

Optimizing Bicoid Signal Extraction

Hossein Hassani* Emmanuel Sirimal Silva† Zara Ghodsi‡

November 15, 2017

Abstract

Signal extraction and analysis is of great importance, not only in fields such as economics and meteorology, but also in genetics and even biomedicine. There exists a range of parametric and nonparametric techniques which can perform signal extractions. However, the aim of this paper is to define a new approach for optimising signal extraction from bicoid gene expression profile. Having studied both parametric and non-parametric signal extraction techniques, we identified the lack of specific criteria enabling users to select the optimal signal extraction parameters. Exploiting the expression profile of *bicoid* gene, which is a maternal segmentation coordinate gene found in *Drosophila melanogaster*, we introduce a new approach for optimising the signal extraction using a nonparametric technique. The underlying criteria are based on the distribution of the residual, more specifically its skewness.

Keywords: Signal extraction; Optimisation; Residual distribution; Bicoid.

1 Introduction

Signal extraction is an important and challenging task in the field of time series analysis and forecasting. Signals can take various forms with the most common being trends and seasonal fluctuations. Trend extraction in particular enables analysts to smooth out a time series and remove the seasonal and cyclical variations - thereby enabling the determination of the long-run behaviour of the underlying data. A trend is formally defined as a smooth additive component which contains information relating to the global change in a time series [1], and the term ‘smooth’ is a vital characteristic of any given signal. In the field of genetics and gene expression studies, signal extraction and noise reduction are crucial as genetic data is often characterised by the existence of considerable noise [2].

*Research Institute of Energy Management and Planning, University of Tehran, No. 13, Ghods St., Enghelab Ave., Tehran, Iran, email: hassani.stat@gmail.com.

†Fashion Business School, London College of Fashion, University of the Arts London, 272 High Holborn, London, WC1V 7EY, UK, email: e.silva@fashion.arts.ac.uk.

‡Translational Genetics Group, Bournemouth University, Fern Barrow, Poole, BH125BB, UK, e-mail: zghodsi@bournemouth.ac.uk.

32 Our interest in this topic is motivated by the findings in [2] where the au-
 33 thors evaluated a variety of parametric and nonparametric signal processing
 34 techniques for extracting the signal in bicoid (*bcd*)¹, which is a morphogen lo-
 35 calised at the anterior end of the egg. After fertilisation, the distribution of
 36 Bcd along the embryo (the signal under study in this paper) determines the
 37 cell’s destiny in a concentration-dependent mode. Here, the authors found that
 38 a nonparametric approach produced the most efficient extraction of the Bcd
 39 signal [2].As Ghodsi et al. [2] point out, the Bcd signal extraction process is
 40 complex because the data associates with both observational and biological
 41 noise, and the extracted residual is not normally distributed as required by
 42 parametric techniques. Figure 1 below shows an example of a typical noisy
 43 Bcd. As noted in [2], the distribution of Bcd follows an exponential trend, and
 44 the high volatility seen in the profile ensures that the extraction of this signal
 45 remains an arduous task.

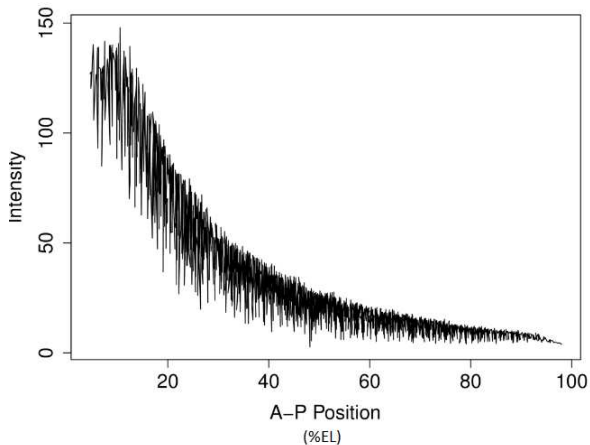


Figure 1: A typical example of noisy Bcd [3].

46 Accordingly, the aim of this paper is to introduce and define a new approach
 47 for optimising Bcd signal extraction. At present, there exists no definitive crite-
 48 rion to aid researchers and scientists interested in extracting the Bcd signal for
 49 analysis. Since the Bcd signal defines what positional information is available
 50 for morphogen readout, studying the characteristics of this signal can improve
 51 our knowledge on several critical developmental processes such as embryogen-
 52 esis, regional specification and canalisation. It is noteworthy that the criteria
 53 which follow are tailored for the sole purpose of extracting an accurate Bcd sig-
 54 nal based on the knowledge disseminated through the work in [2] with regard to
 55 the distribution of the residual following Bcd signal extraction. Therefore, the
 56 criteria presented herewith may not be directly suitable for other applications.

57 In addition to covering the main aim of this paper, we also present readers
 58 with two other interesting concepts related to Bcd expression profile. These are

¹In what follows, the italic lower-case *bcd* presents either the gene or the mRNA and Bcd refers to the protein.

59 sequential and hybrid signal extraction processes which are explained in Section
60 2. Accordingly, this paper is able to present readers with three different
61 approaches for Bcd signal extraction based on their requirements and interests.
62 The first approach is suitable for those who wish to rely on a single model for
63 Bcd signal extraction. We have tailored the criteria presented in this paper
64 to enable a swift and accurate Bcd signal extraction using the nonparametric
65 approach identified as best in [2]. Should the extracted signal appear to have
66 captured some unnecessary fluctuations, then the sequential process described
67 can be applied on the original signal to generate a refined and smoother signal.
68 Even though the findings in [2] suggests that the Bcd residual is skewed, we
69 appreciate that statisticians who subscribe to classical methods would find it
70 difficult to agree with such outcomes. Therefore, as a second approach, we
71 propose a hybrid parametric signal extraction process which can ensure that
72 the residual is in fact white noise. Finally, for those who wish to exploit hybrid
73 modelling from a purely nonparametric perspective with the possibility of capturing
74 the maximum variation via a smooth signal line, we present the hybrid
75 nonparametric approach and show that it can produce far better results when
76 combined with the optimized signal extraction criteria presented herewith. The
77 above three approaches also represent the core contributions of this research.

78 The remainder of this paper is organised such that Section 2 focuses on
79 optimising Bcd signal extraction with Section 3 presenting the empirical results.
80 This is followed by an interesting discussion in Section 4, and the paper
81 concludes in Section 5.

82 2 Optimising Bcd Signal Extraction

83 2.1 Singular Spectrum Analysis

84 The Singular Spectrum Analysis (SSA) technique is a nonparametric filtering
85 technique that is dependent upon its choice of Window Length L and the
86 number of eigenvalues r [4]. SSA was successfully introduced for Bcd signal
87 extraction in [5] and exploited in more detail in [2]. In [2] it was found that
88 the residual following signal extraction in Bcd is not normally distributed or
89 stationary, and also showed that the residual itself has a complex pattern which
90 adds further to the difficulty in smoothing and signal extraction. However, SSA
91 is unique as it can extract several signals for any given time series depending on
92 the chosen value of L . In fact, the choice could be any L such that $2 \leq L \leq N/2$,
93 where N is the length of the series. As such, the findings in [2] which show SSA
94 as the best option for Bcd signal extraction (in relation to Synthesis Diffusion
95 Degradation, Exponential Smoothing, Autoregressive Integrated Moving Average
96 (ARIMA), Fractionalized ARIMA, and Neural Networks) falls short of
97 defining the optimal SSA choices for Bcd signal extraction.

98 Through our work, we intend to fill this gap by introducing new criteria
99 which enables optimisation of the Bcd signal extraction process with SSA. The
100 importance of defining such criteria is further evidenced by the fact that SSA
101 has been applied for extracting the Bcd and other segmentation gene's signal
102 since 2006, see for example [2, 5–10]. Therefore, it is clear that researchers

103 and scientists alike can benefit from some formal criteria for the selection of
 104 SSA choices for Bcd signal extraction. Whilst the remainder of this paper
 105 focuses entirely on SSA, we find it pertinent to acknowledge and comment on
 106 the comparative preferability of SSA over other filtering techniques such as
 107 Hilbert-Huang (HH) [11] and Hodrick-Prescott (HP) [12]. Firstly, the SSA
 108 technique (as detailed below) is a Singular Value Decomposition based method
 109 and as such is very effective for noise reduction [13]. Secondly, the HH approach
 110 is closely associated with Empirical Mode Decomposition which is related to
 111 the setting of intrinsic mode functions. Thirdly, the signal process in the HP
 112 filtering approach has two instead of one unit root and is therefore most suitable
 113 for time series with two unit roots [13]. Finally, a direct comparison of both
 114 SSA and HP under equal conditions showed that SSA performs on par with
 115 the HP filter [13].

116 The basic SSA technique consists of two complementary stages referred to
 117 as decomposition and reconstruction, and each of these stages includes two
 118 separate steps [4]. In brief, at the first stage, Bcd is decomposed into the
 119 sum of a small number of independent and interpretable components such as
 120 a slowly varying trend and a structureless noise [2, 4], and at the second stage
 121 the noise free Bcd is reconstructed [2, 14]. It should be noted that the use of
 122 SSA here is for the sole purpose of obtaining the optimal decomposition of Bcd
 123 and then extracting the signal component. Figure 2 summarises the basic SSA
 124 process as a flowchart.

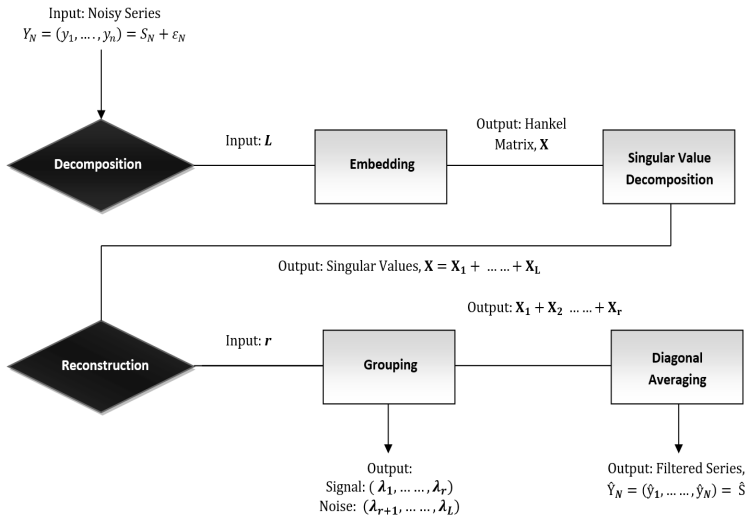


Figure 2: A flowchart of the basic SSA process [4].

125 A more detailed explanation of the steps underlying SSA for bicoid signal
 126 extraction is provided below, and in doing so we mainly follow [2, 4].

127 The first step maps a one dimensional time series $Y_N = (y_1, \dots, y_N)$ into
 128 a multi-dimensional series X_1, \dots, X_K with vectors $X_i = (y_i, \dots, y_{i+L-1})^T \in$
 129 \mathbf{R}^L , where $K = N - L + 1$. Whilst the process itself is referred to as embedding,

130 the vectors X_i are called *L-lagged vectors*. The single choice of the embedding
 131 stage is the Window Length L , which is an integer such that $2 \leq L \leq N/2$.
 132 This step results in the trajectory matrix \mathbf{X} , which is also a Hankel matrix and
 133 takes the form: $\mathbf{X} = [X_1, \dots, X_K] = (x_{ij})_{i,j=1}^{L,K}$.

134 Thereafter, we obtain the singular value decomposition (SVD) of the trajec-
 135 tory matrix and represent it as a sum of rank-one bi-orthogonal elementary
 136 matrices. The eigenvalues of $\mathbf{X}\mathbf{X}^T$ are denoted by $\lambda_1, \dots, \lambda_L$ in decreasing
 137 order of magnitude ($\lambda_1 \geq \dots \lambda_L \geq 0$) and by U_1, \dots, U_L the orthonormal system.
 138 Then, we set

$$d = \max(i, \text{ such that } \lambda_i > 0) = \text{rank } \mathbf{X}.$$

139 If we denote $V_i = \mathbf{X}^T U_i / \sqrt{\lambda_i}$, then the SVD of the trajectory matrix can be
 140 written as:

$$\mathbf{X} = \mathbf{X}_1 + \dots + \mathbf{X}_d, \quad (1)$$

141 where $\mathbf{X}_i = \sqrt{\lambda_i} U_i V_i^T$ ($i = 1, \dots, d$). The matrices \mathbf{X}_i are elementary matrices
 142 as they have rank 1, U_i and V_i denotes the left and right eigenvectors of the
 143 trajectory matrix. The collection $(\sqrt{\lambda_i}, U_i, V_i)$ is called the i -th eigentriple of
 144 the matrix \mathbf{X} , $\sqrt{\lambda_i}$ ($i = 1, \dots, d$) are the singular values of the matrix \mathbf{X} and
 145 the set $\{\sqrt{\lambda_i}\}$ is called the spectrum of the matrix \mathbf{X} . The expansion (1) is
 146 said to be uniquely defined if all the eigenvalues have a multiplicity of one.
 147 The process of splitting the elementary matrices \mathbf{X}_i into several groups and
 148 summing the matrices within each group is called grouping, and transfusing
 149 each resultant matrix from the grouping step to a less noisy series is called
 150 diagonal averaging.

151 As specifically noted in [2], when using SSA, in general the first eigenvalue
 152 corresponds to the trend of a given time series. In order to illustrate this
 153 more clearly to the reader, we show a couple of examples in Figures 3 and 4.
 154 Moreover, in [14, 15] the authors extract and illustrate the trend for tourist
 155 arrivals using SSA based decomposition and the first eigenvalue. Thus, we
 156 extract the first eigenvalue alone and consider the remainder as noise, and
 157 then perform diagonal averaging to transform the matrix containing the first
 158 eigenvalue into a series which will now provide the extracted signal from Bcd.

159 2.1.1 New Approach for Optimising Bcd Signal with SSA

160 In this section we present the new approach for optimising Bcd signal extrac-
 161 tion with SSA and provide justification for the process. The proposed criteria
 162 are developed as follows.

163
 164 **1)** The extracted Bcd trend must be smooth. This is in accordance with the
 165 widely accepted definition of a trend which states that it must be a ‘smooth’
 166 additive component [1].

167
 168 **2)** Setting L sufficiently large enables the first eigenvalue, i.e. $r = 1$ (in some
 169 cases, $r = 1, 2$) to extract a smooth signal for a given series. However, the
 170 value of L must not be too small or too large. By theory, L must lie between

171 $2 \leq L \leq N/2$ [4]. Yet, when it comes to Bcd signal extraction, setting L at
 172 $N/2$ can have negative implications, as with setting L too small.

173 For example, let us first consider the scenario in Figure 3 whereby in a series
 174 with length 301, we consider SSA choices of $L = 2$ and $r = 1$ for Bcd signal
 175 extraction. Notice how the extracted signal fails to meet the ‘smooth’ criteria
 176 as per the definition of a signal in [1]. Accordingly, it is evident that setting L
 177 too small fails to achieve an optimal signal extraction with SSA for Bcd.

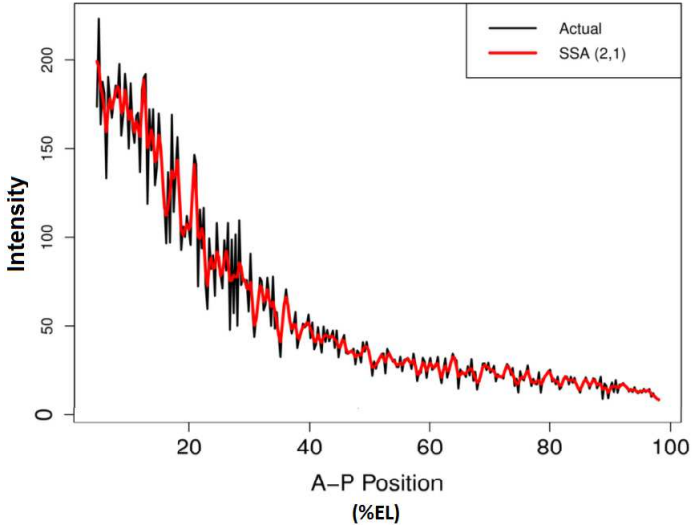


Figure 3: signal extraction from noisy Bcd with SSA choices of $L = 2$ and $r = 1$.

178 Secondly, let us consider what happens when we set L too large for the
 179 same data set. Here, the maximum possible value of L is 150. As such, we set
 180 $L = 150$ and seek to extract the signal in our data. Figure 4 shows the resulting
 181 outcome. In this case, notice how the signal line is smooth (confirming that
 182 setting L large can provide a smoother line) but the extracted signal fails to fit
 183 well to the actual data, especially towards the tail of the series.

184 **3)** Based on points 1) and 2), we suggest the following threshold for the selec-
 185 tion of L for Bcd signal extraction purposes. The window length L should be
 186 some value between $10 \leq L \leq N/4$. Whilst this assumption helps restrict the
 187 selection of L , on its own it fails to provide the researcher with an exact value
 188 for L . Therefore, we call upon the nonparametric nature of SSA to provide the
 189 final closing argument for the criteria.

190

191 **4)** As a nonparametric technique, the SSA residual can be skewed. Based on the
 192 findings in [2] which was an extensive study into signal extraction in Bcd, the
 193 residual from the process was in fact found to be skewed. As such, we propose
 194 using the skewness statistic as an indicator, and finding L which corresponds
 195 to the minimum skewness for a given Bcd series within the threshold $10 \leq L \leq$
 196 $N/4$ and coupling this with $r = 1$ or $r = 1, 2$ as appropriate for optimal Bcd
 197 signal extraction with SSA.

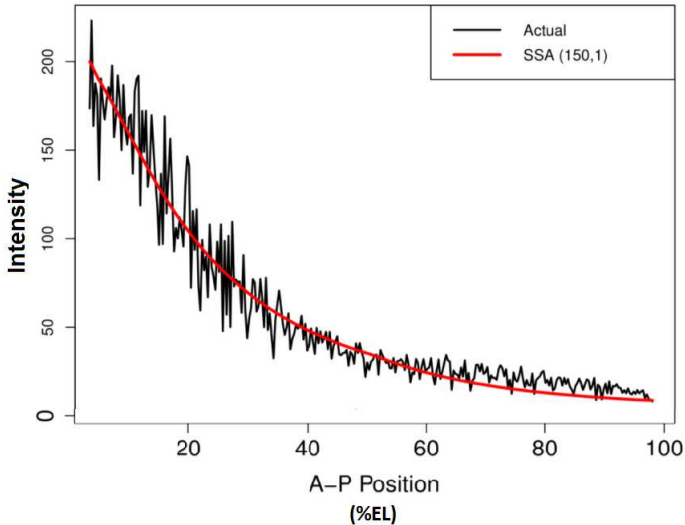


Figure 4: signal extraction from noisy Bcd with SSA choices of $L = 150$ and $r = 1$.

2.2 Sequential and Hybrid Signal Extraction

Section 4 in this paper is dedicated to a discussion which focuses on the exploitation of Sequential SSA and a hybrid signal extraction process for Bcd signal extraction. In what follows we present the ideas that are evaluated later on with empirical data.

2.2.1 Nonparametric Approach

Signal extraction in Bcd data can be an arduous task owing to the complex structure portrayed by the data [2]. Sequential SSA is a relatively new concept which is of great benefit when faced with weak separability between signal and noise as a result of such complexities. For example, when faced with problems in separating a signal of complex form and seasonality, Sequential SSA can be exploited to obtain a more accurate decomposition from the residual after signal extraction [16]. Whilst historically, Sequential SSA was performed on a residual, in this paper we suggest the use of Sequential SSA for refining the Bcd signal further.

The basic idea underlying Sequential SSA is to perform a second round of SSA based decomposition and reconstruction on data that has already undergone an initial round of SSA, with the aim to refine the signal of interest further. Suppose that we exploit the optimised Bcd signal extraction algorithm explained above and extract some signal line. However, if the Bcd data in question has a highly complex structure, it is possible to end up with a signal line that is not as smooth as required. In such instances, we suggest exploiting Sequential SSA, not on the residual, but on the extracted signal to smooth it further and obtain a new and refined signal curve. This approach is greatly

222 beneficial to those who wish to rely on a single model for Bcd signal extraction
223 and enjoy the benefits of using a nonparametric technique.

224 **2.2.2 Hybrid Signal Extraction with SSA**

225 It is possible that some statisticians may not be convinced or used to subspace-
226 based methods such as SSA. Therefore, we find it pertinent to present the
227 possibility of obtaining a hybrid signal extraction process which will combine
228 the optimised SSA signal extraction algorithm for Bcd with other automated
229 signal processing techniques from both parametric and nonparametric back-
230 grounds.

231 The basic idea underlying the hybrid signal extraction process is as follows:

- 232 1. Extract the Bcd signal via the optimised SSA signal extraction algorithm.
233
- 234 2. Fit a different time series model to the residuals following SSA signal
235 extraction and obtain the fitted values.
236
- 237 3. Add the fitted values to the original SSA signal to create the Hybrid SSA
238 signal.
239

240 **2.2.2.1 Hybrid SSA Signal: Parametric Approach** The idea underly-
241 ing the hybrid SSA signal with a parametric approach is to combine the non-
242 parametric SSA signal with the fitted values on residuals from a parametric
243 signal processing model. As most classical statisticians welcome and subscribe
244 to the ARIMA model, here we choose an automated ARIMA model as pro-
245 vided via the forecast package in R [17]. It is important to note that in this
246 paper, we do not rely on ARIMA for its forecasting capabilities. Instead, we
247 consider ARIMA as a tool for extracting any hidden signals within the residual
248 following the initial filtering with SSA. This in turn enables one to ensure that
249 the residual is indeed white noise, as required by parametric models.

250 The modelling equations for ARIMA relevant to this study can be described
251 by following [18]. A non-seasonal ARIMA model may be written as:

$$(1 - \phi_1 B - \dots - \phi_p B^p)(1 - B)^d y_t = c + (1 + \theta_1 B + \dots + \theta_q B^q) e_t, \quad (2)$$

252 where B is the backshift operator, c is a constant, p is the order of the au-
253 toregressive part, q is the degree of first differencing, d is the order of the
254 moving average part of the model, and e_t is white noise [18]. In the R software,
255 the inclusion of a constant in a non-stationary ARIMA model is equivalent to
256 inducing a polynomial signal of order d in the forecast function.

257 **2.2.2.2 Hybrid SSA Signal: Nonparametric Approach** Whilst the
258 underlying idea remains the same, in this instance, as opposed to relying on a
259 parametric time series analysis model, we can combine the nonparametric and

260 optimised SSA signal with fitted values on residuals from a nonparametric time
261 series analysis model in order to obtain the hybrid SSA signal. The benefits of
262 this approach would be that it enables to overcome the parametric restrictions
263 of normality and stationarity of residuals of which the former condition was
264 found to be irrelevant in the case of Bcd data where the residual following signal
265 extraction is skewed [2]. In this case, we rely on the automated Exponential
266 Smoothing (ETS) model found in the forecast package in R. Those interested
267 in the several ETS formula's that are evaluated through the forecast package
268 when selecting the best model to fit the residuals are referred to Chapter 7,
269 Table 7.8 in [18].

270 **3 Empirical Results**

271 **3.1 Data**

272 The evaluation in this study is performed on 17 *Drosophila melanogaster* em-
273 bryos introduced by Alexandrov et al. [10] which was originally obtained from
274 FlyEx database [19,20]. This dataset has been widely used as a valuable source
275 of information for studying the dynamics of segment determination of early
276 *Drosophila development* [21].

277 In FlyEx, the quantitative Bcd data was obtained using the confocal scanning
278 microscopy of fixed embryos immunostained for segmentation proteins [20]. To
279 that aim, A 1024x1024 pixel confocal image with 8 bits of fluorescence data
280 was achieved for each embryo which then transformed into an ASCII table.
281 The ASCII table contains the fluorescence intensity levels attributed to each
282 nucleus of A-P axis. To present the data using a graph, the x-axis shows the
283 anterior to a posterior position along the length of the egg expressed as the
284 percentage, and the y-axis shows the intensity levels which correspond to the
285 amount of expressed *bcd* gene.

286 It is of note that in the study conducted by Alexandrov et al., the out of
287 focus regions were removed by excluding the utmost anterior and posterior
288 areas. After removing the upper and lower values, to get a complete profile
289 along the A-P axis of the embryo, a curve was fitted to the interval of the
290 A-P coordinate between 20 and 80% of egg length (a complete explanation of
291 the method and biological characteristics of this data can be found in [10,22]).
292 However, to introduce a signal processing method capable of both noise filtering
293 and signal extraction, this paper considers the whole data which is unprocessed
294 for any noise reduction methods.

295 **3.2 Signal Extraction**

296 Here, we consider real Bcd data and seek to extract the signal with SSA using
297 the newly proposed criteria as outlined in Section 2.1. Figure 5 below portrays
298 a selection of the actual data and extracted signal with the optimized SSA
299 algorithm, and also outlines the SSA choices which have been used in each
300 case. For the examples in Figure 5, note how the extracted signal is not only
301 smooth, but also well centred around the data, thereby providing the reader

302 with a very accurate outlook for the long term prospects of the Bcd gradient.
303 However, it is evident that on its own, SSA appears to have difficulties in
304 accurately capturing the signal curve initially when it is faced with very high
305 levels of fluctuations as clearly visible within the first few observations of the
306 Bcd profile. We consider this aspect further in the discussion which follows in
307 Section 4.

308 Even though signal extraction is the primary focus of this study, it is no
309 secret that the residual can often enlighten us to crucial information pertaining
310 to any given data set. As such, we follow up the signal extractions with a sound
311 residual analysis.

312 **3.3 Residual Analysis**

313 In order to save space, via Figure 6 we only show the residuals corresponding
314 to the signal extractions shown in Figure 5. A first look at the structure
315 and distribution of the residual over time helps us understand the difficulty in
316 extracting the signal from Bcd profiles. This is largely to do with the highly
317 volatile nature of the data which results in fluctuating amplitudes over time in a
318 particular pattern. In fact, the general patterns appears such that all residuals
319 portray amplitudes which are initially high and then gradually decrease. This
320 in turn means that the techniques adopted for Bcd signal extraction should
321 be able to cope well with such variation and fluctuations in data if it is to
322 accurately perform its task. Moreover, it appears to the naked eye that there
323 is indeed some signal contained within these residuals. Whilst it is expected
324 that a residual following trend extraction would result in capturing the other
325 signals, in some instances there also appears to be a small trend pattern hidden
326 within this data.

327 However, as visual inspections fall short of providing sound evidence, we
328 also consider some statistics for analysing the residuals further. These are
329 reported via Table 1 for all the Bcd data considered in this study. The residuals
330 are initially tested for normality via the Kolmogorov-Smirnov (KS) test for
331 normality. The choice of KS test as opposed to using the popular Shapiro-Wilk
332 (SW) test for normality was because when faced with large samples the KS test
333 is likely to be comparatively more accurate than the SW test [23]. As expected,
334 all residuals failed to pass the normality test reporting probability values of
335 less than 0.001, and thereby leading to a rejection of the null hypothesis of
336 normality. This lets us conclude with 99% confidence that the Bcd residuals
337 following signal extraction are in fact skewed and these results are consistent
338 with the findings in [2].

339 Finally, we go a step further and fit optimal ARIMA models [18] to the
340 residuals. This was done in order to ascertain the randomness of the residu-
341 als following Bcd signal extraction with optimised SSA. Statisticians who rely
342 on classical signal extraction techniques would be overly concerned with the
343 parametric assumptions of normality and stationarity of the residuals. Whilst
344 we have assessed the normality of residuals via the KS test and justified based
345 on [2] that the residuals from this signal extraction exercise should be skewed,
346 fitting of optimal ARIMA models enables us to easily show whether the resid-
347 uals meet the stationary criteria. We fit automated and optimised ARIMA

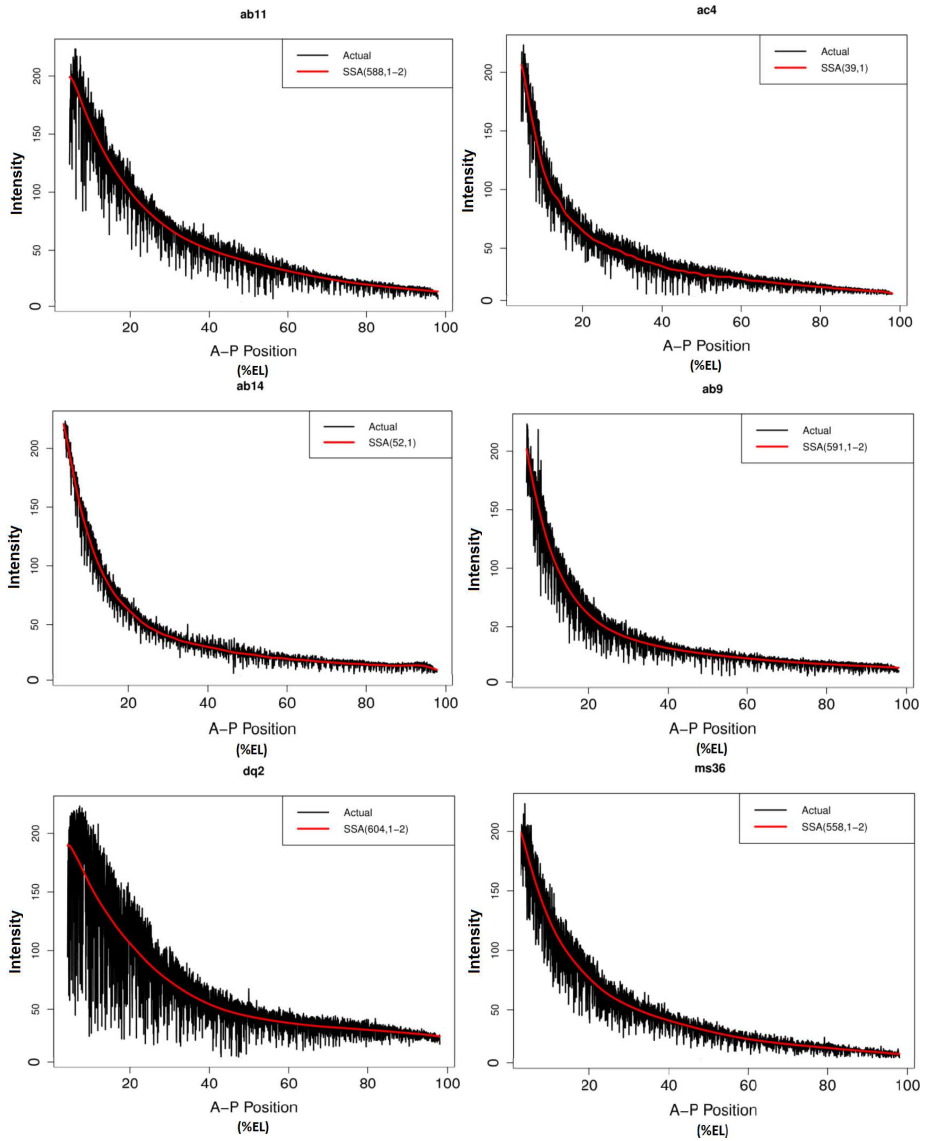


Figure 5: Optimised signal extraction with SSA for a selection of Bcd data.

348 models (as provided via the forecast package in R) on the residuals and report
349 the outcomes in Table 1. A non-seasonal ARIMA model is represented in the
350 form $ARIMA(p, d, q)$ where p indicates the order of the autoregressive parts, d
351 the degree of first differencing and q the order of the moving average part of
352 the model [18]. If the data is non-stationary, then within the $ARIMA(p, d, q)$
353 process the value of $d \geq 1$. If the data is stationary, then no differencing is
354 required, and so $d = 0$. In this case, we notice that $d = 0$ in all instances, and
355 thereby proves that the residuals are indeed stationary.

356 However, the fitting of ARIMA models on the residuals also highlight another
357 interesting point. Notice how for 27 Bcd residuals there have been a variety
358 of 14 different ARIMA models which have been fitted. This in turn indicates
359 the complexity and difficulty associated with the selection of a single technique
360 for extracting Bcd signal, and most certainly highlights the difficulties which
361 any technique would experience when seeking to extract a signal from data
362 with such complex fluctuations. In addition, except for where the model reads
363 $ARIMA(0, 0, 0)$, in all other instances we notice that the residuals are not white
364 noise. We discuss this, and provide a possible solution within the discussion.

Table 1: Residual analysis for Bcd signal extractions.

Embryo	n	SW	ARIMA
ab2	138	<0.001	ARIMA(0,0,1) with zero mean
hz15	85	<0.001	ARIMA(0,0,0) with zero mean
hz28	79	<0.001	ARIMA(2,0,2) with zero mean
ad14	301	<0.001	ARIMA(2,0,5) with zero mean
ad22	294	<0.001	ARIMA(4,0,3) with zero mean
ad23	308	<0.001	ARIMA(1,0,3) with non-zero mean
ab17	485	<0.001	ARIMA(1,0,3) with non-zero mean
ad4	556	<0.001	ARIMA(4,0,4) with zero mean
ad6	566	<0.001	ARIMA(2,0,2) with non-zero mean
ab12	2284	<0.001	ARIMA(4,0,2) with zero mean
ab10	2263	<0.001	ARIMA(1,0,2) with zero mean
ac5	2404	<0.001	ARIMA(4,0,4) with non-zero mean
ab1	2570	<0.001	ARIMA(4,0,4) with zero mean
ac7	2268	<0.001	ARIMA(1,0,2) with zero mean
ad13	2235	<0.001	ARIMA(4,0,2) with non-zero mean
ad29	2193	<0.001	ARIMA(1,0,2) with zero mean
ad32	2183	<0.001	ARIMA(2,0,1) with zero mean
ab7	2346	<0.001	ARIMA(1,0,2) with zero mean
ac3	2356	<0.001	ARIMA(0,0,1) with zero mean
ac9	2215	<0.001	ARIMA(4,0,1) with zero mean
ms14	2305	<0.001	ARIMA(4,0,2) with zero mean
ab11	2355	<0.001	ARIMA(4,0,2) with zero mean
ac4	2383	<0.001	ARIMA(3,0,1) with zero mean
ab14	2218	<0.001	ARIMA(1,0,2) with zero mean
ab9	2369	<0.001	ARIMA(2,0,1) with zero mean
dq2	2423	<0.001	ARIMA(2,0,4) with zero mean
ms36	2239	<0.001	ARIMA(5,0,1) with zero mean

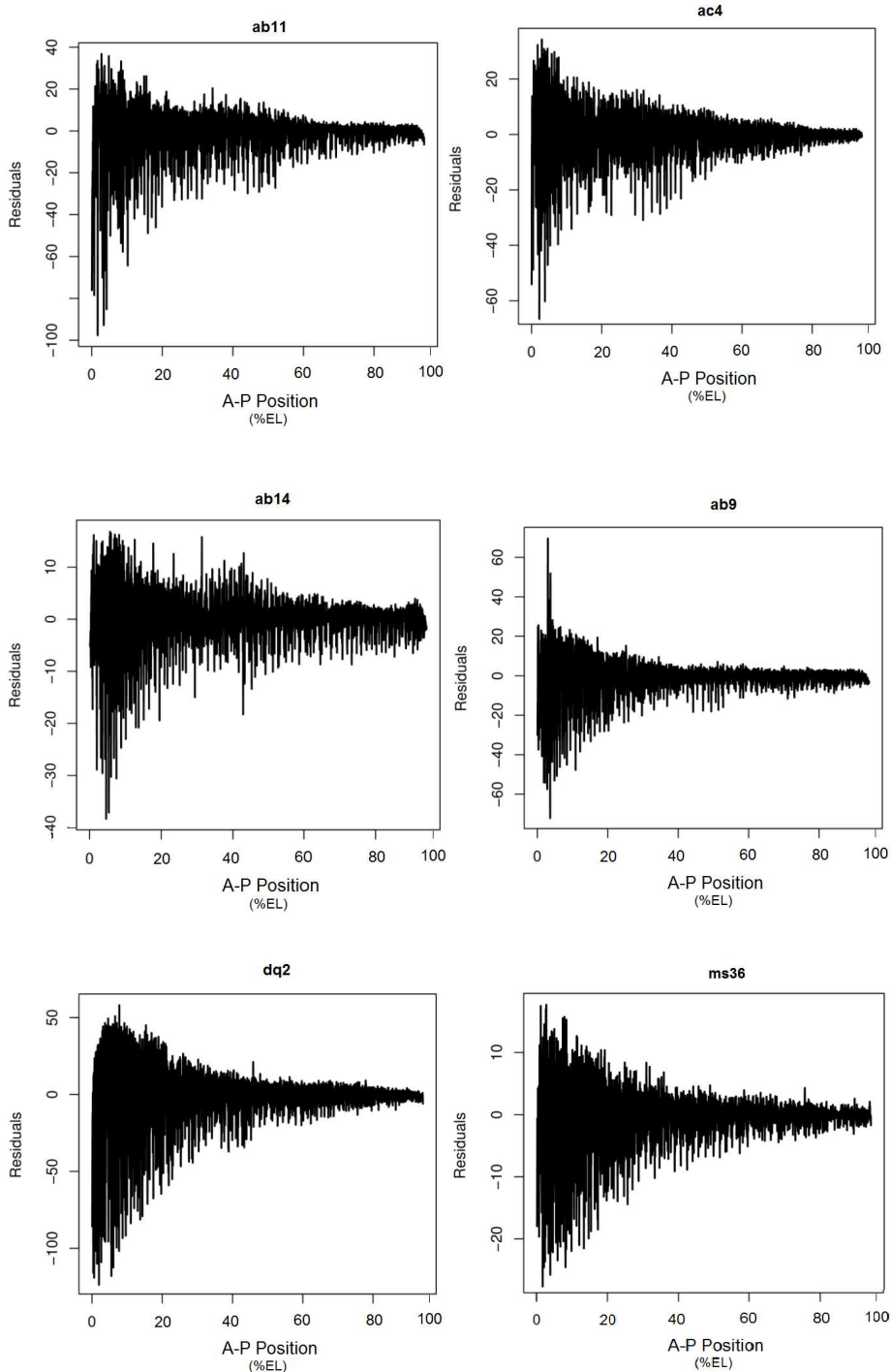


Figure 6: Residuals following optimised signal extraction with SSA for a selection of Bcd data.

4 Discussion

4.1 Sequential SSA on Bcd signal

Note how the signal extraction in ac3, Figure 7, appears to have captured some other fluctuations apart from the signal alone. As such, this extraction, in particular, fails to meet our criteria for a smooth signal. When faced with such situations, we are able to find a solution via sequential SSA. Sequential SSA enables users to take the extracted signal (the signal in our example) and filter same with SSA once more to obtain a more refined output. In what follows we have applied Sequential SSA on the initially extracted Bcd signal.

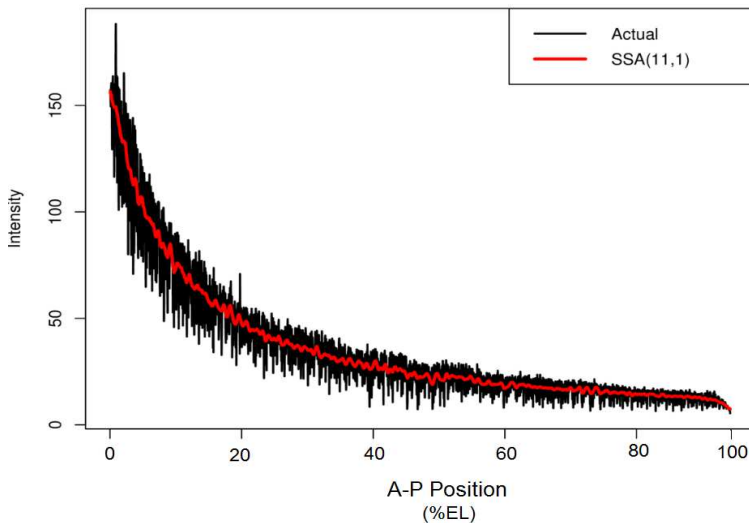


Figure 7: SSA based optimal trend extraction for ac3.

As visible via Figure 8, following sequential SSA we have been able to extract a smoother signal. In this instance, we used the signal extracted via the optimised SSA signal extraction algorithm for Bcd and refined this signal further via Sequential SSA. Here we have used $L = N/2$ and $r = 1$ for signal extraction with Sequential SSA. In line with good practice, the residual was once again tested for normality via the KS test which indicated that the residual is skewed at a 1% significance level, and fitting of an ARIMA model showed that the residual is stationary as well.

4.2 Hybrid SSA Signal Extraction for Bicoid

4.2.1 Hybrid SSA Signal: Parametric Approach

The residual analysis in Table 1 indicates that ARIMA models could be fitted to all but one of the residuals following signal extraction with the optimised SSA signal algorithm. This means that only one of the residuals are pure white noise as it stands. Whilst some might argue that this is acceptable given that

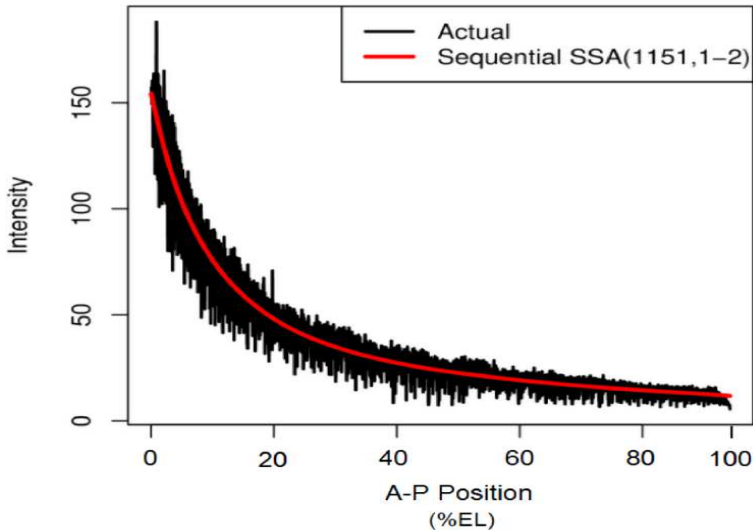


Figure 8: Refined signal extraction with sequential SSA on ac3 signal.

388 the objective is to extract the signal component alone, there may be others who
 389 subscribe to an alternate view along the lines of obtaining a random residual
 390 following signal extraction. The first hybrid SSA signal approach we present is
 391 one which enables users who wish to obtain white noise to achieve this following
 392 Bcd signal extraction with SSA. We begin by fitting the ARIMA models as
 393 identified via Table 1 to the data and extract the fitted values which are then
 394 combined with our original SSA Bcd signal to create a hybrid SSA-ARIMA
 395 signal for Bcd. We consider the examples discussed in text so far and generate
 396 the following results. Figure 9 shows the hybrid SSA-ARIMA signals for Bcd
 397 data. In comparison to the optimised SSA signals in Figure 5, the hybrid
 398 SSA signal with ARIMA fit fails to meet the smooth criteria. As such, it is
 399 evident that on its own, the hybrid SSA-ARIMA approach is only beneficial
 400 for those who wish to capture all the signal in the data whilst ensuring that
 401 the residual following Bcd signal extraction is white noise. It clearly comes at
 402 a high cost of lost smoothness in signal curves. However, it is of note that as
 403 previously mentioned, noise in gene expression data enters not only from the
 404 data acquisition and processing procedures [24] but also the fluctuations seen
 405 in an expression pattern can be a consequence of biological noise which may
 406 also introduce error into the data [25]. Therefore, the source of the natural
 407 biological variability is different from the experimental noise [25]. Biological
 408 noise arises from the active molecular transport, compartmentalization, and the
 409 mechanics of cell division [26]. Therefore, the hybrid SSA with the ARIMA
 410 model can be applied in studies such as segmentation network analysis where
 411 the combination of Bcd signal with its biological noise needs to be considered
 412 as an input to the system.

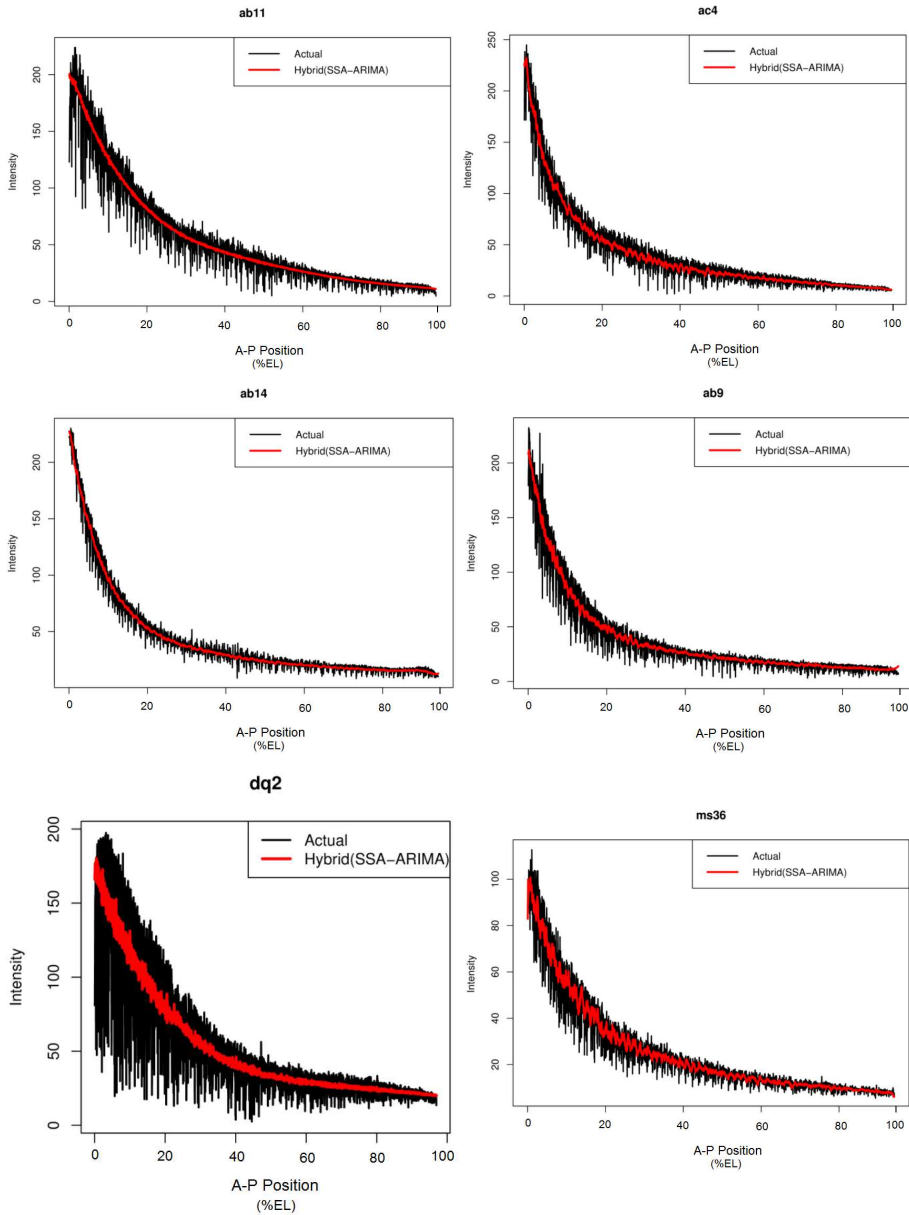


Figure 9: Hybrid SSA signal with ARIMA fit for Bcd data.

4.2.2 Hybrid SSA signal: Nonparametric Approach

Here, we apply the same process as above, but instead of ARIMA, we rely on the nonparametric time series analysis model of ETS. This enables the entire hybrid SSA signal approach to remain nonparametric in nature. The resulting hybrid SSA signals with ETS fit are shown via Figure 10.

There is an interesting point to note here. In comparison to the parametric hybrid signal extraction approach, it is clear that the nonparametric hybrid approach has resulted in much smoother signal curves as one would expect and like to see following a signal extraction exercise. As such, out of the two hybrid approaches, for the purposes of Bcd signal extraction, it is likely that users will prefer the nonparametric approach over the parametric approach.

5 Conclusion

This paper begins with the core aim of introducing new criteria for optimising Bcd signal extraction. Motivated by the findings in [2], we opt to tailor the new Bcd signal extraction criteria for use with the Singular Spectrum Analysis technique which Ghodsi et al. [2] found to be the best option for Bcd signal extraction in relation to SDD, ARIMA, ETS, ARFIMA and NN models. In line with our aim, we initially produce an algorithm for optimising the Bcd signal extraction process with SSA. In brief, the algorithm is optimised based on minimising the skewness statistic for the SSA residual. We suggest that setting L equal to the minimum skewness within the threshold $10 \geq L \geq N/4$ and combining this SSA choice with $r = 1$ or $r = 1, 2$ as appropriate, will enable users to obtain the optimal Bcd signal extraction with SSA.

Through this research, we have succeeded in presenting several contributions to the field of Bcd signal extraction. The first and most important of which deals with the application of the newly proposed algorithm to 27 real Bcd data to show that it can enable researchers to select the appropriate SSA choices to extract a smooth and accurate Bcd signal quickly and easily without the need to spend an increased amount of time for the selection of L for decomposing the data. However, we notice that given the highly complex nature of the Bcd data, on one occasion the SSA algorithm fails to extract an entirely smooth signal. As a solution, we introduce the concept of Sequential SSA on signals, as the second contribution from this research. Via this approach, we are able to refine and smoothen the initial signal which had previously captured some of the observational and biological noise in Bcd data.

In line with good practice, in addition to evaluating the signal extractions alone, this study also pays attention to the residuals. The analysis of the residuals motivated us to introduce hybrid SSA based signal extraction processes for Bcd. In brief, when extracting the trend alone from any given data set, one would reasonably expect other signals to end up within the noise component. However, this would mean that the residual is no longer random and some statisticians could find it difficult to accept such techniques. Accordingly, the first hybrid SSA signal process (and the third contribution from this research) is focussed on providing a Bcd signal extraction procedure which will ensure

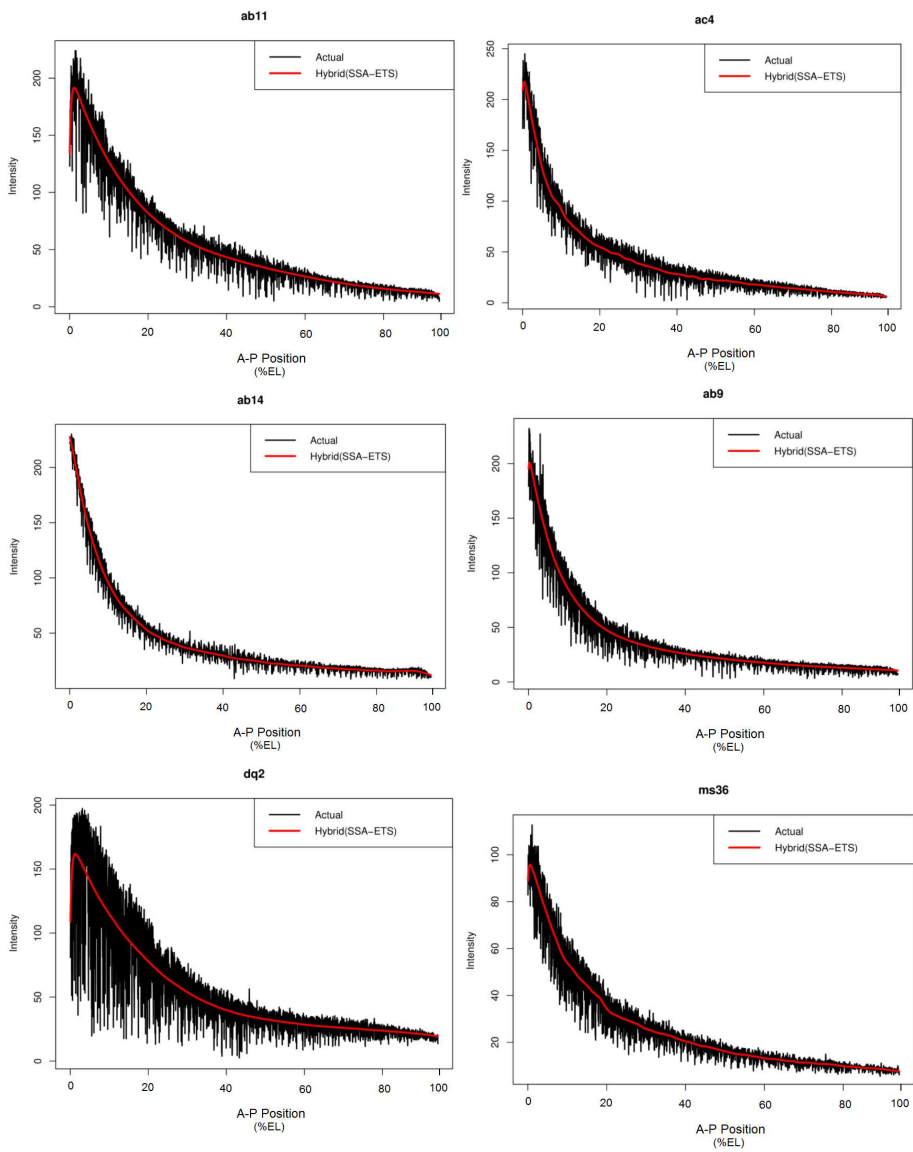


Figure 10: Hybrid SSA signal with ETS fit for bicoid data.

457 the residual is white noise. This was achieved by combining the optimised SSA
458 signal with optimised ARIMA models being fitted to the residuals. Whilst the
459 results did provide the necessary outcomes in terms of residuals with white
460 noise, it comes at a cost - i.e., a loss in the smoothness of the extracted signal.

461 The SSA-ARIMA hybrid approach is a combination of parametric and non-
462 parametric techniques. For those who wish to rely on nonparametric techniques
463 alone so that one is not restricted by the parametric assumptions, we present the
464 SSA-ETS hybrid Bcd signal extraction approach. This process also produces
465 the fourth and second most important contribution of this research, as we find
466 a solution to the problem of accurately modelling the initial curve in Bcd data
467 which was not only experienced in this paper when we employed the optimised
468 SSA signal extraction process, but also experienced in [2]. Accordingly, we are
469 able to present the hybrid SSA-ETS process, which is a combination of the
470 optimised SSA signal extraction algorithm with an optimised ETS algorithm,
471 as the most efficient approach for Bcd signal extraction.

472 We believe that the findings of this research and the information contained
473 within this paper opens up several avenues for future research. For example,
474 future research should evaluate the possibility of optimizing the SSA signal ex-
475 traction process based on different criteria in order to determine whether a more
476 improved signal extraction can be produced. For example, as we are seeking
477 to introduce a novel approach for optimizing Bicoid signal extraction, in this
478 paper we have relied on a binary decomposition. However, future studies could
479 consider the Colonial Theory based approach to decomposition as presented
480 in [27]. In addition, more extensive research into hybrid signal extraction pro-
481 cesses are likely to result in positive, vital and interesting outcomes as clearly
482 shown via this paper. Researchers should evaluate a variety of different signal
483 extraction techniques within the hybrid framework proposed in this paper to
484 ascertain whether outcomes could be further improved.

485 References

- 486 [1] Alexandrov, T. (2009). A method of trend extraction using Singular Spec-
487 trum Analysis. *REVSTAT*, **7**(1), 1–22.
- 488 [2] Ghodsi, Z., Silva, E. S., and Hassani, H. (2015). *Bicoid* Signal Extraction
489 with a Selection of Parametric and Nonparametric Signal Processing Tech-
490 niques. *Genomics Proteomics Bioinformatics*, **13**, 183–191.
- 491 [3] Hassani, H., and Ghodsi, Z. (2014). Pattern Recognition of Gene Expression
492 with Singular Spectrum Analysis. *Medical Sciences*, **2**(3), 127–139.
- 493 [4] Sanei, S., and Hassani, H. (2015). *Singular Spectrum Analysis of Biomedical*
494 *Signals*. CRC Press.
- 495 [5] Holloway D. M., Harrison, L. G., Kosman, D., Vanario Alonso, C. E., and
496 Spirov, A. V. (2006). Analysis of Pattern Precision Shows That *Drosophila*
497 Segmentation Develops Substantial Independence From Gradients of Mater-
498 nal Gene Products. *Developmental Dynamics*, **235**, 2949–2960.

- 499 [6] Golyandina, N. E., Holloway, D. M., Lopes, F. J. P., Spirov, A. V., Spirova,
500 E. N., and Usevich, K. D. (2012). Measuring gene expression noise in early
501 *Drosophila* embryos: nucleus-to-nucleus variability. *International Conference*
502 *on Computational Science*, **9**, 373–382.
- 503 [7] Spirov, A. V., Golyandina, N. E., Holloway, D. M., et al. (2012). Measuring
504 Gene Expression Noise in Early *Drosophila* Embryos: The Highly Dynamic
505 Compartmentalized Micro-environment of the Blastoderm Is One of the Main
506 Sources of Noise. *Evolutionary Computation, Machine Learning and Data*
507 *Mining in Bioinformatics*, **7246**, 177–188.
- 508 [8] Holloway, D. M., Lopes, F. J. P., da Fontoura Costa, L., Travenolo, B. A.
509 N., Golyandina, N., Usevich, K., and Spirov, A. V. (2011). Gene Expression
510 Noise in Spatial Patterning: hunchback Promoter Structure Affects Noise
511 Amplitude and Distribution in *Drosophila* Segmentation. *PLoS Computa-*
512 *tional Biology*, **7**(2), <https://doi.org/10.1371/journal.pcbi.1001069>.
- 513 [9] Surkova, S., Kosman, D., Kozlov, K., Manu., Myasnikova, E., Samsonova,
514 A. A., Spirov, A., Vanario-Alonso, C. E., Samsonova, M., and Reinitz,
515 J. (2008). Characterization of the *Drosophila* segment determination mor-
516 phome. *Developmental Biology*, **313**(2), 844–862.
- 517 [10] Alexandrov, T., Golyandina, N., and Spirov, A. (2010). Singular spec-
518 trum analysis of gene expression profiles of early *drosophila* embryo:
519 Exponential-in-distance patterns. *Research Letters in Signal Processing*,
520 **2008**, <http://dx.doi.org/10.1155/2008/825758>.
- 521 [11] Huang, N. E., Long, S. R., and Shen, Z. (1996). The Mechanism for Fre-
522 quency Downshift in Nonlinear Wave Evolution. *Advances in Applied Me-*
523 *chanics*, **32**, 59–111.
- 524 [12] Hodrick, R., and Prescott E. C. (1997). Postwar U.S. Business Cycles: An
525 Empirical Investigation. *Journal of Money, Credit and Banking*, **29** 1–16.
- 526 [13] Hassani, H., and Thomakos, D. (2010). A review on singular spectrum
527 analysis for economic and financial time series. *Statistics and its Interface*,
528 **3**, 377–397.
- 529 [14] Hassani, H., Webster, A., Silva, E. S., and Heravi, S. (2015). Forecasting
530 U.S. Tourist arrivals using optimal Singular Spectrum Analysis. *Tourism*
531 *Management*, **46**, 322–335.
- 532 [15] Hassani, H., Silva, E. S., Antonakakis, N., Filis, G., and Gupta, R. (2017).
533 Forecasting accuracy evaluation of tourist arrivals. *Annals of Tourism Re-*
534 *search*, **63**, 112–127.
- 535 [16] Golyandina, N., and Shlemov, A. (2013). Variations of Singular Spectrum
536 Analysis for Separability Improvement: Non-Orthogonal Decompositions of
537 Time Series. *arxiv.org*. Available via: [https://arxiv.org/pdf/1308.4022.](https://arxiv.org/pdf/1308.4022.pdf)
538 [pdf](https://arxiv.org/pdf/1308.4022.pdf). [Accessed: 25.10.2016].

- 539 [17] Hyndman, R. J., and Khandakar, Y. (2008). Automatic Time Series Fore-
540 casting: The forecast Package for R. *Journal of Statistical Software*, **27**:1–22.
- 541 [18] Hyndman, R. J., and Athanasopoulos, G. (2013). *Forecasting: principles*
542 *and practice*. OTexts, Australia. Available via: www.OTexts.com/fpp.
- 543 [19] Kozlov, K., Myasnikova, E., Samsonova, M., Reinitz, J., and Kosman,
544 D. (2000). Method for spatial registration of the expression patterns of
545 *Drosophila* segmentation genes using wavelets. *Computational Technologies*,
546 **5**, 112–119.
- 547 [20] Pisarev, A., Poustelnikova, E., Samsonova, M., and Reinitz, J. (2009).
548 FlyEx, the quantitative atlas on segmentation gene expression at cellular
549 resolution. *Nucleic Acids Research*, **37**(suppl-1), D560–D566.
- 550 [21] Poustelnikova, E., Pisarev, A., Blagov, M., Samsonova, M., and Reinitz,
551 J. (2004). A database for management of gene expression data in situ. *Bioin-*
552 *formatics*, **20**(14), 2212–2221.
- 553 [22] Surkova, S., Kosman, D., Kozlov, K., Myasnikova, E., Samsonova, A. A.,
554 Spirov, A., Vanario-Alonso, C. E., Samsonova, M., and Reinitz, J. (2008).
555 Characterization of the *Drosophila* segment determination morphome. *De-*
556 *velopmental Biology*, **313**(2), 844–862.
- 557 [23] Silva, E. S., Ghodsi, M., Hassani, H., and Abbasirad, K. (2016). A quanti-
558 tative exploration of the statistical and mathematical knowledge of university
559 entrants into a UK Management School. *International Journal of Manage-*
560 *ment Education*, **14**(3), 440–453.
- 561 [24] Wu, Y. F., Myasnikova, E., and Reinitz, J. (2007). Master equation simu-
562 lation analysis of immunostained Bicoid morphogen gradient. *BMC Systems*
563 *Biology*, **1**(1), pp. 52.
- 564 [25] Myasnikova, E., Surkova, S., Panok, L., Samsonova, M., and Reinitz, J.
565 (2009). Estimation of errors introduced by confocal imaging into the data on
566 segmentation gene expression in *Drosophila*. *Bioinformatics*, **25**(3), 346–352.
- 567 [26] Spirov, A. V., Golyandina, N. E., Holloway, D. M., Alexandrov, T.,
568 Spirova, E. N., and Lopes, F. J. (2012). Measuring gene expression noise
569 in early *Drosophila* embryos: the highly dynamic compartmentalized micro-
570 environment of the blastoderm is one of the main sources of noise. In: *Euro-*
571 *pean Conference on Evolutionary Computation, Machine Learning and Data*
572 *Mining in Bioinformatics*, (pp. 177-188). Springer, Berlin, Heidelberg.
- 573 [27] Hassani, H., Ghodsi, Z., Silva, E. S., and Heravi, S. (2016). From nature to
574 maths: Improving forecasting performance in subspace-based methods using
575 genetics Colonial Theory. *Digital Signal Processing*, **51**, 101–109.