# Evaluation of BIC-based algorithms for audio segmentation

Mauro Cettolo [a,*], Michele Vescovi [b], Romeo Rizzi [b]

[a] *ITC-irst, Centro per la Ricerca Scientifica e Tecnologica, Via Sommarive, 18 Povo do Trento, I-38050 Povo, Trento, Italy*
[b] *Università degli Studi di Trento, Facoltà di Scienze MM.FF.NN. I-38050 Povo, Trento, Italy*

## Abstract

The Bayesian Information Criterion (BIC) is a widely adopted method for audio segmentation, and has inspired a number of dominant algorithms for this application. At present, however, literature lacks in analytical and experimental studies on these algorithms. This paper tries to partially cover this gap.

Typically, BIC is applied within a sliding variable-size analysis window where single changes in the nature of the audio are locally searched. Three different implementations of the algorithm are described and compared: (i) the first keeps updated a pair of sums, that of input vectors and that of square input vectors, in order to save computations in estimating covariance matrices on partially shared data; (ii) the second implementation, recently proposed in literature, is based on the encoding of the input signal with cumulative statistics for an efficient estimation of covariance matrices; (iii) the third implementation consists of a novel approach, and is characterized by the encoding of the input stream with the cumulative pair of sums of the first approach.

Furthermore, a dynamic programming algorithm is presented that, within the BIC model, finds a globally optimal segmentation of the input audio stream.

All algorithms are analyzed in detail from the viewpoint of the computational cost, experimentally evaluated on proper tasks, and compared.

© 2004 Published by Elsevier Ltd.

* Corresponding author. Tel.: +39-461-214-551; fax: +39-461-314-591.
*E-mail addresses:* cettolo@itc.it (M. Cettolo), vescovi@kirk.science.unitn.it (M. Vescovi), romeo@science.unitn.it (R. Rizzi).

## 1. Introduction

In the last years, efforts have been devoted amongst the research community to the problem of audio segmentation. The number of application of this procedure is considerable: from the extraction of information from audio data (e.g., broadcast news, recording of meetings) to the automatic indexing of multimedia data, or the improvement of accuracy for recognition systems. Typically, these tasks are performed by complex systems, consisting of a number of modules, some of them computationally expensive. Although in these systems the audio segmentation does not represent a computational bottleneck, the response time of the segmentation module can become an important issue under specific requirements. For example, ITC-irst delivered to RAI (the national Italian broadcasting company) an off-line system for the automatic transcription of broadcast news programs. Part of the requirement was also to supply a real-time version of the transcription station, able to guarantee adequate performance. Given these constraints, the speed of each component needs to be maximized without affecting its accuracy.

The segmentation problem has been handled in different ways that can be roughly grouped in three classes:

`energy` based methods: each silence occurring in the input audio stream is detected either by using an explicit model for the silence or by thresholding the signal energy. Segment boundaries are then located in correspondence of detected silences.

`metrics` based methods: the segmentation of the input stream is achieved by evaluating its ''distance'' from different segmentation models. Distances can be measured by the Hotelling's $T^2$ test (Wegmann et al., 1999), the Kullback Leibler distance (Kemp et al., 2000; Siegler et al., 1997), the generalized likelihood ratio (Gish et al., 1991), the entropy loss (Kemp et al., 2000), and the Bayesian Information Criterion (BIC) (Schwarz, 1978).

`explicit models` based methods: models are built for a given set of pre-determined acoustic classes – e.g., female and male speakers, music, noise, etc. Typically, the input data stream is classified from the maximum likelihood principle, through a dynamic programming decoding. The time indexes where the classification changes from one class to another are assumed to be the segment boundaries (Hain et al., 1998). Alternatively, in (Lu et al., 2001) Support Vector Machines are employed to learn class boundaries (Scheirer and Slaney, 1997); compares the Gaussian Maximum A-Posteriori estimator, the Gaussian Mixture Model, the Nearest-Neighbor classifier and a spatial partitioning scheme based on K-d trees.

The main limitations of the first and the third approach are evident. Regarding the energy based methods, there is only a partial correlation between changes in the nature of the audio and silences. On the other hand, these methods are simple to implement and can perform their (limited) task in linear time and with sufficient precision – provided the absence of significant variations in the background acoustic conditions. With explicit models of acoustic classes, changes occurring within the same class are undetectable; for example, if only one model for female voices is employed, no change can be detected within a dialogue between female speakers. Moreover, it is required both to know in advance the classes of interest and the availability of suitable data for their training. On the other hand, these methods can reach very high accuracy rates with linear time cost complexity.

This paper concentrates on a specific class of metric-based methods. Metric-based methods do not require any prior knowledge, nor a training stage. They are efficient (linear in time if some

70 approximations are introduced – as shown later), simple to implement and are able to give good
71 results. This explains why in a lot of laboratories much attention has been devoted to this type of
72 methods, and in particular to the BIC (Cettolo, 2000; Cettolo and Vescovi, 2003; Chen and
73 Gopalakrishnan, 1998; Delacourt et al., 1999; Harris et al., 1999; Sivakumaran et al., 2001;
74 Tritschler and Gopinath, 1999; Vescovi et al., 2003; Wellekens, 2001). Typically, the BIC "dis-
75 tance" is used locally, within a shift variable-size window, where single changes in the nature of
76 audio are searched. In this paper, we deal with metric-based audio segmentation algorithms based
77 on the BIC.

78     In the research community it is customary to combine different (almost orthogonal) solutions to
79 exploit advantages and circumvent limitations. This is also the case of audio segmentation.
80 Usually, real systems include a partitioning module that segments, classifies and clusters the audio
81 stream via a combination of different algorithms. Among the most successful systems we find
82 indeed highly integrated partitioning modules, as for example in Gauvain et al. (1998) and in Hain
83 et al. (1998). Even if this trend has allowed the development of very competitive systems, this has
84 led to the situation where it is difficult to have a clear insight on how to improve the state of the
85 art.

86     As a matter of fact, literature regarding analytical and experimental studies of segmentation
87 algorithms remains – somewhat – limited. The only known exception is Sivakumaran et al. (2001),
88 where an efficient approach to the local BIC-based segmentation algorithm is proposed and its
89 computational cost briefly analyzed. However, only a very limited comparison with other algo-
90 rithms is given. In that work, the input audio stream is progressively encoded by cumulative
91 statistics, and the encoding is exploited to avoid redundant operations in the computation of BIC
92 values.

93     In this paper, we try to cover part of the studies on segmentation algorithms by proposing: (i)
94 an innovative method for the implementation of the BIC-based local algorithm – this guar-
95 antees (to the best of our knowledge) the lowest computational cost; (ii) a global algorithm
96 capable of finding the optimal BIC input segmentation – this represents the performance upper
97 bound for the class of the BIC-based segmentation algorithms. The algorithms are described
98 and compared, both analytically and experimentally, with the method given in Sivakumaran
99 et al. (2001).

100     The paper is organized as follows. In Section 2, a general definition of the audio segmentation
101 problem is given.

102     Section 3 presents the corpora used for experiments, together with the evaluation measures; a
103 brief description of the signal processing front-end is also given.

104     In Section 4, the BIC-based local algorithm as proposed in Delacourt et al. (1999) is described,
105 building on the general framework provided in Section 2. The computational cost of the local
106 algorithm presented in Sivakumaran et al. (2001) is analyzed in detail (Section 4.1.2) and com-
107 pared, both in theory and experimentally, with the cost of two other possible approaches: one
108 more direct but more expensive, which represents a reference (Section 4.1.1); and a novel method
109 that combines the good ideas of the other two (Section 4.1.3). On a test set of broadcast news
110 programs, the new approach is 35% faster than the one proposed in Sivakumaran et al. (2001)
111 (Section 4.2).

112     Section 5 presents, within the general framework of Section 2, a Dynamic Programming (DP)
113 algorithm which uses the BIC method to find the globally optimal segmentation of the input

114 audio. The algorithm with its computational costs are described in detail. It is also compared to
115 the most efficient implementation of the local algorithm. On the 2000 NIST Speaker Recognition
116 Evaluation test set, the global algorithm outperforms the local one by 2.4% (relative) $F$-score in
117 the detection of changes, but is 38 times slower.
118   A summary ends the paper.

## 2. The segmentation problem and the BIC

120   Segmenting an audio stream consists of detecting the time indexes corresponding to changes in
121 the nature of the audio, this to isolate segments that are acoustically homogeneous. This can be
122 seen as a particular instance of the more general problem of partitioning data into distinct ho-
123 mogeneous regions (Baxter, 1996). The data partitioning problem arises in all applications that
124 require partitioning data into chunks, e.g., image processing, data mining, text processing, etc.
125   The problem can be formulated as follows. Let $\mathcal{O} = o_1, o_2, \ldots, o_N$ be an ordered sequence of
126 observations, called the *sample*, in which each $o_i$ is a vector in $\mathbb{R}^d$. We assume that the sample is
127 generated by a Gaussian process with a certain number of transitions. The problem of segmen-
128 tation is that of detecting all the transition points in the data set. In the ambit of acoustic seg-
129 mentation, transitions correspond to changes in the nature of the audio, represented by the
130 sample $\mathcal{O}$.
131   A particular model of the input data is characterized by a specific number $c$ of changes and a
132 specific set $\{1 \leqslant t_1, t_2, \ldots, t_c < N\}$ of time indexes corresponding to these changes. In such a way,
133 the input data are modeled by a set of $c + 1$ Gaussian distributions, each one generating an
134 observation subsequence bounded by two consecutive time indexes out of $\{1, t_1, \ldots, t_c, N\}$.
135   Among all possible models of the sample, the "best" one has to be selected, and its time
136 indexes will define the segmentation of the input stream. The best model is the one that better
137 fits the observations. The application of the maximum likelihood principle would however in-
138 variably lead to choosing the model with the maximum number of changes ($c_{\max} = N - 1$), as
139 this model has the highest number of free parameters. In order to take into account the notion
140 of "dimension" of the model, the following extension to the maximum likelihood principle was
141 first proposed by Akaike (1977). The AIC (Akaike's Information Criterion) suggests to maxi-
142 mize the likelihood for each model $M_i$ separately, obtaining say $L_{M_i} = L_{M_i}(\mathcal{O})$, and then choose
143 the model for which $(\log L_{M_i} - F_{M_i})$ is largest, where $F_{M_i}$ is the number of free parameters of the
144 model $M_i$.
145   Several model selection criteria that can be applied to Akaike's framework of model selection
146 have been proposed in the literature (see Cettolo and Federico, 2000 for a review). In general, each
147 criterion proposes the introduction of a penalty function $P$ that takes into account the dimension
148 of the model.
149   The BIC is the penalty function proposed in Schwarz (1978), and is defined as

$$P_M = \frac{F_M}{2} \log N. \tag{1}$$

151 As an example, let us consider a sample of three observations $\mathcal{O} = o_1, o_2, o_3$. The possible models
152 of $\mathcal{O}$ to be compared are:

YCSLA 259
DISK / 30/6/04

ARTICLE IN PRESS

No. of pages: 24
DTD 4.3.1/SPS

*M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx* 5

153 • a 1-Gaussian process (let it be $M_1$):

$$o_1, o_2, o_3 \sim_{iid} N_d(\mu, \Sigma),$$

155 • two 2-Gaussians processes ($M_2$ and $M_3$):

$$o_1, o_2 \sim_{iid} N_d(\mu_A, \Sigma_A),$$

$$o_3 \sim_{iid} N_d(\mu_B, \Sigma_B)$$

158 and:

$$o_1 \sim_{iid} N_d(\mu_A, \Sigma_A),$$

$$o_2, o_3 \sim_{iid} N_d(\mu_B, \Sigma_B),$$

161 • a 3-Gaussians process ($M_4$):

$$o_1 \sim_{iid} N_d(\mu_A, \Sigma_A),$$

$$o_2 \sim_{iid} N_d(\mu_B, \Sigma_B),$$

$$o_3 \sim_{iid} N_d(\mu_C, \Sigma_C).$$

165 Once $(\log L_{M_i}(\mathcal{O}) - P_{M_i})$ has been computed for each $i = 1, 2, 3, 4$, the highest value gives the
166 winner model. If it is, for example, $M_2$, the best segmentation of $\mathcal{O}$ under the BIC consists of the
167 following two segments: $(o_1 o_2)$ and $(o_3)$.

168 *2.1. Search*

169 Given the set $\{M_i\}$ of possible models of the input sequence, the segmentation problem now
170 reduces to solve the following maximization:

$$i_{\max} = \arg\max_i \log L_{M_i} - P_{M_i}. \tag{2}$$

172 Instead of explicitly listing the set of all possible models, whose size is $O\{2^N\}$, polynomial search
173 algorithms can be employed. A widely adopted search algorithm uses a sliding variable-size
174 analysis window, where the maximization given above is used to detect single changes. That is,
175 models compared in Eq. (2) have only one transition or none at all. The algorithm and three
176 implementations with different computational costs are described in Section 4. The algorithm
177 main limitation is that the search is done locally – hence the objective function that determines the
178 most plausible interpretations of the input audio stream is not globally maximized. In principle,
179 this could affect the effectiveness of the method, especially if the audio stream includes many
180 changes close to each other – as in the case of human–human conversations. In Section 5, a global
181 algorithm able to find the solution of Eq. (2) in the BIC framework, is presented.

182 **3. Data sets and evaluation measures**

183 Two corpora were selected for experiments, namely the Italian Broadcast News Corpus (IBNC)
184 and the 2000 NIST Speaker Recognition Evaluation corpus. The two collections are substantially

6                        *M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx*

185 different in the nature of audio recordings, allowing the evaluation of algorithms under very
186 different working conditions.

187 *3.1. The IBNC corpus*

188   The IBNC corpus (Federico et al., 2000), collected by ITC-irst and available through ELDA,
189 the Evaluations and Language resources Distribution Agency (IBNC, 2000), is a speech corpus of
190 radio broadcast news in Italian. It consists of 150 recordings, for a total of about 30 h, covering
191 radio news of several years. Data were provided by RAI. The corpus presents variations of topics,
192 speakers, channel band (i.e., studio versus telephone), speaking mode (i.e., spontaneous versus
193 planned), etc. It has been manually transcribed, segmented and labeled. Speaker gender and, when
194 possible, identity are also annotated.
195   For testing purposes, six programs (about 75 min of audio) were selected, where 212 changes
196 occur, between either different speakers or different acoustic classes (music, speech, noise, etc.).
197   In broadcast news, acoustically homogeneous segments are typically long. This makes IBNC
198 data suitable for assessing the ability of segmentation algorithms in the detection of changes
199 around which much homogeneous data is available.

200 *3.2. The 2000 NIST speaker recognition evaluation*

201   The 2000 NIST Speaker Recognition Evaluation (NIST, 2000) included a segmentation task,
202 where systems were required to identify speech segments corresponding to each of two unknown
203 speakers. The test set consists of 1000 telephone conversations, lasting about 1 min each, included
204 in the Disc r65_6_1. 777 Speakers, of both genders, pronounced a total of 46K turns/segments.
205 Fig. 1 shows the distribution of segment length. The longest segment is 26.5 s; the mean length is
206 1.3 s, while the median is less than 1 s; it is worth noticing that 81% of segments has length lower



Fig. 1. Distribution of length of NIST test set segments.

YCSLA 259
DISK / 30/6/04

ARTICLE IN PRESS

No. of pages: 24
DTD 4.3.1/SPS

*M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx* 7

207 than 2 s. This means that the NIST data are suitable to make evident if segmentation algorithms
208 are able to detect changes bounding very short segments.
209    The NIST test set was also employed in this work. For evaluation, we relied on a reference
210 segmentation obtained by merging the two reference segmentations provided by NIST for each
211 speaker. It was observed that often the end of a sentence of a speaker does not exactly correspond
212 to the beginning of the successive sentence of the other speaker: this is due either to small im-
213 precision of the manual annotation, or to a short overlap of the sentences, or to a brief silence
214 occurring between the two sentences. In such cases, short segments created when merging were
215 removed by an automatic procedure.

### 3.3. Evaluation measures

217    Performance of automatic change detection is calculated with respect to a set of target changes.
218 Tolerances in the detection can be introduced, e.g., $\pm0.5$ s: in such a way the target is an interval
219 rather than a single time index.
220    For comparing target and hypothesized changes, we adopt the *precision P* and the *recall R*
221 measures:

$$P = \frac{c}{c + i} \times 100,$$

$$R = \frac{c}{c + d} \times 100,$$

224 where $c$ is the number of target intervals which contain at least one hypothesized change, $i$ is the
225 number of hypothesized changes that do not fall inside any target interval (insertions), and $d$ is the
226 number of target intervals which contain no hypothesized change (deletions). In other words,
227 precision is related to the number of false alarms, i.e., the number of wrongly split segments, while
228 recall measures the deletion rate of correct changes, i.e., the number of wrongly merged segments.
229    The evaluation of the segmentation quality is made in terms of *F-score*, a combined measure of
230 $P$ and $R$ of change detection. *F*-score is a measure proposed by van Rijsbergen (Frakes and Baeza-
231 Yates, 1992) and is defined as

$$F = \frac{(1 + \beta^2)PR}{\beta^2 P + R},$$

233 where $\beta$ is a measure of the relative importance, to a user, of precision and recall. For example, $\beta$
234 levels of 0.5, indicating that a user was twice as interested in precision as recall, and 2, indicating
235 that a user was twice as interested in recall as precision, might be used. The most popular value
236 corresponds to $\beta = 1$, for which $F$ reduces to

$$F = \frac{2PR}{P + R}.$$

238 In this paper $\beta$ was set to 1.
239    Concerning the NIST data set, after the official evaluation, NIST made available both the
240 reference segmentations of the test files and a scoring script. The scoring script computes the
241 segmentation error rate, i.e., the percentage of speech from one speaker wrongly assigned to the

8                    *M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx*

242 other speaker. The measure is suitable for the NIST task which, by definition, is a classification
243 task. On the contrary, the algorithm presented here detects speaker/spectral changes, and no
244 attempt is done to classify segments in terms of speakers. The NIST scoring script does not
245 therefore fit with our evaluation requirements.

246    In fact, precision and recall (or the corresponding false alarm and miss detection rates) are
247 metrics that better assess the performance of the change detection algorithm. For this reason, the
248 two metrics and their *F*-score are also used in the evaluation of the algorithms with the NIST data
249 set.

250 *3.4. Audio processing*

251    Multivariate observations derive from a short time spectral analysis, performed over 20 ms
252 Hamming windows at a rate of 10 ms. For every window, 12 Mel scaled Cepstral coefficients and
253 the log-energy are evaluated.

## 4. Local search: an approximated algorithm

255    First of all, let us recall some basic results. Given a sample $\mathcal{O} = o_1, o_2, \ldots, o_N$ of observations
256 $o_i \in \mathbb{R}^d$, the likelihood function $L_{N_d(\mu,\Sigma)}(\mathcal{O})$ achieves its maximum value (Seber, 1984) in $\mu = \bar{o}$, the
257 sample mean, and

$$\Sigma = \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^{N} (o_i - \bar{o})(o_i - \bar{o})^{\text{tr}}, \tag{3}$$

259 the maximum-likelihood estimate of the covariance matrix. Moreover

$$L_{N_d(\bar{o},\hat{\Sigma})}(\mathcal{O}) = (2\pi)^{-Nd/2} |\hat{\Sigma}|^{-N/2} e^{-Nd/2}. \tag{4}$$

261 From here on, sample means and estimates of covariance matrices will be denoted simply by $\mu$ and
262 $\Sigma$: it should be clear from the context if symbols refer to variables or evaluated statistical esti-
263 mators.

264    The number of free parameters in a multivariate normal distribution is equal to the dimension
265 of the mean plus the number of variances and covariances to be estimated. For the full covariance
266 matrix case, it is

$$F_{N_d(\mu,\Sigma)} = d + d \frac{(d+1)}{2}. \tag{5}$$

268 Assuming that the sequence $o_1, \ldots, o_N$ contains at most one change, and that data are generated
269 by a Gaussian process, the maximization of Eq. (2) regards the following $N$ different statistical
270 models:

271 • $N - 1$ two-segments models $M_i$ ($i = 1, \ldots, N - 1$), where model $M_i$ states:

$$o_1, \ldots, o_i \sim_{iid} N_d(\mu_1, \Sigma_1),$$

$$o_{i+1}, \ldots, o_N \sim_{iid} N_d(\mu_2, \Sigma_2),$$

274 • one single-segment model $M_N$ which states

$$o_1, o_2, \ldots, o_N \sim_{iid} N_d(\mu, \Sigma).$$

276 The selection of the best model can then rely on the following decision rule:

277 • Look for the best two-segments model $M_{i_{\max}}$ for the data

$$i_{\max} = \arg\max_{i=1,\ldots,N-1} \log L_{M_i} - P_{M_i}.$$

279 • Take the one-segment model function

$$\log L_{M_N} - P_{M_N}. \tag{6}$$

281 • Choose to segment the data at point $i_{\max}$ if and only if

$$(\log L_{M_{i_{\max}}} - \log L_{M_N}) - (P_{M_{i_{\max}}} - P_{M_N}) > 0. \tag{7}$$

283 If the BIC penalty term (1) is adopted, from Eqs. (4) and (5) it is easy to show that the above
284 described decision rule is equivalent to computing for each $i = 1, \ldots, N-1$ the quantity

$$\Delta\text{BIC}_i = \frac{N}{2} \log|\Sigma| - \frac{i}{2} \log|\Sigma_1| - \frac{(N-i)}{2} \log|\Sigma_2| - \lambda P_{N_d(\mu,\Sigma)}, \tag{8}$$

286 where $\Sigma$, $\Sigma_1$ and $\Sigma_2$ are the maximum likelihood covariance estimates on $o_1 \ldots o_N$, $o_1 \ldots o_i$ and
287 $o_{i+1} \ldots o_N$, respectively, and to hypothesize a change in $i_{\max}$, if $i_{\max}$ is the index that maximizes

```
init_window(1, N_min)

while(not end stream)

    (ΔBIC_imax, imax) ← compute_ΔBIC(δ_l)

    while(ΔBIC_imax ≤ 0 and
            current_win_size < N_max and
            not end stream)
      growth_win(ΔN_grow)
      (ΔBIC_imax, imax) ← compute_ΔBIC(δ_l)

    while(ΔBIC_imax ≤ 0 and not end stream)
      shift_win(ΔN_shift)
      (ΔBIC_imax, imax) ← compute_ΔBIC(δ_l)

    if(ΔBIC_imax > 0) then
      center_win(imax, min(current_win_size, N_second))
      (ΔBIC_ichange, ichange) ← compute_ΔBIC(δ_h)
      if(ΔBIC_ichange > 0) then
        output(ichange)
        init_window(ichange + 1, N_min)
      else
        init_window(imax − N_margin + 1, N_min)
```

Fig. 2. Pseudocode of the local sliding window algorithm.

288  $\Delta BIC_i$ and $\Delta BIC_{i_{\max}} > 0$. $\lambda \in \mathbb{R}$ is a weight, which allows to tune the sensitivity of the method to
289  the particular task under consideration.

290      In order to apply the search of single changes to an arbitrary large number of potential changes,
291  we implemented the local algorithm depicted in Fig. 2, inspired by that proposed in Delacourt et
292  al. (1999). The main idea is to have a shifting variable-size window for the computation of $\Delta BIC$
293  values. The algorithm dynamically adapts the window size in such a way that no more than one
294  change falls inside the window, ensuring the computation of reliable statistics and bounding the
295  computational cost. Moreover, to save computations, $\Delta BIC$ values are not computed on all
296  observations, but at a lower resolution, namely once every $\delta$ observations. The resolution is
297  successively increased if a potential change is detected, in order to validate it and to refine its time
298  position.

299      The main steps of the algorithm are:

300      **Search start**. $\Delta BIC$ values are computed only for the first $N_{\min}$ observations. $N_{\min}$ is the mini-
301  mum size of the window. Values are computed with low resolution, e.g., $\delta_l = 30$. In order to have
302  enough observations for computing both $\Sigma_1$ and $\Sigma_2$, $\Delta BIC$ are not computed for the $N_{\text{margin}}$ in-
303  dexes close to the left and right boundaries of the window.

304      **Window growth**. The window is enlarged by including $\Delta N_{\text{grow}}$ input observations until a change
305  is detected, or a maximum size $N_{\max}$ is reached.

306      **Window shift**. The $N_{\max}$-sized window is shifted forward by $\Delta N_{\text{shift}}$ observations.

307      **Change confirmation**. If in one of the three previous steps a change is detected, $\Delta BIC$ values are
308  re-computed with the high resolution, e.g., $\delta_h \approx \delta_l / 5$, centering the window at the hypothesized



Fig. 3. Working scheme of the local sliding window algorithm.

309 change. The current size of the window is kept, unless it is larger than $N_{\text{second}}$ observations, in
310 which case it is narrowed to that value. If a change is detected again, it is output by the algorithm.
311    **Window reset**. After the change confirmation step, the algorithm has to go on resizing the
312 analysis window to the minimum value $N_{\min}$ and locating it in a position dependent on the result
313 of the confirmation step (see Fig. 3).

314 *4.1. Computations*

315    In the following subsections, three possible implementations of the algorithm described above,
316 ordered according to decreasing computational cost, are presented in detail.

317 *4.1.1. The sum approach (SA)*
318    The evaluation of Eq. (8) determines the overall computational cost of the algorithm presented
319 above, since a high number of $\Delta$BIC values have to be computed for each window.
320    An efficient way to compute the determinant of the covariance matrix is based on the Cholesky
321 decomposition which requires $O\{d^3/6\}$ operations. The estimation of the mean vector $\mu$ and the
322 covariance matrix $\Sigma$ on $N$ $d$-sized observations $o_a, \ldots, o_b$:

$$\mu_a^b = \frac{1}{N} \sum_{i=a}^{b} o_i, \tag{9}$$

$$\Sigma_a^b = \frac{1}{N} \left( \sum_{i=a}^{b} o_i \cdot o_i^{\text{tr}} \right) - \mu_a^b \cdot \mu_a^{b\,\text{tr}} \tag{10}$$

325 requires, respectively, $d(N+1)$ and $d(d+1)(N+1.5)$ operations. Typically, the window size $N$ is
326 significantly larger than the vector dimension $d$, hence the computational cost for the evaluation
327 of the covariance matrix determinant could be discarded.
328    In order to reduce the computational cost of estimating likelihoods of the normal distributions
329 required for the computation of $\Delta$BIC values, it is convenient to keep the sums of the input
330 vectors (SV) and that of the square vectors (SQ):

$$\text{SV}_a^b = \sum_{i=a}^{b} o_i, \quad \text{SQ}_a^b = \sum_{i=a}^{b} o_i \cdot o_i^{\text{tr}}.$$

332 In fact, besides the easy computation of the needed parameters:

$$\mu_a^b = \frac{1}{N} \cdot \text{SV}_a^b, \quad \Sigma_a^b = \frac{1}{N} \cdot \text{SQ}_a^b - \mu_a^b \cdot \mu_a^{b\,\text{tr}},$$

334 the use of SV and SQ avoids many redundant operations in the computation of $\Delta$BIC values both
335 within a given window and after a window growth/shift. With reference to the notation in Table 1,
336 the following cases can happen:
337 • **growth** of the window by $\delta$ observations:
338    ○ $\text{SV}_{\tilde{T}} = \text{SV}_T + \sum_{j=n+N+1}^{n+N+\delta} o_j,$
339    ○ $\text{SQ}_{\tilde{T}} = \text{SQ}_T + \sum_{j=n+N+1}^{n+N+\delta} o_j \cdot o_j^{\text{tr}}.$

12                    *M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx*

Table 1
Notation

| | |
|---|---|
| $N$ | Current window size |
| $n$ | Index of the vector that preceeds the first index of the window |
| $T$ | Set of vectors inside the window $\{o_{n+1} \ldots o_{n+N}\}$ |
| $\tilde{T}$ | Set of vectors inside the window after a growth or a shift |
| $A_k$ | Set of the first $k\delta$ vectors of the window $\{o_{n+1} \ldots o_{n+k\delta}\}$ |
| $B_k$ | Set of the last $N - k\delta$ vectors of the window $\{o_{n+k\delta+1} \ldots o_{n+N}\}$ |
| $SV_X$ | Sum of vectors of the set $X$ |
| $SQ_X$ | Sum of square vectors of the set $X$ |
| $\Sigma_X$ | Covariance matrix on the vectors of the set $X$ |
| $\mu_X$ | Mean vector on the vectors of the set $X$ |

340  • **shift** of the window by $\delta$ observations:

341   ○ $SV_{\tilde{T}} = SV_T - \sum_{j=n+1}^{n+\delta} o_j + \sum_{j=n+N+1}^{n+N+\delta} o_j,$

342   ○ $SQ_{\tilde{T}} = SQ_T - \sum_{j=n+1}^{n+\delta} o_j \cdot o_j^{tr} + \sum_{j=n+N+1}^{n+N+\delta} o_j \cdot o_j^{tr}.$

343  • **computation** of $\Delta BIC_i$ (at resolution $\delta$):

344   ○ $SV_{A_i} = SV_{A_{i-1}} + \sum_{j=n+(i-1)\delta+1}^{n+i\delta} o_j,$

345   ○ $SQ_{A_i} = SQ_{A_{i-1}} + \sum_{j=n+(i-1)\delta+1}^{n+i\delta} o_j \cdot o_j^{tr},$

346   ○ $SV_{B_i} = SV_T - SV_{A_i},$

347   ○ $SQ_{B_i} = SQ_T - SQ_{A_i}.$

348 With this approach, in each step of the algorithm the number of operations for computing Eq. (8)
349 is:

350  • **growth** of the window by $\delta$ observations:

$$\underbrace{d(d+1) \cdot \delta}_{SQ_{\tilde{T}}} + \underbrace{d \cdot \delta}_{SV_{\tilde{T}}},$$

352  • **shift** of the window by $\delta$ observations:

$$\underbrace{d(d+1) \cdot 2 \cdot \delta}_{SQ_{\tilde{T}}} + \underbrace{d \cdot 2 \cdot \delta}_{SV_{\tilde{T}}},$$

354  • **computation** of the covariance matrix of the whole window:

$$\underbrace{d}_{\mu_T} + \underbrace{1.5 \cdot d(d+1)}_{\Sigma_T},$$

356  • **computation** of $\Sigma_{A_i}$ and $\Sigma_{B_i}$ required for the evaluation of $\Delta BIC_i$, with resolution $\delta$
357   ($\forall i, \ i = N_{margin}/\delta + 1, \ldots, (N - N_{margin})/\delta - 1$):

$$\underbrace{d \cdot \delta}_{SV_{A_i}} + \underbrace{d(d+1) \cdot \delta}_{SQ_{A_i}} + \underbrace{d}_{SV_{B_i}} + \underbrace{d(d+1)/2}_{SQ_{B_i}} + \underbrace{2 \cdot d}_{\mu_{A_i},\mu_{B_i}} + \underbrace{3 \cdot d(d+1)}_{\Sigma_{A_i},\Sigma_{B_i}} = d(d+1)(\delta + 3.5) + d(\delta + 3).$$

360 *4.1.2. The distribution approach (DA)*
361   In order to further reduce the computational cost of the algorithm, it is possible to evaluate Eq.
362 (8) through the approach proposed in Sivakumaran et al. (2001).

363   Let $\Sigma_N$ and $\mu_N$ be the sample covariance matrix and the mean of a set of $N$ $d$-dimensional
364 observations. If a (sub)set of $\Delta$ observations with covariance matrix $\Sigma_\Delta$ and mean vector $\mu_\Delta$ has to
365 be added or subtracted to that set, the parameters of the updated set of vectors can be computed
366 by:

$$\Sigma_{N\pm\Delta} = \frac{N}{N\pm\Delta}\Sigma_N \pm \frac{\Delta}{N\pm\Delta}\Sigma_\Delta \pm \frac{N\Delta}{(N\pm\Delta)^2}(\mu_N - \mu_\Delta)(\mu_N - \mu_\Delta)^{\mathrm{tr}}, \tag{11}$$

$$\mu_{N\pm\Delta} = \frac{N}{N\pm\Delta}\mu_N \pm \frac{\Delta}{N\pm\Delta}\mu_\Delta. \tag{12}$$

369 This formulation requires only $3 \cdot d(d+1) + d$ and $3 \cdot d$ operations for computing $\Sigma_{N\pm\Delta}$ and $\mu_{N\pm\Delta}$,
370 respectively, instead of $d(d+1)(N\pm\Delta+1.5)$ and $d(N\pm\Delta+1)$ required by the plain definitions.
371   The alternative approach consists in computing from the input audio stream $o_1, o_2, \ldots, o_{N_{\text{audio}}}$
372 the set of triples $(\Sigma_1^n, \mu_1^n, n)$, where $n = \delta_{\text{h}}, 2\delta_{\text{h}}, 3\delta_{\text{h}}, \ldots, N_{\text{audio}}$.
373   The key of this processing is Eqs. (11) and (12) which allow to obtain $(\Sigma_1^n, \mu_1^n, n)$ from
374 $(\Sigma_1^{n-\delta_{\text{h}}}, \mu_1^{n-\delta_{\text{h}}}, n-\delta_{\text{h}})$ and $(\Sigma_{n-\delta_{\text{h}}+1}^n, \mu_{n-\delta_{\text{h}}+1}^n, \delta_{\text{h}})$, where $\Sigma_{n-\delta_{\text{h}}+1}^n$ and $\mu_{n-\delta_{\text{h}}+1}^n$ are computed directly
375 from the vectors $o_{n-\delta_{\text{h}}+1}, o_{n-\delta_{\text{h}}+2}, \ldots, o_n$ through the definitions.
376   Since in this approach the estimation of a new distribution is based on already computed
377 distributions, it will be referred with the name "distribution approach" (DA).
378   By constraining $\delta_{\text{l}}$ and $N_{\text{second}}$ to be integers multiples of $\delta_{\text{h}}$ and by choosing $N_{\text{min}}, N_{\text{max}}, \Delta N_{\text{grow}}$,
379 $\Delta N_{\text{shift}}, N_{\text{margin}}$ to be divisible by $\delta_{\text{l}}$, it is possible to use the cumulative distributions $(\Sigma_1^n, \mu_1^n, n)$ for
380 the evaluation of $\Delta$BIC values and to reduce the cost of the computation. In fact, whatever the
381 step of the algorithm is, the covariance matrices required by Eq. (8) can be estimated from Eqs.
382 (11) and (12):

$$\Sigma_{n+1}^{n+N} = \frac{n+N}{N}\Sigma_1^{n+N} - \frac{n}{N}\Sigma_1^n - \frac{(n+N)n}{N^2}(\mu_1^{n+N}\mu_1^n)(\mu_1^{n+N} - \mu_1^n)^{\mathrm{tr}}, \tag{13}$$

$$\Sigma_{n+1}^{n+i} = \frac{n+i}{i}\Sigma_1^{n+i} - \frac{n}{i}\Sigma_1^n - \frac{(n+i)n}{i^2}(\mu_1^{n+i} - \mu_1^n)(\mu_1^{n+i} - \mu_1^n)^{\mathrm{tr}}, \tag{14}$$

$$\Sigma_{n+i+1}^{n+N} = \frac{n+N}{N-i}\Sigma_1^{n+N} - \frac{n+i}{N-i}\Sigma_1^{n+i} - \frac{(n+N)(n+i)}{(N-i)^2}(\mu_1^{n+N} - \mu_1^{n+i})(\mu_1^{n+N} - \mu_1^{n+i})^{\mathrm{tr}}. \tag{15}$$

386 Clearly, Eq. (13) is evaluated only once for a given window, while Eqs. (14) and (15) have to be
387 evaluated for each time index of interest (depending on the resolution). A scheme of the DA
388 approach is given in Fig. 4.
389   The number of operations required by each step of the algorithm with the DA approach is:
390 • **growth** or **shift** of the window by $\delta$ observations:

$$\underbrace{d(\delta+1)}_{\mu_{n+N+1}^{n+N+\delta}} + \underbrace{d(d+1)(\delta+1.5)}_{\Sigma_{n+N+1}^{n+N+\delta}} + \underbrace{3\cdot d(d+1) + 4\cdot d}_{(\Sigma_1^{n+N+\delta}, \mu_1^{n+N+\delta}, n+N+\delta)} = d(d+1)(\delta+4.5) + d(\delta+5).$$

392 Note that this cost is that of the input stream encoding.

*M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx*



Fig. 4. $\Delta\text{BIC}_i$ computation in the DA approach.

393 • **computation** of the covariance matrix of the whole window:

$$\underbrace{3 \cdot d(d+1) + d}_{\Sigma_{n+1}^{n+N} = \Sigma_T},$$

395 • **computation** of $\Sigma_{A_i}$ and $\Sigma_{B_i}$ required for the evaluation of $\Delta\text{BIC}_i$, with resolution $\delta$
396 ($\forall i,\ i = N_{\text{margin}}/\delta + 1, \ldots, (N - N_{\text{margin}})/\delta - 1$):

$$\underbrace{6 \cdot d(d+1) + 2 \cdot d}_{\Sigma_{n+1}^{n+i\delta} = \Sigma_{A_i}, \Sigma_{n+i\delta+1}^{n+N} = \Sigma_{B_i}}.$$

398 *4.1.3. The cumulative sum approach*
399 In the previous methods, the estimation of the statistics required for the computation of the
400 BIC are based either on the use of the sum and square sum of input vectors that fall inside the
401 analysis window, or on the use of the set of statistics computed only once, as soon as the ob-
402 servations from the input stream are available. A combination of the two basic ideas gives the
403 possibility to implement an even more efficient approach.
404 The idea is to encode the input stream, not through the distributions as in DA, but with the
405 sums of the SA approach, that is with the sequence of triples $(\text{SQ}_1^n, \text{SV}_1^n, n)$ computed at resolution
406 $\delta_{\text{h}}$. In this new approach, the cost of each step of the algorithm, which can be called "cumulative
407 sum approach" (CSA), is reported in the last column of Table 2. This table also includes, as a
408 summary, the SA and DA costs.

Table 2
Cost of each step of the SA, DA, and CSA approaches

| Step | SA | DA | CSA |
|---|---|---|---|
| Growth | $\delta d(d+1) + \delta d$ | $(\delta + 4.5)d(d+1) + (\delta + 5)d$ | $\delta d(d+1) + \delta d$ |
| Shift | $2\delta d(d+1) + 2\delta d$ | $(\delta + 4.5)d(d+1) + (\delta + 5)d$ | $\delta d(d+1) + \delta d$ |
| $\Sigma_T$ | $1.5d(d+1) + d$ | $3d(d+1) + d$ | $2d(d+1) + 2d$ |
| $\Sigma_{A_i}, \Sigma_{B_i}$ | $(\delta + 3.5)d(d+1) + (\delta + 3)d$ | $6d(d+1) + 2d$ | $4d(d+1) + 4d$ |

YCSLA 259
DISK / 30/6/04

ARTICLE IN PRESS

No. of pages: 24
DTD 4.3.1/SPS

*M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx* 15

409 The higher efficiency is given by: (i) the redundant computations in SA are avoided since each
410 input vector is used only once, during the encoding of the input stream; (ii) the new encoding is
411 cheaper than the DA encoding (cf. the grow/shift costs); (iii) the computation of covariance
412 matrices from sums requires less operations than starting from other distributions.

413 *4.2. Experimental evaluation*

414 *4.2.1. Database*
415 The main goal of these experiments is to compare the running times of the three implemen-
416 tations of the local algorithm. Since the variable-size window algorithm performs more operations
417 with long segments than with short ones, the IBNC corpus (see Section 3.1) suits the nature of this
418 kind of comparison.

419 *4.2.2. Costs comparison*
420 Since the computation of $\Delta BIC_i$ values is done approximately $N/\delta$ times in each window, the
421 total cost of the algorithm mainly depends on the cost of that operation, and this is the reason for
422 which the DA and CSA approaches are convenient with respect to the SA approach; for example,
423 the number of operations with the DA approach does not depend on $\delta$, whereas the SA does, and
424 in our case ($d = 13$) it results convenient for $\delta \geqslant 3$.
425 In order to validate the theoretical comparison described in Sections 4.1.1, 4.1.2 and 4.1.3, and
426 in particular the dependence of the overall computational cost from the resolution $\delta$, the three
427 implementations have been run with a simplified setup. We set $N_{min} = N_{max}$, in order to eliminate
428 the window grow step, and the value of $\lambda$ was set high enough that no candidate change was
429 detected, constraining the computations to be done only at resolution $\delta = \delta_l$.

Table 3
Theoretical and experimental costs comparison of SA, DA and CSA approaches

| $\delta$ | # Operations | | | Execution time (s) | | |
|---|---|---|---|---|---|---|
| | SA ($\times 10^6$) | DA/SA (%) | CSA/SA (%) | SA (s) | DA/SA (%) | CSA/SA (%) |
| 1 | 374.0 | 123.9 | 92.1 | 830.0 | 103.8 | 68.0 |
| 5 | 124.4 | 80.6 | 61.5 | 272.0 | 68.3 | 46.5 |
| 10 | 93.0 | 59.0 | 46.3 | 200.6 | 50.4 | 35.5 |
| 25 | 73.6 | 37.6 | 31.2 | 157.5 | 32.0 | 24.5 |

Setup: $N_{min} = N_{max} = 1000$, $\Delta N_{shift} = 200$, $N_{margin} = 50$, $d = 13$.

Table 4
Performance comparison of SA, DA and CSA approaches in their best setup: $N_{min} = 500, N_{max} = 2000, N_{second} = 1500, \Delta N_{grow} = 100, \Delta N_{shift} = 300, N_{margin} = 100, \delta_l = 25, \delta_h = 5, \lambda = 2.175$

| | $F$-score | Execution time | |
|---|---|---|---|
| | | s | ☆/SA (%) |
| SA | 88.4 | 289.4 | 100.0 |
| DA | 89.4 | 101.0 | 34.9 |
| CSA | 89.4 | 65.9 | 22.8 |

430    Given the setup in the caption and setting $N_{\text{audio}} = 50,000$, the total number of operations
431 required by the three approaches are given in the columns "#operations" of Table 3, for different
432 values of $\delta$. The costs include the computation of the covariance matrices determinant ($d^3/6$),
433 since it can no longer be neglected. Again with the setup in the caption, the execution times were
434 measured on a Pentium III 600 MHz on the 75-min IBNC test set.
435    Finally, Table 4 compares the three approaches on the algorithm setup detailed in the caption.
436 The table shows results in terms of *F-score* on the IBNC test set, together with execution times.
437 The slight difference in *F*-score for SA is due to some minor differences in the SA implementation.
438 Concerning the execution times, since $\delta_l$ was set to 25, the ratio between the costs of the three
439 implementations expected from the results of the last row of Table 3 is confirmed.

## 5. Global search: a DP algorithm

441    In order to avoid the approximation inherent in the local search, we developed a DP algorithm
442 able to find the globally optimal segmentation of the input audio stream according to the max-
443 imization (2) and the BIC penalty term (1). Between the $O\{2^N\}$ possible segmentations, it allows
444 to select the best one in $O\{N^3\}$ steps, as described in the following.
445    The BIC value in Eq. (1) can be seen as the difference between the log-likelihood $(\log L_M(\mathcal{O}))$
446 and the penalty term $P_M = \frac{F_M}{2}\log(N)$ that takes into account the complexity (size) of the model $M$.
447 Assuming a Gaussian process, the models $\{M\}$ among which the algorithm has to select the best
448 one, i.e., the one with the highest BIC value, differ in either the number of Gaussians or the way
449 the input stream is partitioned by the Gaussians.
450    Let $\mathscr{G}_k$ be the set of models with $k$ Gaussians. Since the number of free parameters of a $d$-di-
451 mensional Gaussian is $F = d + (d(d+1)/2)$, the penalty term, for each $M \in \mathscr{G}_k$, is

$$P_M = P(k,N) = \frac{F}{2}k\log(N), \quad M \in \mathscr{G}_k.$$

453 This means that all $M \in \mathscr{G}_k$ have the same penalty term. Then, the best way of segmenting the
454 input stream in $k$ homogeneous segments is given by the model $M_k$ such that

$$M_k = \arg\max_{M \in \mathscr{G}_k}\text{BIC}(M \mid \mathcal{O}) = \arg\max_{M \in \mathscr{G}_k}\log L_M(\mathcal{O}).$$

456 The algorithm for searching the optimum BIC segmentation builds the set
457 $\mathscr{M} = \{M_k : k = 1,\ldots,K\}$, with $K \leqslant N$, of the optimum ways of segmenting the input stream in $k$
458 segments, for all possible $k$s.
459    Let $V_{k,t}$ be the following matrix for the DP

$$V_{k,t} := \max\{\log L_M(o_1,\ldots,o_t) : M \in \mathscr{G}_k\}, \qquad k = 1,\ldots,K, \quad t = 1,\ldots,N.$$

461 The maximum number $K$ of segments for each $t$ can be defined as $K = \lfloor t/S_{\min} \rfloor$, where $S_{\min}$ is the
462 minimum allowed duration (size) of a segment whatever the approach we are using for the seg-
463 mentation. The matrix $V_{k,t}$ is filled column by column by means of the following equations:

$$V_{1,t} = A(o_1,\ldots,o_t), \tag{16}$$

$$V_{k,t} = \max_{(k-1)S_{\min} \leqslant t' \leqslant t-S_{\min}}(V_{k-1,t'} + A(o_{t'+1},\ldots,o_t)), \quad k = 2,\ldots,K, \tag{17}$$

466 where $A(o_a, \ldots, o_b)$ denotes the *auto-consistency* of $o_a, \ldots, o_b$, that is the log-likelihood of $G_a^b$ on
467 $o_a, \ldots, o_b$, where $G_a^b$ is the Gaussian estimated on $o_a, \ldots, o_b$. Note that the range of $t'$ is a function
468 of $S_{\min}$. The main role of the parameter $S_{\min}$ is to ensure smoothed differences between covariance
469 matrices, and then auto-consistencies, of segments that differ for only few observations. Indeed, if
470 we allowed segments of a single observation, the resulting segmentation would strongly depend on
471 noise. However, in the context of our DP approach to segmentation, a careful choice of the
472 parameter $S_{\min}$ can also help in reducing the computational requirements of the algorithm. This
473 issue will be analyzed in Section 5.5.2.
474     The best segmentation with $k$ segments of the input up to time $t$ is obtained by searching,
475 through Eq. (17), the best segmentation in $k - 1$ segments up to $t' < t$, and adding to it the $k$th
476 segment containing the observations from $t' + 1$ to $t$.
477     In addition to $V_{k,t}$, another matrix $M_{k,t}$ is needed for recording the time indexes of change:

$$M_{1,t} = 0,$$

$$M_{k,t} = \arg\max_{(k-1)S_{\min} \leqslant t' \leqslant t - S_{\min}} (V_{k-1,t'} + A(o_{t'+1}, \ldots, o_t)), \quad k = 2, \ldots, K.$$

481 The two matrices $V_{k,t}$ and $M_{k,t}$ have the same structure, and their entries store, respectively, the log-
482 likelihood of the best segmentation in $k$ segments of the input up to the index $t$, and the time index
483 of the last spectral change.
484     When the end of the input is reached, each entry of the last column of $V_{k,t}$ contains the log-
485 likelihood of the best segmentation of the whole input stream in $k$ segments. By subtracting the
486 penalty term, it is possible to obtain $k_{\mathrm{opt}}$, the number of segments of the optimum segmentation

$$k_{\mathrm{opt}} = \arg\max_{k=1,\ldots,K} (V_{k,N} - \lambda P(k,N)) \quad \text{with } K = \lfloor N/S_{\min} \rfloor.$$

488 The optimum segmentation is finally obtained by backtracking over the matrix $M_{k,t}$, starting from
489 the entry $M_{k_{\mathrm{opt}},N}$ and going back to the previous rows, at the columns (time indexes of changes)
490 specified in the entries.

### 5.1. Efficient computation of the auto-consistency

492     The algorithm requires the computation of the auto-consistency of $o_{t'+1}, \ldots, o_t$, which by def-
493 inition can be computed by

$$A(o_{t'+1}, \ldots, o_t) = -\frac{t - t'}{2} \left( \log \left| \Sigma_{t'+1}^t \right| + d \log 2\pi + d \right). \tag{18}$$

495 An efficient way of computing the covariance matrix $\Sigma_{t'+1}^t$ for all possible indexes $t$ and $t'$ is the
496 CSA method described in Section 4.1.3. This method allows the computation of $\Sigma_{t'+1}^t$ with a
497 number of operations equals to $2d(d + 1) + 2d$ (see Table 2), by exploiting the encoding of the
498 input signal with cumulative statistics.

### 5.2. Bounding the auto-consistency

500     Computing $A(o_a, \ldots, o_b)$ is costly since it requires the estimation of a covariance matrix and the
501 computation of its determinant (see Eq. (18)). However, in some cases it is possible to avoid those
502 computations. Let $B(o_a, \ldots, o_b)$ be an upper bound on $A(o_a, \ldots, o_b)$, that is

18                        *M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx*

$$A(o_a, \ldots, o_b) \leqslant B(o_a, \ldots, o_b) \quad \forall a, b; \; a \leqslant b, \tag{19}$$

504  where we assume that the computation of $B()$ is cheaper than that of $A()$. If during the compu-
505  tation of $V_{k,t}$, it happens that for a given $t'$

$$V_{k,t} \geqslant V_{k-1,t'} + B(o_{t'+1}, \ldots, o_t) \quad \forall k,$$

507  then, given Eq. (19), we have

$$V_{k,t} \geqslant V_{k-1,t'} + A(o_{t'+1}, \ldots, o_t) \quad \forall k.$$

509  This means that it is not convenient to hypothesize a change in $t'$, whatever the number of seg-
510  ments $k$; therefore, for that $t'$, the computation of $A(o_{t'+1}, \ldots, o_t)$ can be avoided.
511      The problem is now to define such a bound $B()$. Let $\mu_a^b$ and $\Sigma_a^b$ be the parameters of the
512  Gaussian $G_a^b$ estimated on the $d$-dimensional observations $o_a, \ldots, o_b$ with $a \leqslant b$ and $n = b - a + 1$,
513  and let assume that $A(o_a, \ldots, o_b)$ has already been computed. Let us suppose that our goal is to
514  obtain the bound $B(o_c, \ldots, o_{a-1}, o_a, \ldots, o_b)$ for $A(o_c, \ldots, o_{a-1}, o_a, \ldots, o_b)$. It can be shown that
515  $A(\mu_a^b, \ldots, \mu_a^b, o_a, \ldots, o_b) \leqslant A(o_c, \ldots, o_{a-1}, o_a, \ldots, o_b)$; then we can set the bound $B()$ equal to the
516  auto-consistency of the sequence where the new observations $o_c, \ldots, o_{a-1}$ have been substituted
517  with the mean of the old sub-sequence. Moreover, such a bound can be efficiently computed by

$$B(\underbrace{o_c, \ldots, o_{a-1}}_{m}, \underbrace{o_a, \ldots, o_b}_{n}) = A(\mu_a^b, \ldots, \mu_a^b, o_a, \ldots, o_b)$$

$$= \frac{n+m}{n} A(o_a, \ldots, o_b) - \frac{d(n+m)}{2} \log\left(\frac{n}{n+m}\right)$$

519  based on the the value of the previous auto-consistency and the number $m = a - c$ of the new
520  included observations.

521  *5.3. Further reductions of the algorithm cost*

522      The complexity of the DP algorithm can be further reduced by introducing some reasonable
523  and quite obvious approximations listed below.
524      $K_{\max}$, the maximum number of searched segments. It bounds the number of rows of the DP
525  matrices $V$ and $M$, with a number that is independent from the input audio size $N$.
526      $S_{\max}$, the maximum size of a segment. Actually, this approximation allows to greatly reduce the
527  number of operations of the algorithm, but it can be too dangerous, unless $S_{\max}$ is set to a very
528  high value, thus losing the benefit of its use. Then, to keep the advantage without losing accuracy,
529  the possibility of hypothesizing a segment larger than $S_{\max}$ is kept in few specific circumstances.
530      *Resolution $\delta$*, for encoding the audio stream. Like in the local algorithm described in Section 4,
531  $\delta$ establishes the resolution of the computation of triples $(SQ_1^t, SV_1^t, t)$ and of the entries of ma-
532  trices $V_{k,t}, M_{k,t}$; that is, those values are computed only for $t = \delta, 2\delta, 3\delta, \ldots, N$. The resolution $\delta$
533  reduces the total cost of the algorithm of a factor which is a function of the square of $\delta$.

534  *5.4. Computational costs*

535      The main stages of the DP algorithm without any approximation have the following costs,
536  including the $d^3/6$ operations for the determinant evaluation of the covariance matrices:

YCSLA 259
DISK / 30/6/04

ARTICLE IN PRESS

No. of pages: 24
DTD 4.3.1/SPS

*M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx* 19

537 • input stream encoding: $N(d(d + 1) + d)$,
538 • matrices initialization: $\text{O}\{N \cdot (d^3/6 + 2d(d + 1) + 2d)\}$,
539 • matrices filling: $\text{O}\{N^2/2 \cdot (N/S_{\min} + d^3/6 + 2d(d + 1) + 2d)\}$,
540 • selection of the winner model: $\text{O}\{N/S_{\min}\}$,
541 where the most expensive step is obviously the filling of the DP matrices. The overall complexity
542 of the algorithm is then $\text{O}\{N^3/S_{\min} + N^2 d^3\}$.
543    By introducing the approximations described in Section 5.3, the step costs become:
544 • input stream encoding: $N(d(d + 1) + d)$,
545 • matrices initialization: $\text{O}\{(N/\delta) \cdot (d^3/6 + 2d(d + 1) + 2d)\}$,
546 • matrices filling: $\text{O}\{(NS_{\max}/\delta^2) \cdot (K_{\max} + d^3/6 + 2d(d + 1) + 2d)\}$,
547 • selection of the winner model: $\text{O}\{K_{\max}\}$.
548    Also in this case, the DP step is the most expensive. The complexity of the algorithm is
549 $\text{O}\{(NS_{\max}/\delta^2) \cdot (K_{\max} + d^3)\}$.
550    The benefit given by the bounding $B()$ introduced in Section 5.2 has been experimentally
551 quantified: depending on the value of $S_{\max}$, it allowed 25–36% reduction of the execution time.

### 5.5. Experimental evaluation

#### 5.5.1. Database

554    The main goal of these experiments is to compare the segmentation accuracy of the global and
555 local algorithms. Since it is known that short segments are not well handled by the local algo-
556 rithm, the NIST data are suitable to make evident the gain, if any, given by the global algorithm
557 over the local one. The evaluation has mainly been done on the NIST test set described in Section
558 3.2. Nevertheless, the global algorithm has also been tested on the IBNC test set (Section 3.1), and
559 compared with the local algorithm performance given in Section 4.2.2.

#### 5.5.2. Results

561    First of all, the global algorithm has been compared with the CSA implementation of the local
562 one (see Section 4.1.3) on the NIST test set. The set-ups of the two algorithms are reported in
563 Tables 5 and 6.
564    Table 7 shows both performance and execution times [1] measured for the two algorithms. The
565 improvement in terms of $F$-score given by the global algorithm is 2.4% relative and the cost is 38
566 times higher than the local algorithm. It is worth noticing the benefit in running time (36%) given
567 by the use of the bound $B()$ (Section 5.2).
568    We now present the results of the experiments carried out to measure the impact of the ap-
569 proximations introduced in the exact DP algorithm.
570    Figs. 5–8 plot $F$-score and execution time as functions of $S_{\min}$, $K_{\max}$, $S_{\max}$ and $\delta$, respectively –
571 each computed keeping fixed all the other parameters of the algorithm. In the following con-
572 siderations, the size of the parameters are sometimes in seconds, and not in number of obser-
573 vations as usual: in these experiments 1 s of audio is represented by 100 observation vectors.

---

[1] Experiments were performed on a Pentium III 1 GHz, 1 Gb RAM.

20                       *M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx*

Table 5
Set-up of the global algorithm for the segmentation of the NIST data

| | | |
|---|---|---|
| $S_{min} = 75$ | $S_{max} = 1500$ | $K_{max} = 80$ |
| $\delta = 5$ | $\lambda = 0.60$ | |

Table 6
Set-up of the local (CSA) algorithm for the segmentation of the NIST data

| | | |
|---|---|---|
| $N_{min} = 200$ | $\Delta N_{grow} = 50$ | $\delta_l = 25$ |
| $N_{max} = 500$ | $\Delta N_{shift} = 100$ | $\delta_h = 5$ |
| $N_{second} = 400$ | $N_{margin} = 50$ | $\lambda = 0.85$ |

Table 7
Performance of the local and global algorithms on the NIST test set (0.3 s is the tolerance admitted in the change detection)

| | Performance | | | Execution time (s) | |
|---|---|---|---|---|---|
| | Precision | Recall | *F*-score | | Without $B()$ |
| Local | 54.4 | 75.7 | 63.3 | 367.7 | – |
| Global | 54.7 | 79.7 | 64.8 | 13,930.5 | 21,796.8 |



Fig. 5. *F*-score and execution time as functions of $S_{min}$.

574    $S_{min}$: It results that the value of the minimum window is critical, and the best value is not equal
575 to the minimum size of test segments as expected (half a second, after the merging stage men-
576 tioned in Section 3.2), but is slightly larger (0.75 s). Perhaps, that value could be – somehow –
577 related to the tolerance used in the automatic evaluation procedure (0.3 s).

578    $K_{max}$: The limit on the number of searched segments also affects accuracy when its value is lower
579 than a threshold (about 50). Since the average length of segments in the test set is 1.3 s, and the
580 audio files lasted about 1 min, the expected number of segments in each file is about 45, which is

YCSLA 259
DISK / 30/6/04

ARTICLE IN PRESS

No. of pages: 24
DTD 4.3.1/SPS

*M. Cettolo et al. / Computer Speech and Language xxx (2004) xxx–xxx*          21



Fig. 6. *F*-score and execution time as functions of $K_{max}$.



Fig. 7. *F*-score and execution time as functions of $S_{max}$.

581 compatible with the threshold observed in the experiments. This means that $K_{max}$ must be care-
582 fully chosen, taking into account the characteristics of audio under processing.

583 $S_{max}$: It can be noted that accuracy does not change if $S_{max}$ is reduced from 60 s down to 10 s,
584 halving the execution time. Below 10 s, the number of insertions increases considerably, affecting
585 the overall accuracy. The nature of the test set (only 0.1% of segments is longer than 10 s) allows
586 to reduce $S_{max}$ down to 10 s, but experiments not reported here on a different test set (news
587 programs) show the importance of introducing some corrections in order to relax, at least in some
588 circumstances, the $S_{max}$ constraint.

589 *Resolution δ*: The use of a δ slightly higher than 1 for the processing allows to considerably
590 reduce the amount of time required by the algorithm. Furthermore, $δ = 5$ results in a better
591 performance than the maximum resolution, which causes more false alarms.

592 It is worth pointing out that if $S_{min}$, $K_{max}$, $S_{max}$ and δ are set to proper values, their use does not
593 introduce any degradation of the approximated DP algorithm with respect to the exact one.

22                    *M. Cettolo et al. | Computer Speech and Language xxx (2004) xxx–xxx*



Fig. 8. $F$-score and execution time as functions of $\delta$.

Table 8
Performance of the local (using both diagonal and full covariance matrices) and global algorithms on broadcast news programs (tolerance $= 0.5$ s)

| | Performance | | | Execution time (s) | |
| --- | --- | --- | --- | --- | --- |
| | Precision | Recall | $F$-score | | Without $B()$ |
| Local (diagonal $\Sigma$) | 82.5 | 84.4 | 83.5 | 16.3 | – |
| Local | 89.2 | 89.6 | 89.4 | 42.2 | – |
| Global | 92.9 | 86.8 | 89.8 | 4021.3 | 4771.7 |

For completeness, the global algorithm and the CSA implementation of the local one have also been compared on the IBNC test set. Results are shown in Table 8 (rows `local` and `global`). In this case, the benefit of the global search over the local one is negligible, as it was expected given the audio characteristics of broadcast news recordings. The first row of the table refers to the CSA implementation of the local search with the use of diagonal covariance matrices, instead of full matrices as in all the experiments reported so far. The rationale behind this trial is that it is generally recognized that Cepstral coefficients are reasonably uncorrelated and can be satisfactorily modeled with diagonal-covariance Gaussians, which dramatically reduce computational requirements of our algorithms. Unfortunately, the impressive reduction in time costs (16.3 s vs. 42.2) is paid with a 6.6% relative $F$-score decrease.

## 6. Summary

In this work, three different approaches to the implementation of the well-known local BIC-based audio segmentation algorithm have been beforehand analyzed: (i) a simple method that uses only a sum and a square sum of the input vectors, in order to save computations in estimating covariance matrices on partially shared data; (ii) the approach proposed in Sivakumaran et al. (2001), that encodes the input signal with cumulative distributions; (iii) an innovative approach

that encodes the input signal in cumulative pair of sums. The two latter approaches provide a better use of the typical approximation made in the algorithm, which is the computation of $\Delta$BIC values not on all observations, but at a lower resolution.

The three approaches have been compared both theoretically and experimentally: the proposed new approach has turned out to be the most efficient.

The local algorithm is heuristic, hence the quality of its solutions is only a lower bound on the quality of the solutions achievable by the BIC criterion. In order to discover if there is any room of improvement within or outside the heuristics, we have developed a DP algorithm that, within the BIC model, finds a globally optimal segmentation of the input audio stream. In the paper, it is described, analyzed, and experimentally compared with the local BIC-based algorithm.

The global algorithm yields a small but consistent improvement of performance, with respect to the local one, while its time cost is definitely higher. Its computational complexity was reduced without affecting accuracy, by introducing some reasonable approximations, yielding a less than real time cost.

Summarizing, results show that no much further improvement is possible under the BIC criterion, fact that is important to be aware of. On the other side, experiments make evident that the local algorithm is able to segment an audio stream almost as well as the global algorithm is. That is, the sliding window approach is effective as an heuristics towards the BIC criterion objective function. This encouraged us in proposing improvements to the implementation of the local search, providing a benchmark for possible future experimental work on segmentation.

## Acknowledgements

## References

Akaike, H., 1977. On entropy maximization principle. In: Krishnaiah, P.R. (Ed.), Applications of Statistics. North-Holland, Amsterdam, Netherlands, pp. 27–41.

Baxter, R.A., 1996. Minimum message length inference: theory and applications. Ph.D. Thesis, Department of Computer Science Monash University, Clayton, Victoria, Australia.

Cettolo, M., Federico, M., 2000. Model selection criteria for acoustic segmentation. In: Proceedings of the ISCA Automatic Speech Recognition Workshop, Paris, France.

Cettolo, M., Vescovi, M., 2003. Efficient audio segmentation algorithms based on the BIC. In: Proceedings of the ICASSP, vol. VI, Hong Kong, pp. 537–540.

Cettolo, M., 2000. Segmentation, classification and clustering of an Italian broadcast news corpus. In: Proceedings of the Sixth RIAO – Content-Based Multimedia Information Access – Conference, Paris, France.

Chen, S.S., Gopalakrishnan, P.S., 1998. Speaker, environment and channel change detection and clustering via the Bayesian Information Criterion. In: Proceedings of the DARPA Broadcast News Transcription and Understanding Workshop, Lansdowne, VA.

650  Delacourt, P., Kryze, D., Wellekens, C., 1999. Speaker-based segmentation for audio data indexing. In: Proceedings of
651      the ESCA ETRW Workshop Accessing Information in Spoken Audio, Cambridge, UK.
652  Federico, M., Giordani, D., Coletti, P., 2000. Development and evaluation of an Italian broadcast news corpus. In:
653      Proceedings of the Second International Conference on Language Resources and Evaluation (LREC), Athens,
654      Greece.
655  Frakes, W.B., Baeza-Yates, R., 1992. Information Retrieval: Data Structures and Algorithms. Prenctice-Hall,
656      Englewood Cliffs, NJ.
657  Gauvain, J.-L., Lamel, L., Adda, G., 1998. Partitioning and transcription of broadcast news data. In: Proceedings of
658      the ICSLP, Sidney, Australia, pp. 1335–1338.
659  Gish, H., Siu, M., Rohlicek, R., 1991. Segregation of speakers for speech recognition and speaker identification. In:
660      Proceedings of the ICASSP, vol. II, Toronto, Canada, pp. 873–876.
661  Hain, T., Johnson, S.E., Tuerk, A., Woodland, P.C., Young, S.J., 1998. Segment generation and clustering in the HTK
662      broadcast news transcription system. In: Proceedings of the DARPA Broadcast News Transcription and
663      Understanding Workshop, Lansdowne, VA.
664  Harris, M., Aubert, X., Haeb-Umbach, R., Beyerlein, P., 1999. A study of broadcast news audio stream segmentation
665      and segment clustering. In: Proceedings of the EUROSPEECH, vol. III, Budapest, Hungary, pp. 1027–1030.
666  IBNC, 2000. Available from: <http://www.elda.fr/catalogue/en/speech/S0093.html>.
667  Kemp, T., Schmidt, M., Westphal, M., Waibel, A., 2000. Strategies for automatic segmentation of audio data. In:
668      Proceedings of the ICASSP, vol. III, Istanbul, Turkey, pp. 1423–1426.
669  Lu, L., Li, S.Z., Zhang, H.-J., 2001. Content-based audio segmentation using support vector machines. In: Proceedings
670      of the ICME, Tokyo, Japan, pp. 956–959.
671  NIST, 2000. Available from: <www.nist.gov/speech/tests/spk/2000/>.
672  Scheirer, E., Slaney, M., 1997. Construction and evaluation of a robust multifeature speech/music discriminator. In:
673      Proceedings of the ICASSP, Munich, Germany, pp. 1331–1334.
674  Schwarz, G., 1978. Estimating the dimension of a model. The Annals of Statistics 6 (2), 461–464.
675  Seber, G.A.F., 1984. Multivariate Observations. John Wiley & Sons, New York, NY.
676  Siegler, M.A., Jain, U., Raj, B., Stern, R.M., 1997. Automatic segmentation, classification and clustering of broadcast
677      news audio. In: Proceedings of the DARPA Speech Recognition Workshop, Chantilly, VA.
678  Sivakumaran, P., Fortuna, J., Ariyaeeinia, A.M., 2001. On the use of the Bayesian Information Criterion in multiple
679      speaker detection. In: Proceedings of the EUROSPEECH, vol. II, Aalborg, Denmark, pp. 795–798.
680  Tritschler, A., Gopinath, R., 1999. Improved speaker segmentation and segments clustering using the Bayesian
681      Information Criterion. In: Proceedings of the EUROSPEECH, vol. II, Budapest, Hungary, pp. 679–682.
682  Vescovi, M., Cettolo, M., Rizzi, R., 2003. A DP algorithm for speaker change detection. In: Proceedings of the
683      EUROSPEECH, vol. IV, Geneva, Switzerland, pp. 2997–3000.
684  Wegmann, S., Zhan, P., Gillick, L., 1999. Progress in broadcast news transcription at Dragon systems. In: Proceedings
685      of the ICASSP, Phoenix, AZ, pp. 33–36.
686  Wellekens, C., 2001. Seamless navigation in audio files. In: Proceedings of the ITRW on Speaker Recognition, Crete,
687      Greece, pp. 9–12.