

1 **Singular Spectrum Analysis with Conditional**
2 **Predictions for Real-Time State Estimation and**
3 **Forecasting**

4 **H. Reed Ogrosky¹, Samuel N. Stechmann^{2,3}, Nan Chen², and Andrew J.**
5 **Majda^{4,5}**

6 ¹Department of Mathematics and Applied Mathematics, Virginia Commonwealth University, Richmond,
7 VA, USA

8 ²Department of Mathematics, University of Wisconsin-Madison, Madison, WI, USA

9 ³Department of Atmospheric and Oceanic Sciences, University of Wisconsin-Madison, Madison, WI, USA

10 ⁴Department of Mathematics and Center for Atmosphere Ocean Science, Courant Institute of

11 Mathematical Sciences, New York University, New York, NY, USA

12 ⁵Center for Prototype Climate Modeling, NYU Abu Dhabi, Saadiyat Island, Abu Dhabi, United Arab
13 Emirates

14 **Key Points:**

- 15 • Singular spectrum analysis (SSA) and extended empirical orthogonal function (EEOF)
16 methods suffer from endpoint issues.
17 • SSA with conditional predictions (SSA-CP) is presented as a simple modification
18 to improve real-time estimates near endpoints.
19 • Forecasts are also possible, including error estimates, and are optimal for Gaus-
20 sian data and shown to be skillful for non-Gaussian data.

Abstract

Singular spectrum analysis (SSA) or extended empirical orthogonal function (EEOF) methods are powerful, commonly-used data-driven techniques to identify modes of variability in time series and space-time datasets. Due to the time-lag embedding, these methods can provide inaccurate reconstructions of leading modes near the endpoints, which can hinder the use of these methods in real time. A modified version of the traditional SSA algorithm, referred to as SSA with conditional predictions (SSA-CP), is presented to address these issues. It is tested on low-dimensional, approximately Gaussian data, high-dimensional non-Gaussian data, and partially-observed data from a multiscale model. In each case SSA-CP provides a more accurate real-time estimate of the leading modes of variability than the traditional reconstruction. SSA-CP also provides predictions of the leading modes and is easy to implement. SSA-CP is optimal in the case of Gaussian data, and the uncertainty in real-time estimates of leading modes is easily quantified.

1 Introduction

Singular spectrum analysis (SSA) or extended empirical orthogonal function (EEOF) methods are powerful, commonly-used tools available for identifying modes of variability in time series and space-time datasets. SSA's usefulness has been demonstrated in a variety of fields over the last 3-4 decades, including, e.g., nonlinear dynamics (e.g., Broomhead and King, 1986), geoscience (e.g., Weare and Nasstrom, 1982; Vautard and Ghil, 1989; Keppenne and Ghil, 1990; Vautard *et al.*, 1992; Mo, 2001; Kikuchi and Wang, 2008; Roundy and Schreck, 2009), and economics (e.g., Lisi and Medio, 1997; Hassani *et al.*, 2014). Its popularity is due both to its ease of implementation and to its ability to eliminate noise and extract trends, oscillations, and other signals in both univariate and multivariate time series.

As with some other methods for mode identification in space-time data (e.g., Fourier filtering), SSA suffers from endpoint issues; i.e., estimates of leading modes can be inaccurate in real-time without future information. Therefore, SSA may provide inaccurate initial conditions for real-time forecasts. Despite these challenges, it is sometimes used either as a filtering step prior to generating real-time forecasts (e.g., Mo, 2001; Golyandina *et al.*, 2001; Hassani *et al.*, 2014), or in tests of forecast models (e.g., Kang and Kim, 2010; Kondrashov *et al.*, 2013; Chen and Majda, 2015), due to its effectiveness at mode identification.

This motivates the question: Is there a modified version of SSA that (i) is as straightforward to implement as SSA, but that (ii) provides the most accurate real-time state estimation possible of leading modes of variability?

This question, along with the related question of how to best modify SSA for use on datasets with gaps in the data, has motivated the proposal and study of numerous modified versions of SSA. These methods include schemes for modifying incomplete columns of the lag-embedded matrix by weighting known values (Schoellhamer, 2001), iterative SSA methods (Kondrashov and Ghil, 2006; Kondrashov *et al.*, 2010), methods based on linear recurrent formulae (Golyandina and Osipov, 2007), combined recurrent forecasting and hindcasting (Rodrigues and Carvalho, 2013), energy-minimizing reconstructions of principal components (Shen *et al.*, 2014; 2015), and a method utilizing a predicted spatial basis (Chen *et al.*, 2018). Some of these methods will be discussed in Section 5.

Here, we propose and study yet another modification of SSA. This method makes use of conditional mean predictions based on the covariance matrix of the lag-embedded data, and we refer to it as SSA with conditional predictions (SSA-CP). Another appropriate name would be real-time SSA (RT-SSA).

The results of tests shown here suggest that this method is effective at addressing these endpoint issues in a variety of settings. The datasets used in these tests include both univariate datasets and multivariate datasets with small (2-3) or somewhat large (64) number of spatial dimensions; partially observed systems and datasets with all dynamical variables observed; Gaussian and non-Gaussian data; and synthetic time series and time series generated by observational data.

Given these results, there are at least four reasons for using this method. First, it is simple and easy to implement, requiring only small additional steps during the normal SSA algorithm. Second, it provides both state estimation and prediction of leading modes of variability. Third, it provides an optimal reconstruction if the data is Gaussian using the statistics of the first two moments. Fourth, it outperforms many other proposed methods of SSA state estimation for both Gaussian and non-Gaussian data.

The rest of the paper is organized as follows: Section 2 describes the traditional SSA method and the proposed modification. Section 3 lists datasets and models used in tests of this method. Results are presented in Section 4. Discussion of the methods and results is given in Section 5, including a brief comparison of the results with those of other modified SSA methods. Conclusions are given in Section 6.

2 SSA algorithms

A brief review of the traditional SSA algorithm is now given, followed by a description of the proposed modification. When used on multivariate time series, SSA is often referred to as Multichannel SSA (MSSA) in the literature; here SSA will be used to refer to either the univariate or multivariate cases. The theory of SSA, which has been developed over the last several decades, is not discussed here; see, e.g., Aubry *et al.* (1991); Ghil *et al.* (2002); Golyandina *et al.* (2001); Hassani (2007) for discussion of this underlying theory.

2.1 Traditional SSA

We briefly describe the traditional SSA algorithm for a dataset with spatial dimension D ; the traditional univariate SSA algorithm can be reproduced by setting $D = 1$ below.

Let \vec{x}_i be a D -dimensional column vector at time i , with $1 \leq i \leq N$. The four steps of SSA are as follows:

Step 1: Create the time-lagged embedding matrix \mathbf{X} of size $(MD) \times (N-M+1)$:

$$\mathbf{X} = \begin{bmatrix} \vec{x}_1 & \vec{x}_2 & \dots & \vec{x}_{N-M+1} \\ \vec{x}_2 & \vec{x}_3 & \dots & \vec{x}_{N-M+2} \\ \vdots & \vdots & & \vdots \\ \vec{x}_{M-1} & \vec{x}_M & \dots & \vec{x}_{N-1} \\ \vec{x}_M & \vec{x}_{M-1} & \dots & \vec{x}_N \end{bmatrix} \quad (1)$$

where M is the length of the embedding window.

Step 2: Find eigenvalues and eigenvectors of the covariance matrix $\mathbf{C} = \mathbf{X}\mathbf{X}^T / (N-M+1)$. Each eigenvector \vec{v} (sometimes referred to as an empirical orthogonal function, or EOF) is an (MD) -dimensional column vector with corresponding eigenvalue λ :

$$\vec{v} = [\vec{v}_1^T, \dots, \vec{v}_M^T]^T, \quad (2)$$

where \vec{v}_s is a D -dimensional column vector used to denote the lag- s portion of the eigenvector.

107 Step 3: Find the principal component (PC) of each mode by projecting the lag-embedded
 108 data onto the appropriate eigenvector:

$$\vec{\phi} = \mathbf{X}^T \vec{v}. \quad (3)$$

109 The entries of each principal component will be denoted $\vec{\phi} = [\phi_1, \dots, \phi_{N-M+1}]^T$.

110 Step 4: Reconstruct the data corresponding to each mode by calculating the recon-
 111 structed component (RC) $\vec{z}(t)$:

$$\vec{z}(t) = \frac{1}{M_t} \sum_{i=L_t}^{U_t} \phi_{t-i+1} \vec{v}_i \quad (4)$$

112 where (M_t, L_t, U_t) are defined by (see, e.g., Ghil et al., 2002)

$$(M_t, L_t, U_t) = \begin{cases} \left(\frac{1}{t}, 1, t\right), & 1 \leq t \leq M-1 \\ \left(\frac{1}{M}, 1, M\right), & M \leq t \leq N-M+1 \\ \left(\frac{1}{N-t+1}, t-N+M, M\right), & N-M+2 \leq t \leq N \end{cases} \quad (5)$$

113 so that each reconstructed component \vec{z} is a (possibly multivariate) time series of length
 114 N , with each $\vec{z}(t)$ a D -dimensional column vector.

115 Each reconstructed component entry at time t^* depends directly on one embed-
 116 ding window of principal component entries, and each principal component entry depends
 117 on one embedding window of data. As a result, each reconstructed component entry at
 118 time t^* is influenced primarily by the values of \vec{x}_{t^*-M+1} through \vec{x}_{t^*+M-1} ; i.e., two embed-
 119 ding windows worth of data, spanning the window immediately prior to t^* and the
 120 window immediately following t^* , contribute directly to the reconstruction at t^* . For $t^* >$
 121 $N-M$, the embedding window's worth of data immediately following t^* is not entirely
 122 known. The reconstruction process makes use of the known data by averaging over the
 123 available products $\phi_{t-i+1} \vec{v}_i$ in (4), but these final $M-1$ entries of each reconstruction
 124 are only estimates of the state of each mode, and can be expected to change as data be-
 125 comes available at times occurring after the end of the time series. (The same endpoint
 126 issues affect the reconstruction for $t^* < M$.)

127 **2.2 SSA with conditional predictions (SSA-CP)**

128 The primary goal of this section is to present a simple method, SSA with condi-
 129 tional predictions (SSA-CP), that improves the estimates of the final $M-1$ entries of
 130 each reconstructed component, including in particular the current state estimate. In ad-
 131 dition, SSA-CP will provide a prediction of reconstructed components for $t > N$. (The
 132 same procedure may be directly applied to the first $M-1$ entries of each reconstruc-
 133 tion, but for simplicity of presentation, we focus solely on the last $M-1$ entries.)

134 The steps of SSA-CP are as follows:

135 Step 1: Perform steps 1 and 2 of traditional SSA.

136 Step 2: Construct an extended lag-embedded matrix $\tilde{\mathbf{X}}$ of size $(MD) \times N$. The
 137 first N columns of $\tilde{\mathbf{X}}$ are identical to the columns of \mathbf{X} . For the final $M-1$ columns,
 138 those entries which are known from the time series are filled in. The unknown entries
 139 below the diagonal consisting of x_N 's are estimated using their conditional mean pre-

140 diction,

$$\tilde{\mathbf{X}} = \begin{bmatrix} \vec{x}_1 & \dots & \vec{x}_{N-M+1} & \vec{x}_{N-M+2} & \dots & \vec{x}_{N-1} & \vec{x}_N \\ \vec{x}_2 & \dots & \vec{x}_{N-M+2} & \vec{x}_{N-M+3} & \dots & \vec{x}_N & \vec{\mu}_{N+1|N} \\ \vdots & & \vdots & \vdots & & \vdots & \vdots \\ \vec{x}_{M-1} & \dots & \vec{x}_{N-1} & \vec{x}_N & \dots & \vec{\mu}_{N+M-3|N-1,N} & \vec{\mu}_{N+M-2|N} \\ \vec{x}_M & \dots & \vec{x}_N & \vec{\mu}_{N+1|N-M+2,\dots,N} & \dots & \vec{\mu}_{N+M-2|N-1,N} & \vec{\mu}_{N+M-1|N} \end{bmatrix}. \quad (6)$$

141 The calculation of each $\vec{\mu}_{i|N-l,\dots,N}$ in (6) is as follows.

142 Let \vec{y} refer to the k -th column of $\tilde{\mathbf{X}}$, with $N+1 \leq k \leq N+M-1$, and let
 143 \vec{y}_1, \vec{y}_2 refer to the known and unknown portions of $\vec{y} = [\vec{y}_1^T, \vec{y}_2^T]^T$, respectively. If \vec{y} is
 144 a Gaussian random variable with mean $\vec{\mu} = 0$ and covariance matrix \mathbf{C} , then \vec{y}_2 has
 145 a conditional distribution that is Gaussian with mean

$$\vec{\mu}_{2|1} = \mathbf{C}_{21} \mathbf{C}_{11}^{-1} \vec{y}_1, \quad (7)$$

146 where \mathbf{C} can be written as

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{bmatrix} \quad (8)$$

147 with \mathbf{C}_{11} describing the covariance of the known values with themselves, etc. (Kaipio and
 148 Somersalo, 2005). The unknown entries \vec{y}_2 are then filled in with the appropriate entries
 149 of $\vec{\mu}_{2|1}$, where $\vec{\mu}_{N+j|k-M+1,\dots,N}$ in (6) denotes a D -dimensional column vector, i.e. the
 150 j -th set of D entries of the vector $\vec{\mu}_{2|1}$, calculated for the k -th column of $\tilde{\mathbf{X}}$ (with $N+$
 151 $1 \leq k \leq N+M-1$). If necessary, a small amount of noise may be added to the co-
 152 variance matrix in order to evaluate \mathbf{C}_{11}^{-1} in (7).

153 Step 3: Modify step 3 of traditional SSA by replacing \mathbf{X} with $\tilde{\mathbf{X}}$; this change re-
 154 sults in extended principal components $\tilde{\phi} = \tilde{\mathbf{X}}^T \vec{v}$; each extended principal component
 155 is a column vector of length N .

156 Step 4: Modify step 4 of traditional SSA by replacing ϕ with $\tilde{\phi}$ to construct an ex-
 157 tended RC:

$$\tilde{z}(t) = \frac{1}{\tilde{M}_t} \sum_{i=\tilde{L}_t}^{\tilde{U}_t} \tilde{\phi}_{t-i+1} \vec{v}_i \quad (9)$$

158 where $(\tilde{M}_t, \tilde{L}_t, \tilde{U}_t)$ are defined by

$$(\tilde{M}_t, \tilde{L}_t, \tilde{U}_t) = \begin{cases} \left(\frac{1}{t}, 1, t\right), & 1 \leq t \leq M-1 \\ \left(\frac{1}{M}, 1, M\right), & M \leq t \leq N \\ \left(\frac{1}{N-t+M}, t-N+1, M\right), & N+1 \leq t \leq N+M-1 \end{cases} \quad (10)$$

159 so that each extended reconstructed component \tilde{z} is a (possibly multivariate) time se-
 160 ries of length $N+M-1$, with the last $M-1$ entries corresponding to predictions of
 161 the future state of the mode.

162 In the case that the dataset has a Gaussian distribution, the conditional mean pro-
 163 vides an optimal estimate of the missing data (Kaipio and Somersalo, 2005).

164 3 Data and Methods

165 The SSA-CP method will be tested on several datasets and compared to the tra-
 166 ditional SSA reconstruction.

3.1 Data

The first test uses a fifteen year portion of the daily Real-time Multivariate MJO (RMM) indices (Wheeler and Hendon, 2004) from 1 January 1999 through 31 December 2013. The RMM indices have a distribution that is approximately normal with mean and variance approximately 0 and 1, respectively (Chen and Majda, 2015). For this 2-dimensional dataset, $D = 2$ and $N_{tot} = 5479$, with N_{tot} referring to the number of days.

GPCP daily precipitation data (Huffman *et al.*, 2012) are used for the second test. This dataset has a spatial resolution of $1^\circ \times 1^\circ$; the portion from 1 January 1997 through 31 December 2013 is used. Prior to applying SSA, the following steps were taken: (i) a meridional mode truncation to move from $2D(x, y)$ to $1D(x)$, (ii) removal of annual mean and seasonal cycle, and (iii) interpolation to 64 equally-spaced zonal gridpoints. The meridional mode truncation step is a projection of the data onto the leading meridional mode proportional to $e^{-y^2/2}$ where y is proportional to latitude; this step is identical to that used in, e.g., Stechmann and Majda (2015); Stechmann and Ogrosky (2014). Steps (i) and (iii) reduce the number of dimensions to $D = 64$, and the number of times is $N_{tot} = 6209$. Note that these anomalies have a non-Gaussian distribution at each longitude; see the SI for the statistics of these anomalies.

A simulation of a multiscale model (Majda and Harlim, 2012) is used for the third test. The model equations are

$$du_1 = (-\gamma_1 u_1 + F(t)) dt + \sigma_1 dW_1, \quad (11a)$$

$$du_2 = (-\gamma_2 + i\omega_0/\epsilon + ia_0 u_1) u_2 dt + \sigma_2 dW_2, \quad (11b)$$

where $\gamma_1 = \gamma_2 = 0.2$, $\sigma_1 = \sigma_2 = 0.5$, $\omega_0 = a_0 = 1$, $\epsilon = 0.5$, and $F(t) = \sin(t/5)$. An approximate solution was calculated numerically with the Euler-Maruyama method using $dt = 0.005$ and $t_{end} = 2000$. The real part of u_2 was then sampled every 0.5 time units to create a dataset with $D = 1$ and $N_{tot} = 4000$. A portion of this signal can be seen in Figure S2 in the SI.

3.2 Methods

The results of each real-time reconstruction method (SSA-CP and traditional) will be compared with the traditional reconstruction that has knowledge of future data. This is done in two steps.

First, both the traditional SSA and SSA-CP methods were applied to each dataset after removing the final $2M-2$ time entries from the dataset; e.g., using an embedding window of $M = 51$ days for the RMM indices, the methods were applied to the first $N = N_{tot} - 2M + 2 = 5379$ days. The embedding window was chosen to be large enough to be consistent with the intraseasonal timescale of the indices and is similar to that used in Chen and Majda (2015); other choices of this parameter value will be discussed in Section 5. The standard reconstruction $z(t)$ for each mode therefore has $N = 5379$ entries, while the SSA-CP reconstruction $\tilde{z}(t)$ has $N + M - 1 = 5429$ entries. Note that the first $N - M + 1 = 5329$ entries for each reconstruction method are identical to one another; i.e. $z(t) = \tilde{z}(t)$ for $1 \leq t \leq N - M + 1$. Next, the traditional reconstruction method was used again, this time on the full $N_{tot} = 5479$ entries, resulting in a reconstruction $u(t)$ with $N_{tot} = 5479$ entries. The entries of $u(t)$ up to $N_{tot} - M + 1 = 5429$ are taken to be ‘truth’, and each of the methods applied to the shorter time series are compared with this truth.

Second, these tests are repeated for each dataset with decreasing N_{tot} ; i.e., define $N_{tot,i} = N_{tot} - i + 1$, and repeat the test described above but using only the first $N_{tot} = N_{tot,i}$ entries of the dataset, so that $N = N_i := N_{tot,i} - 2M + 2$. For the RMM indices and multiscale model, $i \in I = [1, \dots, 1001]$; for the GPCP data, $i \in I = [1, 6, 11, \dots, 1001]$.

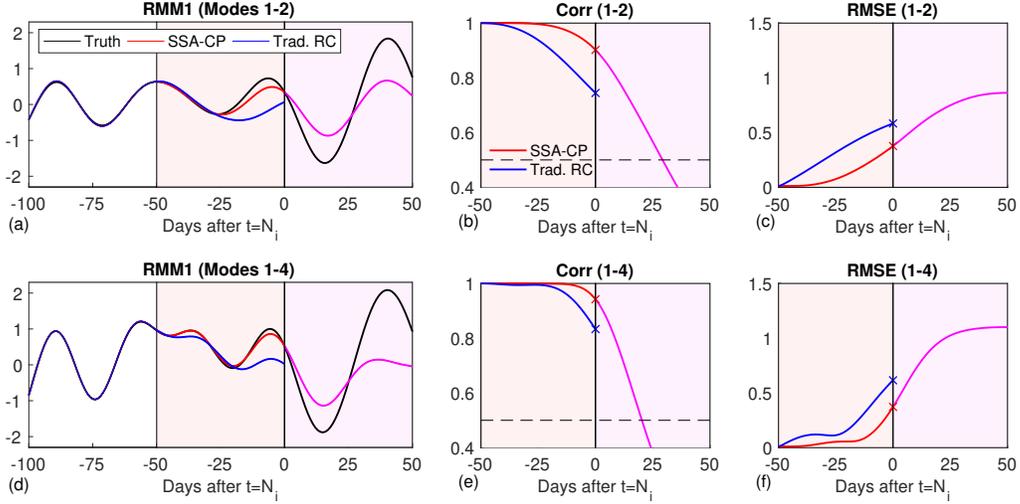


Figure 1. (a) Reconstructed RMM1 using components 1-2 with $t = N_{601} = 4779$ (31-Jan-2012) using (blue) traditional reconstruction, (red/magenta) SSA-CP, and (black) reconstruction using future information. (b,c) Bivariate pattern correlation and RMSE of the (blue) traditional reconstruction and truth as a function of days prior to/after N_i , and (red/magenta) SSA-CP reconstruction and truth using modes 1-2. (d-f) Same as (a-c) but using components 1-4.

213 The pattern correlation and root mean square error (RMSE) are then calculated as a
 214 function of days before or after N_i ; see the SI for details.

215 4 Results

216 We next show results for three tests.

217 4.1 RMM index

218 How well does the method perform on low-dimensional data that is nearly Gaus-
 219 sian?

220 Figure 1(a,d) shows the results of using the SSA-CP or traditional reconstruction
 221 methods on the RMM indices with an embedding window $M = 51$ days. For times away
 222 from the endpoints of the data i.e. $t < N_i - M + 1$, both methods are in agreement
 223 with the truth. For past times near the endpoints, i.e. $N_i - M + 1 < t < N_i$ (light or-
 224 ange shaded region), SSA-CP captures both the phase and amplitude of the RMM1 in-
 225 dex better than the traditional reconstruction. For future times $t > N_i$, SSA-CP is able
 226 to make predictions, with good agreement in phase and an underestimate of the ampli-
 227 tude of the true reconstruction. This underestimate of amplitude is due to using con-
 228 ditional mean predictions which tend to zero as $t \rightarrow \infty$.

229 Figure 1(b,c,e,f) shows that when these tests are repeated, SSA-CP has significantly
 230 improved pattern correlation and reduced error compared to the traditional reconstruc-
 231 tion. As a current state estimation, at $t = N_i$ SSA-CP improves the pattern correla-
 232 tion from 0.74 to 0.90 (0.83 to 0.94) for the 2 (4) leading modes. Likewise, SSA-CP re-
 233 duces the error at $t = N_i$ from 0.58 to 0.38 (0.62 to 0.37). For future times $t > N_i$,
 234 SSA-CP is able to make meaningful predictions for an extended period of time, with pat-
 235 tern correlations exceeding 0.5 out to approximately 29 (20) days when 2 (4) leading modes
 236 are used.

237

4.2 Precipitation data

238

How well does the method perform on large-dimensional, possibly non-Gaussian data?

239

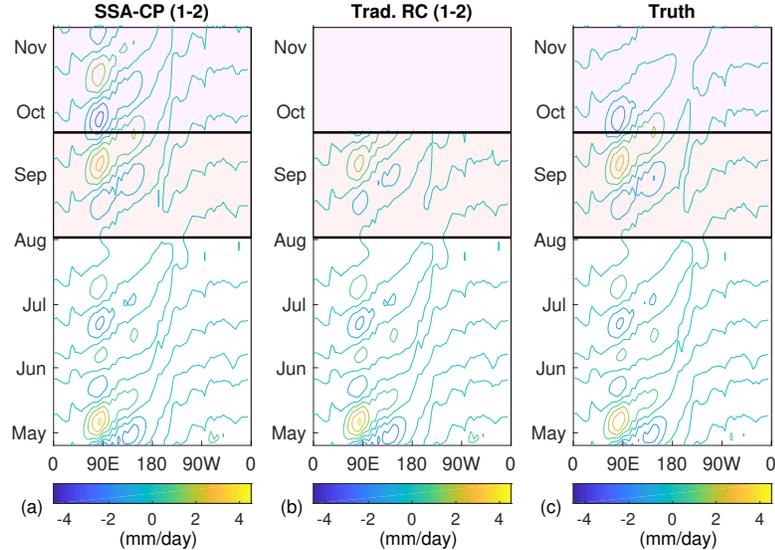


Figure 2. (a) Reconstructed precipitation during 2013 using SSA-CP modes 1-2 with $t_N = 6109$, corresponding to 22 September 2013. (b) Same as (a) but using traditional reconstruction. (c) Reconstructed modes 1-2 using future information.

240

241

242

243

244

245

246

Figure 2 shows reconstructed precipitation anomalies using the 2 leading modes with an embedding window of 51 days. Both methods produce identical reconstructions prior to 2 August 2013. For 3 August 2013 through 22 September 2013, SSA-CP produces a reconstruction with amplitude in much better agreement with the non-real-time reconstruction (truth) than the traditional reconstruction. It also provides a prediction with decaying amplitude throughout October, qualitatively similar to the truth but with slower decay.

247

248

249

250

251

252

253

254

Repeating these tests for various N_i produces the pattern correlation and RMSE shown in Figure 3. For the recent past in time interval $N_i - M + 1 < t < N_i$, SSA-CP produces higher pattern correlation and lower RMSE than the standard reconstruction method. For state estimation at $t = N_i$, the pattern correlation is 0.1-0.2 higher at almost all longitudes when using SSA-CP than when using the standard method. Likewise, the RMSE is lower using SSA-CP than the traditional reconstruction at all longitudes. Note that low pattern correlation values for each method at longitudes like 150W are due to small anomalies in the leading modes.

255

4.3 Partially-observed multiscale model

256

How well does the method perform on partially-observed data?

257

258

259

260

261

Figure S3 in the SI shows the pattern correlation and RMSE for both methods applied to the multiscale model (11). For $N_i - M + 1 < t < N_i$, SSA-CP has significantly higher pattern correlation and lower error than the traditional reconstruction. At $t = N_i$, using SSA-CP improves the pattern correlation from 0.54 to 0.75 for 2 leading modes, and lowers the error from 0.12 to 0.06. For $t > N_i$, predictions using SSA-CP have a

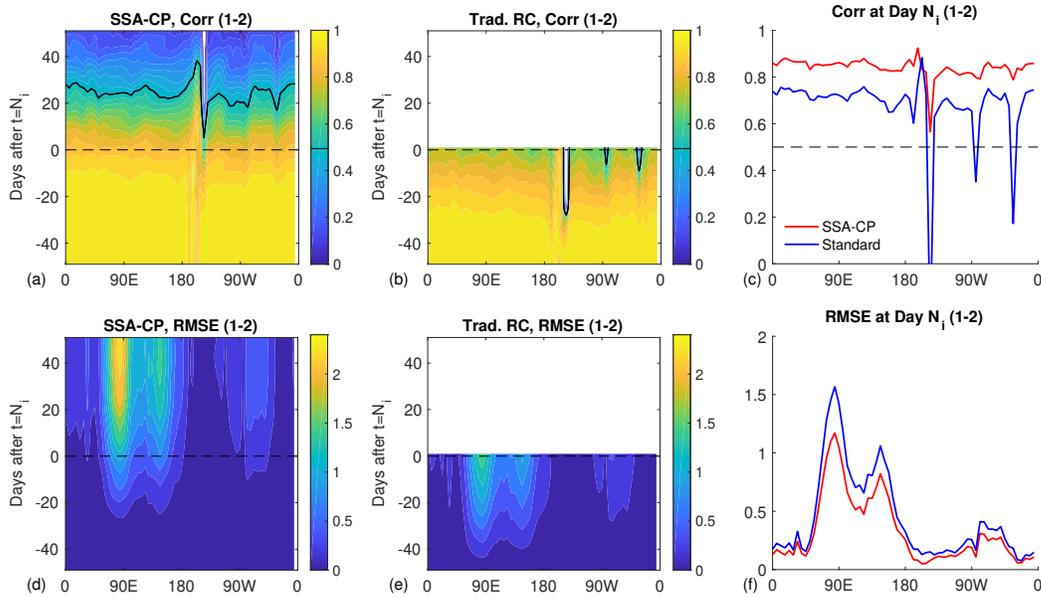


Figure 3. (a) Pattern correlation, using 200 runs of SSA, of reconstructed precipitation components (1-2) using SSA-CP. (b) Same as (a) but for traditional reconstruction. (c) Pattern correlation at Day $t = N_i$ for each method. (d-f) Same as (a-c) but showing RMSE.

262 pattern correlations of 0.5 or higher out to approximately 23 days when 2 leading modes
 263 are used.

264 5 Discussion

265 SSA-CP has been proposed as a method that supplements the mode-identification
 266 ability of SSA with improved estimates of mode reconstructions near the ends of time
 267 series. We note that it is not at all necessarily the best possible data-driven, model-free
 268 prediction method that could be designed. Its effectiveness at identifying modes of vari-
 269 ability in real-time is of course also limited to cases where SSA is effective at identify-
 270 ing modes of interest.

271 How sensitive are the results to changes in the embedding window? As a first step
 272 towards addressing this question, the RMM tests from the previous section were rerun
 273 with an embedding window of $M = 75$ days. Figure 4(a-d) shows that while both SSA-
 274 CP and the traditional reconstruction produce slightly lower pattern correlation at $t =$
 275 N_i than in the previous test with $M = 51$, SSA-CP again results in significantly higher
 276 PC and lower RMSE than the traditional reconstruction. For $t > N_i$, the pattern cor-
 277 relation stays higher than 0.5 for 35 (24) days when the leading 2 (4) modes are used
 278 (not shown). Extensive testing of this sensitivity is left for future work.

279 How does SSA-CP compare with other methods in the literature that have been
 280 proposed for either (i) improving state estimation of reconstructed components near the
 281 endpoints of time series, or (ii) using SSA on datasets with gaps? We briefly examine
 282 this through a comparison of the results of SSA-CP with methods from Schoellhamer
 283 (2001) and Golyandina and Osipov (2007) for the first test from Section 4. Figure 4(a-
 284 d) shows the pattern correlation and RMSE of these two methods along with the tradi-
 285 tional reconstruction and SSA-CP. All of the modified versions of SSA produce higher
 286 pattern correlation than the traditional reconstruction, with SSA-CP having the high-
 287 est. For the leading two modes, all methods produce lower RMSE than the traditional

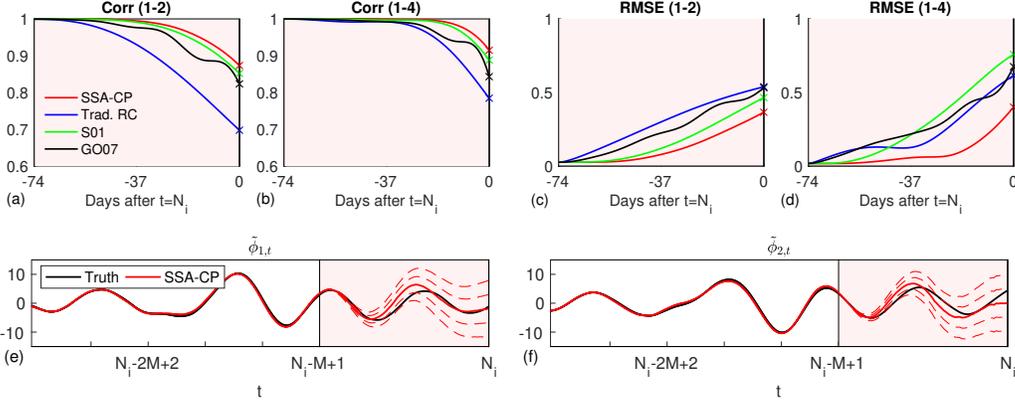


Figure 4. (a-b) Bivariate pattern correlation and RMSE of reconstructed RMM indices using the (blue) traditional reconstruction, (red) SSA-CP reconstruction, (green) weighted reconstruction of Schoellhamer (2001), and (black) Π -projector/simultaneous filling in method of Golyandina and Osipov (2007) as a function of days prior to/after N_i , using modes 1-2; here $M = 75$, $N_i = 4779$. (c-d) Same as (a-b) but for modes 1-4. (e-f) Leading two principal components of SSA-CP; dashed red lines indicate $\pm 1, 2$ standard deviations.

288 reconstruction, but when the leading four modes are used, only SSA-CP outperforms the
 289 traditional reconstruction over each of the final $M - 1$ days.

290 Does SSA-CP outperform the traditional reconstruction on other datasets? In ad-
 291 dition to the tests described here, other tests were conducted using datasets generated
 292 by stochastic processes (complex-valued Ornstein-Uhlenbeck process), deterministic dy-
 293 namical systems (Lorenz 63 model, multiple examples from Golyandina *et al.* (2001)),
 294 other observational data (Kelvin wave calculated using NCEP/NCAR reanalysis data (Kalnay
 295 *et al.*, 1996) and the methods of Ogrosky and Stechmann, 2015; 2016), and numerous
 296 synthetic test signals both with and without noise. SSA-CP significantly outperformed
 297 traditional SSA in almost all of these tests. In cases of deterministic signals of Golyan-
 298 dina *et al.* (2001), both methods produced excellent reconstructions of the leading modes
 299 near the endpoints. In cases like this, the standard reconstruction may be just as desir-
 300 able as SSA-CP or any other modification, as the additional effort of implementing SSA-
 301 CP, though minimal, may not be necessary to provide reasonable initial conditions for
 302 a forecast. In addition, one benefit of the standard reconstruction is its invertibility; if
 303 all modes are reconstructed and summed together, the original dataset is recovered. This
 304 invertibility is not shared by SSA-CP.

305 There are several compelling reasons for using SSA-CP rather than the traditional
 306 reconstruction, however. First, it is nearly as simple to use as traditional SSA. Second,
 307 it is optimal for Gaussian data and is based on well-known theory. Third, it is straight-
 308 forward to quantify the uncertainty in the extended principal components or reconstruction.
 309 For example, the variance of $\tilde{\phi}_{N-l+1}$, where $1 \leq l \leq M - 1$, is given by

$$\text{Var}(\tilde{\phi}_{N-l+1}) = [\tilde{v}_{l+1}^T, \dots, \tilde{v}_M^T] \mathbf{C}_{21} [\tilde{v}_{l+1}^T, \dots, \tilde{v}_M^T]^T \quad (12)$$

310 Figure 4(e,f) shows the two leading principal components of the RMM indices calculated
 311 using SSA-CP with $N_i = 4779$ and $M = 75$. One and two standard deviations from
 312 the extended principal component entries are shown, with the standard deviation calcu-
 313 lated using (12).

314 Finally, since non-Gaussianity leads to a lack of independence between modes in
 315 linear methods like empirical orthogonal functions (EOFs), there is no guarantee that

316 the method will work well on data with strong non-Gaussianity (Monehan *et al.*, 2009).
 317 However, the method works well on the non-Gaussian data used here, perhaps owing to
 318 the somewhat mild deviations from Gaussianity. The method could potentially be ex-
 319 tended to non-Gaussian frameworks with conditional Gaussian or Gaussian mixture struc-
 320 tures (see, e.g. Chen and Majda (2018); Majda (2016)).

321 We note that SSA is just one of many data analysis tools capable of identifying modes
 322 of variability in spatiotemporal datasets (see Crommelin and Majda, 2004, for a discus-
 323 sion of some other linear methods for mode identification). SSA was chosen to be the
 324 focus of the current study due to its linearity, simplicity, and popularity, combined with
 325 the linearity of the proposed modifications. Other mode identification methods, includ-
 326 ing nonlinear methods like Nonlinear Laplacian Spectral Analysis (NLSA), have been
 327 shown to be effective at capturing modes of variability that SSA has difficulty captur-
 328 ing, like modes with pronounced intermittent behavior (Giannakis and Majda, 2012a,b),
 329 and theory supporting both such methods and forecasting techniques of relevance has
 330 been developed in recent years (Comeau *et al.*, 2017; Zhao and Giannakis, 2016). Includ-
 331 ing conditional predictions into such methods is certainly possible, though it is not clear
 332 how much one can expect this linear approach to improve the results of a nonlinear method
 333 like NLSA.

334 6 Conclusions

335 In summary, a modified SSA algorithm, SSA-CP, has been presented and tested.
 336 This modification is proposed to address endpoint issues that arise when using SSA. When
 337 compared with the traditional reconstruction method, SSA-CP results in significantly
 338 improved real-time estimates of leading modes of variability when applied to a variety
 339 of datasets.

340 This method was shown to be useful for providing improved initial conditions for
 341 forecasts. It is derived from well-known theory using Gaussian statistics, and provides
 342 optimal predictions for Gaussian data, but also performs well in tests with non-Gaussian
 343 data. The uncertainty in the real-time estimates may be quantified using the covariance
 344 matrix that is inherently part of the method.

345 While the current study has been primarily focused on applying the method to at-
 346 mospheric science data, this method may prove useful in application areas outside of at-
 347 mospheric science. In addition, it is possible that the ideas used here may be adapted
 348 for other methods of mode identification. These subjects are left for future work.

349 Acknowledgments

350 The GPCP data for this article are available from NOAA/OAR/ESRL PSD, Boulder,
 351 Colorado, USA, from their web site at <http://www.esrl.noaa.gov/psd/>. The RMM in-
 352 dices can be obtained online at <http://www.bom.gov.au/climate/mjo/>. Other data used
 353 are in the figures.

354 The research of S.N.S. is partially supported by a Sloan Research Fellowship from
 355 the Alfred P. Sloan Foundation and a Vilas Associates Award from the University of Wisconsin-
 356 Madison. The research of N.C. is supported by the Office of Vice Chancellor for Research
 357 and Graduate Education (VCRGE) at University of Wisconsin-Madison.

358 References

359 Aubry, N., Guyonnet, R., & Lima, R. (1991). Spatiotemporal analysis of complex
 360 signals: theory and applications. *Journal of Statistical Physics*, *64*, 683–739.

- 361 Broomhead, D. S., & King, G. P. (1986). Extracting qualitative dynamics from
362 experimental data. *Physica D*, *20*, 217–236.
- 363 Chen, N., & Majda, A. J. (2015). Predicting the real-time multivariate Madden-
364 Julian oscillation index through a low-order nonlinear stochastic model.
365 *Monthly Weather Review*, *143*, 2148–2169.
- 366 Chen, N., Majda, A. J., Sabeerali, C. T., & Ravindran, A. J. (2017). Predicting
367 Monsoon Intraseasonal Precipitation using a Low-Order Stochastic Model.
368 *Journal of Climate*, *31*, 4403–4427.
- 369 Chen, N., & Majda, A. J. (2018). Conditional Gaussian Systems for Multiscale
370 Nonlinear Stochastic Systems: Prediction, State Estimation and Uncertainty
371 Quantification. *Entropy*, *20*, 509.
- 372 Comeau, D., Giannakis, D., Zhao, Z., & Majda, A. J. (2018). Predicting regional
373 and pan-Arctic sea ice anomalies with kernel analog forecasting. *Climate Dy-*
374 *namics*. <https://doi.org/10.1007/s00382-018-4459-x>.
- 375 Crommelin, D. T., & Majda, A. J. (2004). Strategies for model reduction: Compar-
376 ing different optimal bases. *Journal of the Atmospheric Sciences*, *61*, 2206–
377 2217.
- 378 Ghil, M., Allen, R. M., Dettinger, M. D., Ide, K., Kondrashov, D., Mann, M. E.,
379 Robertson, A., Saunders, A., Tian, Y., Varadi, F., & Yiou, P. (2002). Ad-
380 vanced spectral methods for climatic time series, *Reviews of Geophysics*, *40*,
381 pp. 3.1–3.41.
- 382 Giannakis, D. & Majda, A. J. (2012). Nonlinear Laplacian spectral analysis for time
383 series with intermittency and low-frequency variability. *Proceedings of the Na-*
384 *tional Academy of Sciences of the United States of America*, *109*, 2222–2227.
- 385 Giannakis, D., & Majda, A. J. (2012). Comparing low-frequency and intermittent
386 variability in comprehensive climate models through nonlinear Laplacian spec-
387 tral analysis. *Geophysical Research Letters*, *39*, L10710.
- 388 Golyandina, N., Nekrutkin, V., & Zhigljavsky, A. A. (2001). *Analysis of Time Series*
389 *Structure: SSA and Related Techniques*, Chapman and Hall, CRC.
- 390 Golyandina, N., & Osipov, E. (2007). The “Caterpillar”-SSA method for analysis of
391 time series with missing values. *Journal of Statistical Planning and Inference*,
392 *137*, 2642–2653.
- 393 Hassani, H. (2007). Singular Spectrum Analysis: Methodology and Comparison.
394 *Journal of Data Science*, *5*, 239–257.
- 395 Hassani, H., Webster, A., Silva, E.S., & Heravi, S. (2014). Forecasting U.S. tourist
396 arrivals using optimal singular spectrum analysis. *Tourism Management*, *46*,
397 322–335.
- 398 Huffman, G.J., Bolvin, D.T., & Adler, R.F. (2012). GPCP Version 2.2 SG Combined
399 Precipitation Data Set. WDC-A, NCDC, Asheville, NC. Data set accessed 12
400 February 2014 at <http://www.ncdc.noaa.gov/oa/wmo/wdcamet-ncdc.html>.
- 401 Kaipio, J., & Somersalo, E. (2005). *Statistical and Computational Inverse Problems*,
402 Springer, New York.
- 403 Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell,
404 M., Saha, S., White, G., Woollen, J., Zhu, Y., Chelliah, M., Ebisuzaki, W.,
405 Higgins, W., Janowiak, J., Mo, K.C., Ropelewski, C., Wang, J., Leetmaa, A.,
406 Reynolds, R., Jenne, R. & Joseph, D. (1996). The NCEP/NCAR 40-year re-
407 analysis project. *Bulletin of the American Meteorological Society*, *77*, 437–471.
- 408 Kang, I.-S., & Kim, H.-M. (2010). Assessment of MJO predictability for boreal
409 winter with various statistical and dynamical models. *Journal of Climate* *23*,
410 2368–2378.
- 411 Keppenne, C. L., & Ghil, M. (1990). Adaptive filtering and prediction of the South-
412 ern Oscillation Index. *Journal of Geophysical Research: Atmospheres*, *97*,
413 20,449–20,454.
- 414 Kikuchi, K., & Wang, B. (2008). Diurnal Precipitation Regimes in the Global Trop-
415 ics. *Journal of Climate*, *21*, 2680–2696.

- 416 Kondrashov, D., & Ghil, M. (2006). Spatio-temporal filling of missing points in
 417 geophysical data sets. *Nonlinear Processes in Geophysics*, *13*, 151–159.
- 418 Kondrashov, D., Shprits, Y., & Ghil, M. (2010). Gap filling of solar wind data by
 419 singular spectrum analysis. *Geophysical Research Letters*, *37*, L15101.
- 420 Kondrashov, D., Chekroun, M. D., Robertson, A. W., & Ghil, M. (2013). Low-order
 421 stochastic model and “past-noise forecasting” of the Madden-Julian Oscilla-
 422 tion. *Geophysical Research Letters*, *40*, 5305–5310.
- 423 Lisi, F., & Medio, A. (1997). Is a random walk the best exchange rate predictor?
 424 *International Journal of Forecasting*, *13*, 255–267.
- 425 Majda, A. J. (2016). *Introduction to Turbulent Dynamical Systems in Complex Sys-*
 426 *tems*, Springer.
- 427 Majda, A. J., & Harlim, J. (2012). *Filtering Complex Turbulent Systems*, Cambridge
 428 University Press.
- 429 Mo, K. C. (2001). Adaptive filtering and prediction of intraseasonal oscillations.
 430 *Monthly Weather Review*, *129*, 802–817.
- 431 Monehan, A. H., Frye, J. C., Ambaum, M. H. P., Stephenson, D. B., & North, G. R.
 432 (2009). Empirical Orthogonal Functions: The Medium is the Message. *Journal*
 433 *of Climate*, *22*, 6501–6514.
- 434 Ogrosky, H. R., & Stechmann, S. N. (2015). Assessing the equatorial long-wave
 435 approximation: asymptotics and observational data analysis. *Journal of the*
 436 *Atmospheric Sciences*, *72*, 4821–4843.
- 437 Ogrosky, H. R., & Stechmann, S. N. (2016). Identifying convectively coupled equato-
 438 rial waves using theoretical wave eigenvectors. *Monthly Weather Review*, *144*,
 439 2235–2264.
- 440 Rodrigues, P. C., & de Carvalho, M. (2013). Spectral modeling of time series with
 441 missing data. *Applied Mathematical Modelling*, *37*, 4676–4684.
- 442 Roundy, P. E., & Schreck III, C. J. (2009). A combined wave-number–frequency and
 443 time-extended EOF approach for tracking the progress of modes of large-scale
 444 organized tropical convection. *Quarterly Journal of the Royal Meteorological*
 445 *Society*, *135*, 161–173.
- 446 Schoellhamer, D. H. (2001). Singular spectrum analysis for time series with missing
 447 data. *Geophysical Research Letters*, *28*, 3187–3190.
- 448 Shen, Y., Li, W., Xu, G., & Li, B. (2014). Spatiotemporal filtering of regional GNSS
 449 network’s position time series with missing data using principal component
 450 analysis. *Journal of Geodesy*, *88*, 1–12.
- 451 Shen, Y., Peng, F., & Li, B. (2015). Improved singular spectrum analysis for time
 452 series with missing data. *Nonlinear Processes in Geophysics*, *22*, 371–376.
- 453 Stechmann, S. N., & Majda, A. J. (2015). Identifying the skeleton of the Madden-
 454 Julian oscillation in observational data. *Monthly Weather Review*, *143*, 395–
 455 416.
- 456 Stechmann, S. N., & Ogrosky, H. R. (2014). The Walker circulation, diabatic heat-
 457 ing, and outgoing longwave radiation. *Geophysical Research Letters*, *41*, 9097–
 458 9105.
- 459 Vautard, R., & Ghil, M. (1989). Singular spectrum analysis in non-linear dynamics,
 460 with applications to paleoclimatic time series. *Physica D*, *35*, 395–424.
- 461 Vautard, R., Yiou, P., & Ghil, M. (1992). Singular spectrum analysis: A toolkit for
 462 short noisy chaotic signals. *Physica D*, *58*, 95–126.
- 463 Weare, B. C., & Nasstrom, J. S. (1982). Examples of Extended Empirical Orthogo-
 464 nal Function Analyses. *Monthly Weather Review*, *110*, 481–485.
- 465 Wheeler, M. C., & Hendon, H. H. (2004). An all-season real-time multivariate
 466 MJO index: development of an index for monitoring and prediction. *Monthly*
 467 *Weather Review*, *132*, 1917–1932.
- 468 Zhao, Z., & Giannakis, D. (2016). Analog forecasting with dynamics-adapted ker-
 469 nels. *Nonlinearity*, *29*, 2888–2939.