l(m|\mathbf{d}, \mathbf{M})$ only L values can be used. The information needed for the localization of change-point is included in time-progress of any value of $\tilde{p}_l(m|\mathbf{d}, \mathbf{M}), m = \text{const}$. Best choice is $m = N/2$ giving the sequence $\tilde{p}_l(N/2|\mathbf{d}, \mathbf{M}), l=1, \dots, L$. This choice ensures the best numerical stability and the lowest sensitivity to disturbances comparing with other asymmetric choices $m \neq N/2$; it also gives the best results for the model order selection due to fact that the “left” and “right” parts of data vector have the same lengths.

2.1.1 Summary of recursive algorithms

The RBACDN and RBSCDN algorithms can be summarized as follows:

- a) Initialization of $D, \Phi_1, \mathbf{g}_1, \Delta_1, \Delta$ for $l=1$ (first data segment) and $m = N/2 \rightarrow \tilde{p}_l(N/2|\mathbf{d}, \mathbf{M})$
- b) Recursion=data and position update
 - a. Data update for a new sample for (2), (3)
 - b. Removing old sample (2), (3)
 - c. Position update for (2)
 - $\rightarrow \tilde{p}_l(N/2|\mathbf{d}, \mathbf{M}), l = 2, \dots$
- c) Post-processing of $\tilde{p}_l(N/2|\mathbf{d}, \mathbf{M}), l = 2, \dots$
 - a. Searching for two adjacent minima of smoothed $\tilde{p}_l(N/2|\mathbf{d}, \mathbf{M})$ (low-pass filter, cut-off $\approx p/100$) \rightarrow stationary points.
 - b. Searching for one maximum of $\tilde{p}_l(N/2|\mathbf{d}, \mathbf{M})$ between two stationary points \rightarrow final change-points.

Notes: Searching for minima of signal is relatively simple because of their almost the same values (which are close to zero) due to normalization of BACD or BPCD by the Bayesian evidence. Searching for one maximum decreases the number of false alarms and excludes the need of smoothing $\tilde{p}_l(N/2|\mathbf{d}, \mathbf{M})$. The logarithm allows using longer window length (even 3000 samples) without any numerical instability. The choice of the window length depends on signal parameters.

Details of step b)

A) Data update for function (4), (8) and (9) are given by

I) Adding new data row \mathbf{x}

$$\hat{D}_l = D_l + d[N+1] \mathbf{d}[N+1];$$

$$\hat{\mathbf{g}}_{1,l} = \mathbf{g}_{1,l} + d[N+1] \mathbf{x}; \mathbf{W}_l = \Phi_l \mathbf{x}^T; l = 1 + \mathbf{x} \mathbf{W}_l;$$

$$\hat{\Delta}_{1,l} = l \Delta_{1,l}; \hat{\Phi}_{1,l} = \Phi_{1,l} - \mathbf{W}_l \mathbf{W}_l^T / l, \text{ where}$$

$$\mathbf{x} = \begin{bmatrix} 0 & \mathbf{I}_{M_1} & 0 \\ \mathbf{1} & \mathbf{d}[N] & \mathbf{d}[N-1] & \mathbf{L} & \mathbf{d}[N+1-M_2] \end{bmatrix}$$

for the RBACDN;

$$\mathbf{x} = \begin{bmatrix} 0 & \mathbf{I}_{P_1+1} & 0 \\ \mathbf{1} & t^0[N+1] & t^1[N+1] & \mathbf{L} & t^{P_1}[N+1] \end{bmatrix}$$

for the RBPCDN;

$$\mathbf{x} = [d[N] \quad d[N-1] \quad \mathbf{L} \quad d[N+1-M]]$$

for the *BE* normalizing BACD;

$$\text{and } \mathbf{x} = [t^0[N+1] \quad t^1[N+1] \quad \mathbf{L} \quad t^P[N+1]]$$

for the *BE* normalizing BPCD.

II) Removing old data row \mathbf{z}

$$D_{l+1} = \hat{D}_l - d[l] \mathbf{d}[l];$$

$$\hat{\mathbf{g}}_{1,l} = \hat{\mathbf{g}}_{1,l} - d[l] \mathbf{z}; \mathbf{W}_l = \hat{\Phi}_{1,l} \mathbf{z}^T; l = 1 - \mathbf{z} \mathbf{W}_l;$$

$$\hat{\Delta}_{1,l} = l \Delta_{1,l}; \hat{\Phi}_{1,l} = \hat{\Phi}_{1,l} + \mathbf{W}_l \mathbf{W}_l^T / l$$

where

$$\mathbf{z} = [d[l-1] \quad d[l-2] \quad \mathbf{L} \quad d[l-M_1] \quad \begin{bmatrix} 0 & \mathbf{I}_{M_2} & 0 \end{bmatrix}]$$

for the RBACDN;

$$\mathbf{z} = [t^0[l] \quad t^1[l] \quad \mathbf{L} \quad t^P[l] \quad \begin{bmatrix} 0 & \mathbf{I}_{P_2+1} & 0 \end{bmatrix}]$$

for the RBPCDN;

$$\mathbf{z} = [d[l] \quad d[l-1] \quad \mathbf{L} \quad d[l+1-M]]$$

for the *BE* normalizing BACD;

$$\text{and } \mathbf{z} = [t^0[l] \quad t^1[l] \quad \mathbf{L} \quad t^P[l]]$$

for the *BE* normalizing BPCD.

B) Position update for (2) and $m = N/2$

I) Replacing of $m+1$ row \mathbf{r}_1 of \mathbf{G}_1 with row of zeros

$$\hat{\mathbf{g}}_{1,l} = \mathbf{g}_{1,l} - d[m+1] \mathbf{r}_1; \mathbf{W}_l = \Phi_{1,l} \mathbf{r}_1^T;$$

$$l = 1 - \mathbf{r}_1 \mathbf{W}_l; \hat{\Delta}_{1,l} = l \Delta_{1,l};$$

$$\hat{\Phi}_{1,l} = \hat{\Phi}_{1,l} + \mathbf{W}_l \mathbf{W}_l^T / l$$

$$\mathbf{r}_i = \mathbf{r}_A = \begin{bmatrix} 0 & \mathbf{I}_{M_1} & 0 \\ \mathbf{L} & & \end{bmatrix} \begin{bmatrix} d[m] & d[m-1] & \mathbf{L} & d[m+1-M_2] \end{bmatrix}$$

for the RBACDN,

$$\mathbf{r}_i = \mathbf{r}_P = \begin{bmatrix} 0 & \mathbf{I}_{P+1} & 0 \\ \mathbf{L} & & \end{bmatrix} \begin{bmatrix} t^0[m+1] & t^1[m+1] & \mathbf{L} & t^{P+1}[m+1] \end{bmatrix}$$

for the RPCDN.

II) Replacing of $m+1$ row \mathbf{r}_i of \mathbf{G}_1 with new data \mathbf{q}_i

$$\mathbf{g}_{1,l+1} = \bar{\mathbf{g}}_{1,l} + d[m+1] \mathbf{q}_i; \mathbf{W}_i = \bar{\Phi}_{1,l} \mathbf{q}_i^T;$$

$$l = l + 1 + \mathbf{q}_i \mathbf{W}_i;$$

$$\Delta_{1,l+1} = l \bar{\Delta}_{1,l}; \Phi_{1,l+1} = \bar{\Phi}_{1,l} - \mathbf{W}_i \mathbf{W}_i^T / l$$

$$\mathbf{q}_i = \mathbf{q}_A = \begin{bmatrix} d[m] & d[m-1] & \mathbf{L} & d[m+1-M_1] & 0 & \mathbf{I}_{M_2} & 0 \end{bmatrix}$$

for the RBACDN,

$$\mathbf{q}_i = \mathbf{q}_P = \begin{bmatrix} t^0[m+1] & t^1[m+1] & \mathbf{L} & t^{P+1}[m+1] & 0 & \mathbf{I}_{P_2} & 0 \end{bmatrix}$$

for the RBPCDN.

The model order selection needed for the proper algorithm performance is dependent on signal parameters and it is not included yet in this algorithm.

3 Experiments and results

The described algorithms were preliminary validated by experiments with synthetic and real signals. Extensive simulations for signals modelled by various autoregressive and polynomial models were evaluated. Both, the RBACDN and the RBPCDN have shown performance close to the BACD and the BPCD performance for one change-point and window length greater than 200 samples. The RBACDN and the RBPCDN behaviour for real signals are illustrated in Figs. 1 to Figs. 3 for various signal changes.

Fig. 1 shows the segmentation process using the RBACDN for a non-stationary oboe signal. The average error rate evaluated on several realizations of oboe signal is about 18%. Results show the tendency for change-point omissions rather than for change-point insertions (false alarms). Cepstral distances (for definition and using for the classification of signal changes see [12]) vary from 2.5 db (very small, almost inaudible change) to 8 db (strong audible change). As it can be seen from the spectrogram the changes are gradual rather than abrupt.

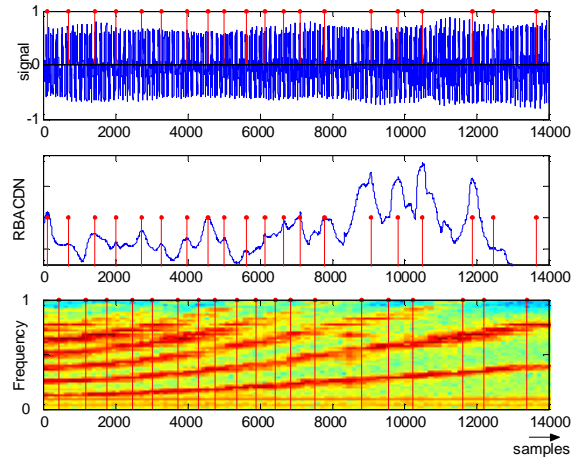


Fig.1 Oboe signal segmentation. From top to bottom: waveform, RBACDN output (11), and spectrogram with estimated change-points. Model order $M_1=M_2=20$, sliding window 2000 samples, cepstral distances are from 2.5 db to 8 db.

Another example of abrupt and strong changes in drum signal illustrates Fig. 2. This type of signal requires polynomial rather than autoregressive modeling because cepstral changes are smaller than 1 db. This small values indicate that the autoregressive model is inadequate. The polynomial model zero order is used in this case because there are abrupt changes especially in signal level rather than in signal slope. The average error rate evaluated on several realizations of the oboe signal is about 10%. Results show the tendency rather for change-point insertions than for change-point omissions.

Fig. 3 illustrates speech segmentation which seems to be the most problematic application of the suggested detectors. Results of speech segmentation show a high sensitivity to the choice of window length and the proper model order selection. The described RBACDN is not yet modified for an automatic model order selection using BE and for the use of different model orders for “left” and “right” parts of signal. Moreover, speech is a highly non-stationary process requiring a small length of sliding window. But the short window generates inconsistent change-point estimates causing extra change-point insertions. When a long window is used the RBACDN detects more change-points, especially in consonant parts of speech. Also consonant-vowels parts are not reliable detected in this case.

The segmentation experiments with real signals were validated by the multiple model algorithm [11] but above all by the inspection of segmented spectrograms (see bottom of all figures), and by listening the isolated segments containing tones or parts of speech sentences.

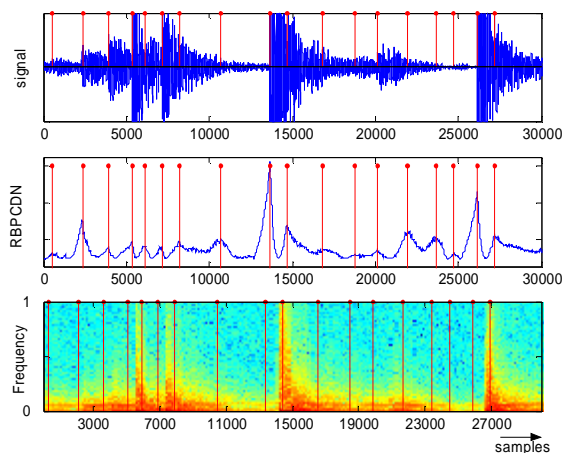


Fig.2 Drum signal segmentation. From top to bottom: waveform, RBPCDN output (11), and spectrogram with estimated change-points. Model order $M1=M2=0$, sliding window 2000 samples.

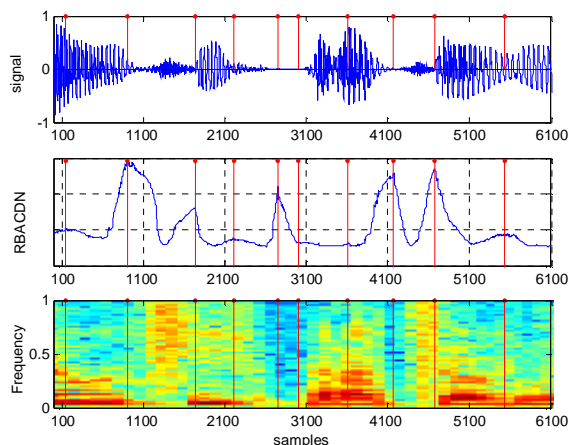


Fig.3 Speech signal segmentation. From top to bottom: waveform, RBACDN output (11), and spectrogram with estimated change-points. Model order $M1=M2=8$, sliding window 1000 samples.

4 Conclusion

Two new sliding window segmentation algorithms based on the normalization of the probability density function by the Bayesian evidence was suggested. The reduction of computational costs were described and verified by experiments and illustrated on real signals. Further research will be focused on an automatic model order selection using the Bayesian evidence, and the optimization of RBACDN for the automatic speech segmentation.

Acknowledgement:

Theoretical part of this work has been supported by the research program Transdisciplinary Research in

Biomedical Engineering MSM210000012 of the Czech University in Prague while the experimental part including evaluation results by the grant GA 102/02/0124 Voice Technologies for Support of Information Society.

References:

- [1] F. Gustafsson, *Adaptive filtering and change detection*. J. Wiley New York, 2000.
- [2] J. J. K. Ó Ruanaidh and W. J. Fitzgerald, *Numerical Bayesian methods applied to signal processing*. Springer-Verlag New York, 1996.
- [3] A. Procházka, J. Uhlíř, J., P.J.W. Rayner, N.G. Kingsbury (eds.), *Signal Analysis and Prediction*. Birkhauser, Boston, 1998.
- [4] M. J. K. Basseville and I. V. Nikoiforov, *Detection of abrupt changes: Theory and applications*. Prentice-Hall, Inc. New York, 1999.
- [5] J.J.K.O'Ruanaidh, W.J.Fitzgerald and K.J.Pope, Recursive Bayesian location of a discontinuity in time series, in *Proc. International Conference on Acoustics, Speech and Signal Processing*, Adelaide, Australia, 1994.
- [6] J-Y. Tournert, M Doisy, and M. Lavielle, "Bayesian off-line detection of multiple change-points corrupted by multiplicative noise; application to SAR image edge detection," *Signal Processing*, vol. 83, pp. 1871–1887, 2003.
- [7] E.Punskaya, C. Andrieu, A. Doucet, and W. J. Fitzgerald, "Bayesian curve fitting using MCMC with applications to signal segmentation," *IEEE Trans. on Signal Processing*, vol. 50, pp. 747–758, Mar. 2002.
- [8] S. Cheng and H. Wang, A sequential metric-based audio segmentation method via the Bayesian information criterion, in *Proc. European Conference on Speech Communication and Technology (Eurospeech2003)*, Geneva, Switzerland, Sept 2003.
- [9] S. J. Godsill and J. W. Rayner, *Digital audio restoration*. Springer-Verlag New York, 1998.
- [10] G. Tzanetakis and P. Cook, Multifeature audio segmentation for browsing and annotation, in *Proc. 1999 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, New Paltz, New York, 1999.
- [11] K. Jensen and D. Murphy, Segmentation melodies into notes, in *Proceedings of the DSAGM*, Copenhagen, Denmark, 2001.
- [12] R. Cmejla, and P. Sovka: The using of Bayesian detector in signal processing, in *Proc. IASTED ISIP'99*, Bahamas, 1999.
- [13] P. Anderson, "Adaptive forgetting in recursive identification through multiple models," *International Journal of Control*, vol. 42, pp. 1175–1193, 1985.