# Sparse eigenbasis approximation: multiple feature extraction across spatiotemporal scales with application to coherent set identification

Gary Froyland, Christopher P. Rock, and Konstantinos Sakellariou

School of Mathematics and Statistics

University of New South Wales

Sydney NSW 2052, Australia

December 18, 2018

## Abstract

The output of spectral clustering is a collection of eigenvalues and eigenvectors that encode important connectivity information about a graph or a manifold. This connectivity information is often not cleanly represented in the eigenvectors and must be disentangled by some secondary procedure. We propose the use of an approximate sparse basis for the space spanned by the leading eigenvectors as a natural, robust, and efficient means of performing this separation. The use of sparsity yields a natural cutoff in this disentanglement procedure and is particularly useful in practical situations when there is no clear eigengap. In order to select a suitable collection of vectors we develop a new Weyl-inspired eigengap heuristic and heuristics based on the sparse basis vectors. We develop an automated eigenvector separation procedure and illustrate its efficacy on examples from time-dependent dynamics on manifolds. In this context, transfer operator approaches are extensively used to find dynamically disconnected regions of phase space, known as almost-invariant sets or coherent sets. The dominant eigenvectors of transfer operators or related operators, such as the dynamic Laplacian, encode dynamic connectivity information. Our sparse eigenbasis approximation (SEBA) methodology streamlines the final stage of transfer operator methods, namely the extraction of almost-invariant or coherent sets from the eigenvectors. It is particularly useful when used on domains with large numbers of coherent sets, and when the coherent sets do not exhaust the phase space, such as in large geophysical datasets.

# Contents

# 1   Introduction

Spectral clustering has found broad application in areas such as network analysis, manifold learning (e.g. diffusion maps), Lagrangian dynamics, and stochastic processes. The output of a spectral method is the spectrum and eigenvectors of some matrix or linear operator, typically either of Laplace type such as a graph Laplacian or Laplace-Beltrami operator on a manifold, or of Markov type[1]. The eigenvectors carry geometric structure to e.g. help discern the topology of an unknown manifold, cluster point data, or analyse nonlinear dynamics.

For example, the Laplace-Beltrami operator arises as the generator of heat flow on a manifold. The four leading eigenfunctions of this operator[2] (corresponding to the largest four (real) eigenvalues $0 = \lambda_1 \geq \cdots \geq \lambda_4$) are shown in the upper row of Figure 1. These



Figure 1: Upper: the leading four eigenfunctions of the Laplace-Beltrami operator on the manifold shown. Lower: a sparse basis output by Algorithm 3.1, separating the four "blobs" at the periphery of the central disk. Dark blue represents zero value. The eigenfunctions were computed using a finite-element mesh with 10774 nodes.

eigenfunctions – beyond the first "trivial" constant eigenfunction corresponding to $\lambda_1 = 0$ and shown in Figure 1 (upper left) – describe signed "heat modes" that decay most slowly under heat flow, at rates $\sim e^{\lambda_i t}$, $i = 2, \ldots, 4$. Because the four "blobs" at the periphery of the main disk have narrow channels connecting them to the main disk, the slowest decaying mode (second from the left in the upper row of Figure 1) has a lot of heat (red) in the bottom blob and lower part of the main disk and is cool (blue) in the upper blob and upper part of the main disk. The four "blobs" are correctly detected in the upper row of Figure 1, but they are mixed together in the leading four eigenfunctions. We wish to separate these four features and this is what is achieved in the lower row of Figure 1, by finding an approximate sparse basis through orthogonal rotation of the four leading eigenfunctions.

Arguably the most popular method of performing this separation is the embedding of (discrete versions of) the second to fourth eigenfunctions in $\mathbb{R}^3$ and applying a hard clus-

---

[1] often there is a "spectral mapping" relationship between these two types, of the form $M = \exp(L)$, where $M$ is Markov type and $L$ is Laplace type.

[2] with homogeneous Neumann boundary conditions.

Figure 2: Left: clustering an embedding of the second to fourth eigenvectors from Figure 1 in $\mathbb{R}^3$ (10774 points in $\mathbb{R}^3$) and applying k-means clustering, requesting four clusters. Right: As per left, but requesting five clusters.

tering algorithm such as k-means; see e.g. [34] for a detailed description and the historical development of k-means. In this example, for the k-means approach to at least partially isolate the peripheral blobs in the example of Figure 1, one has to request a clustering of *five* clusters because k-means *partitions* the dataset/domain, and requesting four clusters groups each of the peripheral blobs with approximately one quarter of the large central disk; see Figure 2. A key advantage of our approach is that we do *not necessarily* separate by partitioning, but instead may classify a large part of the dataset/domain as "unclustered". Such an idea has been partly implemented in [31] by adding one more cluster to k-means as a "background" cluster, but we show in Sections 5.2 and 5.3 that this is not a good approach for more challenging dynamics. Other separation approaches are discussed in Section 2.1.

Perhaps the most similar precursor to our approach in the dynamical systems literature, applied in the context of Markov chains, is the search for a linear transformation of the leading eigenvectors [11] to make them as close as possible to a collection of indicator functions (vectors taking a small number of distinct values) spanning a similar subspace, under constraints of non-negativity and small intersection of support. The optimisation in [11] can be expensive for large vectors and large numbers of vectors, and will always return a full partition of the domain. An LU decomposition approach [4] has been effective in the case where the eigenvectors are already close to linear combinations of indicator vectors. This is a somewhat idealised situation and, in such a setting, all of the above methods should perform very well.

In the context of general Markov processes and Lagrangian dynamics, there are many situations where it is not instructive to *partition* all nodes or the entire manifold into almost-disconnected pieces according to the relevant eigenvectors. Rather, it is more informative to know that some parts of the network or manifold have a natural almost-disconnected structure, while the remainder (possibly a large remainder) of the network or manifold is well connected. The sparsity built in to our new sparse basis approach will automatically identify those nodes or subsets of a manifold that do not belong to any almost-disconnected subset; such nodes or subdomains will be zeroed out by the sparsification.

In Section 5 we illustrate our approach on spectral algorithms designed to identify almost-

invariant or coherent sets in dynamical systems on manifolds. The term *almost-invariant set* refers to a set that is approximately invariant under autonomous dynamics. The term *coherent set* refers to a time-dependent family of sets that are approximately equivariant under time-dependent or nonautonomous dynamics. In contrast to an almost-invariant set, which is fixed in phase space, the family of coherent sets are in general mobile in phase space as time evolves; see [27] for a short overview. Almost-invariant and coherent sets in the phase space of complicated dynamics are important because they are the most predictable features in the medium term. These concepts have been used in molecular dynamics to identify conformations to aid drug design [10], in atmospheric dynamics to identify atmospheric vortices [28], in oceanography to identify ocean gyres [26] and eddies [18, 19], and in fluid dynamics [48]. Spectral approaches to the identification of coherent features rely on the dominant eigenvectors of transfer operators or related operators, such as the dynamic Laplacian [16, 23]. There are various categories of constructions, based on the problem to be solved; see Table 1. The final stage of the extraction of almost-invariant or coherent sets from eigenvectors or singular vectors has previously proceeded either by time-consuming hand-selected thresholds or by common clustering methods such as k-means.

A main goal of this research is to describe simple and robust procedures for automatically separating features represented in the eigenvectors or singular vectors arising from transfer operator methods. Our separation procedures streamlines the use of these methods on domains with large numbers of coherent sets, particularly when the coherent sets do not exhaust the phase space. Although the dynamical context is our primary motivation, we believe our approach will translate well to other application domains where disentangling the output of spectral clustering is important. We have thoroughly tested our approach on a variety of nonlinear dynamics and have proposed heuristics that have consistently performed well in our tests. We illustrate our new approach by separating coherent features in idealised fluid flow dynamics and in large-scale ocean currents in the North Atlantic derived from satellite altimetry.

Our main contributions include a refinement of the standard eigengap criterion for determining suitable numbers of features. Our Weyl-inspired refinement strongly highlights and appropriately scales gaps in the spectrum, emphasising natural time scales. We also develop new vector-based heuristics for determining suitable numbers of features; these help to determine good choices of dimension for the subspaces and therefore the corresponding spatial scales of features. We develop a "reliability" heuristic for the sparse basis vectors, including feature rejection criteria. These vector-based heuristics supplement the eigengap criteria and we have found them to be particularly useful in realistic situations where there is no clear eigengap. Finally, we propose automatic methods of hard thresholding of the sparse basis to provide hard separation of features.

An outline of the paper is as follows. Section 2 provides a brief background on hard separation of spectral clustering output, sparse principal components analysis, and transfer operator and dynamic Laplace operator constructions in nonlinear dynamics. Section 3 describes in detail the variant of sparse basis approximation we will use, based on the SPCArt algorithm [33]. Section 4 discusses the use of our sparse basis algorithm, including

new eigenvalue and vector based heuristics for selecting the size of the basis and assigning reliability to the sparse basis output, and simple thresholding procedures. We present our numerical examples in Section 5 and conclude in Section 6.

# 2 Background

## 2.1 Separation of spectral clustering output by hard clustering

Denote the eigenvalue (resp. eigenvector) output of a spectral clustering method by ordered lists of eigenvalues $\lambda_1 \geq \lambda_2 \geq \cdots$ and their corresponding eigenvectors $v_1, v_2, \ldots$. Suppose that we wish to separate the features encoded in the vectors $v_1, \ldots, v_r$. One of the most popular separation approaches is to perform hard clustering on the $p \times r$ data array $V = [v_1|\cdots|v_r]$, treating each row of this array as a Euclidean point in $\mathbb{R}^r$. Figure 2 illustrates this procedure with $r = 3$ using k-means for the hard clustering. In the context of nonlinear dynamics, the k-means algorithm [34] has been applied to eigenvector data embedded in Euclidean space, for example [31, 3, 46, 43, 21] to identify coherent sets. A "soft" variant, fuzzy c-means, has earlier been used in the dynamics context on eigenvector data [17, 14] to identify almost-invariant sets. Other variants include k-medians [35], k-medoids [39], and kernel k-means [47].

Other hard separation approaches for spectral clustering include [53], who try to find a rotation of the eigenvectors that brings them as close as possible to a collection of binary vectors, which represent a feature partition. The work of [41] builds on [53] by proposing rotating the eigenvectors to approximately achieve a collection of nonnonegative vectors, followed by a maximum likelihood assignment. Related to this is [32] who search for non-negative vectors "nearby" the eigenvectors via a penalty term. Another recent contribution includes [42], who search for sparse orthogonal vectors that maximise the trace of the usual bilinear form involving the Laplace matrix, followed by an application of k-means. All of the above approaches produce a *partition* of the dataset/domain (as does k-means), and this is something we wish to avoid through the use of sparse bases. Moreover, the explicit imposition of nonnegative vectors removes what we have found to be an important indicator of the quality of the separation, as we show in Section 4. There are several other related methods that involve manipulating the weights in the affinity matrix, but we omit discussing these because in our dynamical systems applications the Markov and Laplace operators that arise have additional special properties we wish to preserve.

## 2.2 Sparse principal component analysis

Principal component analysis (PCA) is an unsupervised learning technique often used for dimension reduction and feature selection. PCA may be viewed as the process of finding a linear projection from a high-dimensional space (occupied by the dataset) onto a low-dimensional space that preserves as much of the variance of the original dataset as possible. One drawback of PCA is that every variable in the lower-dimensional space generally depends

on *all* of the variables in the high-dimensional space.

This shortcoming motivated sparse principal component analysis (SPCA), which modifies PCA by imposing constraints (resp. adding a penalty) to force (resp. encourage) each variable in the lower-dimensional space to depend on only a subset of the original variables. In 2003, Joliffe et al. proposed the first sparse principal component analysis technique, called SCoTLASS [36], and several versions of SPCA have appeared since then. Many such versions, termed *deflation methods*, find one sparse principal component at a time, then deflate the original data in that direction and find the next orthogonal principal component, e.g. [36, 5]. To avoid sequential deflation finding a non-global optimum, so-called *block methods* have been proposed, which optimise all the subspaces at once [37, 6, 54, 33].

Because of its simplicity, efficiency, and robustness, in this paper we will use a tailored version of the sparse principal component analysis by rotation and truncation (SPCArt) algorithm [33], which is a block method. Letting $\mathcal{V} = \text{span}\{v_1, \ldots, v_r\}$ we implicitly assume that the eigenvectors from our spectral clustering satisfy the following assumption.

**Assumption A.** *There is a $\mathcal{V}$-approximating, $r$-dimensional subspace $\mathcal{S} \subset \mathbb{R}^p$ with a non-negative sparse basis $s_1, \ldots, s_r$ whose supports have little or no intersection.*

Assumption A should hold when there exist pairwise disjoint subsets of the graph or manifold that are sufficiently strongly disconnected, provided suitable eigenvectors $\{v_1, \ldots, v_r\}$ of a Laplace-type or Markov-type operator or matrix are chosen. Figure 1 (lower row) and Figure 3 (right column) illustrate the type of sparse basis vectors we will obtain.

## 2.3 Markov- and Laplace-type operators arising in nonlinear dynamics

Our examples in Section 5 will be drawn from nonlinear dynamics and we provide a brief overview here, referring the reader to original papers for further details. One has a smooth domain $M$, typically a smooth $d$-dimensional manifold or subset of $\mathbb{R}^d$, on which the nonlinear dynamics acts.

There are several related linear operators that are intimately connected with the nonlinear evolution on phase space. These operators fall into two broad classes, which we will call *Markov-type* and *Laplace-type*. The former have their spectrum contained in $\{z \in \mathbb{C} : |z| \leq 1\}$, the closed unit disk in the complex plane, and the spectrum of the latter is contained in $\{z \in \mathbb{C} : \Re(z) \leq 0\}$, the closed half plane with non-positive real part.

In the most straightforward case, the dynamics is generated by repeated application of a transformation $T : M \to M$ or by solving a differential equation $\dot{x} = F(x)$. Both of these mechanisms generate *autonomous* dynamics because the underlying governing dynamics does not depend on time. In the former case one uses left[3] eigenvectors of the transfer operator (Markov type) to identify subsets of the domain that are approximately invariant under repeated application of the dynamics; see Construction 1 in Table 1. In the latter case,

---

[3]the description of left and/or right eigenvectors has been chosen to match the original papers. For Markov-type operators, multiplying vectors on the left of matrices is consistent with most Markov chain texts.

| Construction | Autonomous/ Nonauton. Dynamics | Finite/ Infinite Time | Closed/ Open Dynamics | Uses time derivative | Objects Identified |
|---|---|---|---|---|---|
| 1. Left eigenvectors of transfer operator [8] | Auton. | $\infty$ | closed | no | almost-invariant sets |
| 2. Left eigenvectors of generator [22, 23] | Auton. or Periodic | $\infty$ | closed | yes | almost-invariant sets |
| 3. Right eigenvectors of transfer operator (resp. generator) [40, 29] | Auton. | $\infty$ | open | no (resp. yes) | Basins of attraction |
| 4. Right eigenvectors of symmetrised transf. op. [14, 27] | Auton. | finite | closed | no | finite-time almost-inv. sets |
| 5. Singular vectors of transfer operator [28, 15] | Nonauton. | finite | closed | no | finite-time coherent sets |
| 6. Eigenvectors of dynamic Laplacian [16, 24] | Nonauton. | finite | closed | no | finite-time coherent sets |
| 7. Oseledets vectors of transfer operator cocycle [25] | Nonauton. | $\infty$ | closed | no | coherent sets |

Table 1: Summary of transfer operator and dynamic Laplace operator constructions to identify almost-invariant and coherent sets, and the dynamical settings in which each construction is used.

one uses left eigenvectors of the generator (Laplace type) to detect almost-invariant sets in continuous time without trajectory integration; this is Construction 2 in Table 1. When there is open dynamics created by the existence of a "hole(s)" in phase space, one is interested in the basins of attraction; this is Construction 3 in Table 1, covering both discrete and continuous time. In many instances, there is a natural timescale of interest and one wishes to find subsets that are approximately invariant over this *finite* timescale; this is addressed in Construction 4 in Table 1.

*Nonautonomous* dynamics arises when the generating law varies with time. In discrete time one applies a sequence of different transformations $T_n \circ T_{n-1} \circ \cdots \circ T_1$, and in continuous time, one solves a differential equation with time-dependent right-hand-side $\dot{x} = F(x, t)$. Because the underlying dynamics changes over time, one wishes to find subsets that remain "coherent" over the finite time duration; this is addressed in Constructions 5 and 6 in Table 1. Finally, identifying coherent sets in time-dependent dynamics over an infinite time period is addressed in Construction 7 in Table 1. While we have described the dynamics above as deterministic, each of the Constructions 1–7 can also be applied to stochastic dynamics arising from transformations with additional noise, or stochastic differential equations.

Each of the Constructions 1–7 in Table 1 is a type of spectral clustering; Constructions 1, 3, 4, 5, and 7 produce Markov-type operators and Constructions 2, 3, and 6 produce Laplace-type operators. The eigenvectors arising from Constructions 4, 5, and 6 are orthogonal in a suitable inner product space, while the vectors produced by Constructions 1, 2, 3, and 7

need not be. For the latter constructions, one applies e.g. QR factorisation to orthogonalise the vectors before applying Algorithm 3.1.

There are various ways in which the constructions in Table 1 can be numerically implemented. A thin spline implementation of Construction 5 is introduced in [52], and a spectral implementation in continuous time of Construction 5 is presented in [9]. A radial basis function implementation of Construction 6 is contained in [20] and a finite element implementation is introduced in [21]; the latter is the numerical implementation we will use for all dynamic Laplacian eigenvector computations in this paper. A graph-based method similar to Construction 6 is described in [31]. Diffusion map implementations of Constructions 5 and 6 are described in [3].

In section 5 we will apply our method to output from Constructions 5 and 6; the former is of Markov type and the latter is of Laplace type. While the different constructions in Table 1 solve different dynamical problems, the properties of the eigenvectors from all constructions are strongly related because they arise from spectral approaches. We expect our methods to perform equally well on all constructions in Table 1.

Figure 3: Left: the 8 leading eigenvectors $v_1, \ldots, v_8$ (upper to lower) of the dynamic Laplacian constructed from the Bickley jet (see Section 5.1). Right: an approximate sparse basis $s_1, \ldots, s_8$ output by Algorithm 3.1 separating the 8 main coherent sets; small negative values of have been removed from the vectors $s_7, s_8$.

10

# 3 Sparse eigenbasis approximation

Let $v_1, \ldots, v_r \in \mathbb{R}^p$ be a set of linearly independent eigenvectors or singular vectors; typically $r \ll p$. These vectors form a basis for a subspace $\mathcal{V} \subset \mathbb{R}^p$. We wish to transform these vectors to a basis of sparse vectors $s_1, \ldots, s_r \in \mathbb{R}^p$ for a subspace $\mathcal{S} \subset \mathbb{R}^p$ with $\mathcal{V} \approx \mathcal{S}$. Without loss of generality, we can assume that the vectors $v_1, \ldots, v_r$ are an orthonormal basis; if not, we perform a QR factorisation of the matrix $V = [v_1 | \cdots | v_r]$ and extract the orthonormal columns. If the subspace $\mathcal{V}$ satisfies Assumption A, we can expect the sparse vectors $s_1, \ldots, s_r$ to have supports with small overlaps, and therefore be close to orthogonal. Thus, we wish to find an orthogonal "rotation" matrix $R$ that transforms $v_1, \ldots, v_r$ into $s_1, \ldots, s_r$ such that (i) the $r$-dimensional space $\mathcal{S} \subset \mathbb{R}^p$ spanned by $s_1, \ldots, s_r$ is close to $\mathcal{V} \subset \mathbb{R}^p$ (approximate subspace preservation) and (ii) the vectors $s_1, \ldots, s_r$ are sparse.

The approach we take is based on the SPCArt algorithm (sparse PCA via truncation and rotation [33]). We present below a simplified "full rank" version of SPCArt, tuned to the setting of Assumption A, which we call *Sparse Eigenbasis Approximation* (SEBA). Let $\mathfrak{S}^r$ denote the Stiefel manifold $\{A \in \mathbb{R}^{r \times r} : A^\top A = I_r\}$ of $r \times r$ orthogonal matrices, and let $\mathfrak{U}^{p,r} = \{A \in \mathbb{R}^{p \times r} : \text{each column of } A \text{ has } \ell_2 \text{ norm } 1\}$. Define $\|A\|_F := \sqrt{\sum_{i=1}^p \sum_{j=1}^r A_{ij}^2}$ to be the Frobenius norm and $\|A\|_{1,1} := \sum_{i=1}^p \sum_{j=1}^r |A_{ij}|$ to be the $\ell_{1,1}$ matrix norm. Given some small positive sparsity parameter $\mu$, we wish to solve the following nonconvex optimisation problem for $S \in \mathfrak{U}^{p,r}$ and $R \in \mathfrak{S}^r$:

$$\underset{\substack{S \in \mathfrak{U}^{p,r} \\ R \in \mathfrak{S}^r}}{\arg\min} \frac{1}{2}\|V - SR\|_F^2 + \mu\|S\|_{1,1}, \tag{1}$$

The $r$ columns of $S$ define a sparse basis $\{s_1, \ldots, s_r\}$ for a subspace $\mathcal{S}$ close to $\mathcal{V}$. The first term in (1) measures how close[4] the rotated columns of $V$, namely $VR^T$, are to the columns of $S$ (recall $\|VR^\top - S\|_F = \|V - SR\|_F$ for orthogonal $R$). The second term measures the sparsity of the basis formed by the columns of $S$ using the common sparsity-inducing [12] $\ell_1$ penalty.

Because of the nonconvexity of the problem (1), finding a global optimum $(S, R) \in \mathfrak{U}^{p,r} \times \mathfrak{S}^r$ is difficult. Hu *et al.* [33] proposed alternately fixing $R$ and optimising $S$, and fixing $S$ and optimising $R$ to find a local minimum $(S, R)$. Each of these individual optimisation problems are fast to solve exactly.

1. **Fixed $R$:** The optimisation for $S$ in (1) is exactly solved by "soft thresholding" [37]. Define a thresholding transformation $C_\mu : \mathbb{R} \to \mathbb{R}$ by $C_\mu(z) = \text{sign}(z) \max\{|z| - \mu, 0\}$; $C_\mu$ is applied to vectors elementwise. For $j = 1, \ldots, r$, set the $j^{th}$ column of $S$ to $S_j = C_\mu((VR^\top)_j)/\|C_\mu((VR^\top)_j)\|$.

2. **Fixed $S$:** The optimisation for $R$ becomes a Procrustes problem $\min_{R \in \mathfrak{S}^r} \frac{1}{2}\|V - SR\|_F^2$. This problem may be efficiently exactly solved by polar decomposition [54]. Set $R =$

---

[4]Note that if we define $S := VR^\top$ for some $r \times r$ orthogonal (in fact, $R$ nonsingular is sufficient) matrix $R$, then the columns of $S$ will span the same subspace as the columns of $V$.

Polar$(S^\top V)$, where Polar$(\cdot)$ is the orthonormal component of the polar decomposition; that is, $S^\top V = RH$, where $R$ is orthogonal and $H$ is symmetric and positive definite. If $S^\top V$ has singular value decomposition $S^\top V = PDQ^\top$, then $R = PQ^\top$. In MATLAB, we use `[P,~,Q]=svd(S'*V,0);`

One initialises with $R = I_r$ and alternately applies steps 1 and 2 above until the change in $R$ is below a specified tolerance[5]. The parameter $\mu$ should be chosen less than $1/\sqrt{p}$ as $C_\mu$ sends the constant unit vector to the zero vector for $\mu \geq 1/\sqrt{p}$. Hu *et al.* [33] suggest using $\mu = 1/\sqrt{p}$ and with extensive experimentation we have found that this works well, as do values between 75% and 100% of $1/\sqrt{p}$. In our experiments we have used $\mu = 0.99/\sqrt{p}$ because in many cases the constant unit vector is part of our initial basis $\mathcal{V}$. Alternative thresholding methods were presented in [33], but in our experiments we found soft thresholding was the most robust and we use this in all computations. In summary, our sparse eigenbasis approximation (SEBA) algorithm is:

**Algorithm 3.1** ((SEBA) – Input orthonormal $p \times r$ matrix $V$; output sparse $p \times r$ matrix $S$).

1. Set $\mu = 0.99/\sqrt{p}$ and $R = I_r$.

2. For $j = 1, \ldots, r$, set the $j^{th}$ column of $S$ to $S_j = C_\mu((VR^\top)_j)/\|C_\mu((VR^\top)_j)\|$.

3. Set $R = \text{Polar}(S^\top V)$.

4. If the matrix 2-norm of the difference between the revised $R$ in step 3 and the previous $R$ is larger than some tolerance[5], go to step 2; otherwise go to step 5.

5. For $j = 1, \ldots, r$, set $S_j \to \text{sign}(\sum_{i=1}^{p} S_{ij})S_j$.

6. For $j = 1, \ldots, r$ set $S_j \to S_j/\max_{1 \leq i \leq p} S_{ij}$.

7. Reorder the columns of $S$ so that $m_j := \min_i S_{ij}$ is in decreasing order, $j = 1, \ldots, r$.

In Step 5 we possibly change the sign of the columns of $S$ to ensure they are predominantly nonnegative and in Step 6 we scale the columns of S so that their maximum value is 1; we will later interpret $S_{ij}$ as a likelihood of membership of index $i$ in the $j^{th}$ feature. Step 7 orders the columns of $S$ in terms of "reliability"; see Section 4.1 for a discussion.

As noted in [33], SPCArt (steps 1–4 of Algorithm 3.1) produces vectors $s_1, \ldots, s_r$ that should span a similar subspace to $\mathcal{V}$ because of the first term of (1). Because the input vectors $v_1, \ldots, v_r$ are orthonormal and the vectors $s_1, \ldots, s_r$ are arrived at by orthogonal rotation and truncation, the latter should also be close to orthogonal. The sparsity should be reasonably balanced because of the global optimisation of $R$ in step 3 and uniform thresholding across each vector in step 2. A short MATLAB code listing to implement Algorithm 3.1 is included in Appendix A.1.

---

[5]We used a tolerance of $10^{-14}$ in our experiments, and we observed there is little, if any, difference if the tolerance is increased by a few orders of magnitude.

# 4 Using the SEBA Algorithm 3.1

In this section we describe the "reliability" ordering in Step 7 of Algorithm 3.1, and introduce some rules of thumb for selecting an appropriate basis size for input to Algorithm 3.1 and how to extract a (sub)partition from the output of Algorithm 3.1. We also detail how to (optionally) apply weights to vector entries. Code is provided in Appendices A.2 and A.3.

## 4.1 Ordering the sparse vectors in terms of "reliability"

Under Assumption A, if Algorithm 3.1 is performing well we should be able to find a subspace $\mathcal{S}$ that is a good approximation of $\mathcal{V}$ and which is spanned by a nonnegative sparse basis $s_1, \ldots, s_r$. Nonnegativity of $s_1, \ldots, s_r$ is not guaranteed by Algorithm 3.1 and we have found that if a sparse vector $s_j$ has one or more large negative entries, this is an indication that the feature $s_j$ has not been cleanly separated from other features. We propose a "reliability ordering" of the sparse output vectors: in Step 7 of Algorithm 3.1 we order the vectors $s_1, \ldots, s_r$ so that $m_j := \min_{1 \le i \le p} s_{ij}$ is in descending order. That is, $j$ for which $m_j = 0$ means that $s_j$ should highlight a very reliably coherent feature and appear early in the ordering, while those $s_j$ toward the end of the ordering are potentially increasingly spurious. In the case of the Bickley jet with $r = 8$, we found $m_j = 0$, $j = 1, \ldots, 6$, with the remaining $m_7, m_8 > -0.02$; see the eight sparse vectors shown in Figure 3 (right column), ordered by decreasing $m_j$. In Section 4.2.2 we will adapt this heuristic to select the number of input eigenvectors.

## 4.2 Selecting an appropriate number of input vectors

Consider a set $v_1, \ldots, v_r \in \mathbb{R}^p$ of linearly independent eigenvectors as in Section 3. As one increases $r$, the $r^{\text{th}}$ eigenvector becomes increasingly oscillatory, corresponding to increasingly rapidly decaying modes. The "spatial scale" of the features encoded in the $r^{\text{th}}$ eigenvector therefore decreases with increasing $r$. Similarly, the $r^{\text{th}}$ eigenvalue describes the exponential decay rate of the $r^{\text{th}}$ eigenvector, and therefore a "temporal scale" for the corresponding features. In the dynamical systems context, the number $r$ controls the smallest *spatial scale* at which eigenvectors identify almost-invariant or coherent sets, while the $r^{\text{th}}$ eigenvalue provides a lower bound on the *temporal scale* at which each of the almost-invariant or coherent sets mix with the rest of phase space. In particular, gaps in the spectrum indicate a jump in temporal scales, rather than a jump in spatial scales. The reliability ordering of Section 4.1 will guide the choice of a number $k \le r$ to indicate that Algorithm 3.1 has produced $k$ reliable features at the spatiotemporal scale governed by $r$.

### 4.2.1 Heuristics based on the spectrum

In spectral clustering, the eigengap heuristic (see [50] and references therein) is a common method for selecting an appropriate number of vectors for further analysis. The eigengap heuristic suggests to look for large gaps in the spectrum, i.e. if $|\lambda_{r+1} - \lambda_r|$ is large compared

to $|\lambda_{i+1} - \lambda_i|$ for $i = 1, \ldots, r - 1$, one should truncate the collection of eigenvectors to $v_1, \ldots, v_r$. A potential issue with this rule (apart from the absence of an unambiguous spectral gap) is that depending on the dimension of the underlying manifold, the natural asymptotic behaviour of $\lambda_k$ as a function of $k$ given by Weyl's law [51] is not necessarily linear. We therefore propose a new variant of the eigengap heuristic by taking this asymptotic into account. If our phase space is a $d$-dimensional Riemannian manifold, denote by $N(\lambda)$ the number of eigenvalues of the Laplace(-Beltrami) operator $\Delta$ on $M$ no greater than $\lambda$ (applying either homogeneous Neumann or Dirichlet boundary conditions). Weyl's law states that $\lim_{\lambda \to \infty} N(\lambda)/\lambda^{d/2} = -(2\pi)^{-d} \omega_d |M|$ where $\omega_d$ is the volume of the $d$-dimensional unit ball and $|M|$ is the volume of $M$. Thus, ordering the eigenvalues of $\Delta$ as $0 \geq \lambda_1 \geq \lambda_2 \geq \cdots$, we have $\lambda_r \sim -Cr^{2/d}$ for a constant $C$ as $r \to \infty$.

**Laplace-type matrix/operator, Neumann boundary conditions:** For example, a graph Laplacian, a Laplace-Beltrami operator, or the dynamic Laplacian[6] [16, 23] of Sections 5.1 and 5.3 with homogeneous Neumann boundary conditions. In this setting, the leading eigenvalue is $\lambda_1 = 0$ and fitting the asymptotic $\lambda_r \sim -Cr^{2/d}$ to this condition suggests plotting $\lambda_r/(r-1)^{2/d}$ vs. $r$ for $r = 2, \ldots$, and looking for the largest drops from $r$ to $r + 1$; see Figure 4.



Figure 4: Left: plot of $\lambda_r$ vs $r$ for $r = 1, \ldots, 20$ for the Bickley jet of Section 5.1. Right: plot of $\lambda_r/(r-1)$ vs $r$ for $r = 2, \ldots, 20$. The eigengaps in the left image at $r = 2$ and $r = 8$ are highlighted even more strongly in the Weyl rescaling in the right image. The eigengap at $r = 15$ in the left image has been de-emphasised in the right image relative to $r = 2$ and $r = 8$.

**Laplace-type operator, Dirichlet boundary conditions:** For example, a Laplace-Beltrami operator or the dynamic Laplacian with Dirichlet boundary conditions. The leading eigenvalue $\lambda_1$ is strictly negative (and unknown) so we plot $\lambda_r/r^{2/d}$ vs. $r$ for $r = 1, \ldots$, and look for the largest drops from $r$ to $r + 1$.

---

[6]which is not necessarily a Laplace-Beltrami operator [38].

**Markov-type operator:** For example, a matrix or operator arising from diffusion maps [3], the various transfer operator constructions in Table 1, or the normalised finite-time transfer operator [28, 15] of Section 5.2. In these examples $\lambda_1 = 1$ and the rest of the spectrum is contained in the unit circle in the complex plane. Taking logs and retaining only the real part, we obtain a Laplace-type spectrum with Neumann boundary conditions. Thus, we plot $\Re(\log \lambda_r)/(r-1)^{2/d}$ vs. $r$ for $r = 2, \ldots$, and look for the largest drops from $r$ to $r + 1$. In the case of transfer operators arising from continuous-time flows, this logarithmic transformation can be made rigorous via the spectral mapping theorem, see [22, 23].

The above heuristics are applied to spectra arising from both transfer operators and the dynamic Laplacian in Section 5.

### 4.2.2 A heuristic based on sparse vectors

Following the arguments of Section 4.1, for the $p \times r$ matrix $S$ produced by Algorithm 3.1, we create a cumulative minimum value quantity $\text{Min}(S) := \sum_{j=1}^{r} -m_j = -\sum_{j=1}^{r} \min_{1 \leq i \leq p} S_{ij}$. Denoting by $S^{(r)}$ the sparse array produced by Algorithm 3.1 with $r$ input vectors, one may then plot $\text{Min}(S^{(r)})$ vs. $r$ and select those $r$ for which there has been a large drop from $r - 1$; i.e. $\text{Min}(S^{(r)}) - \text{Min}(S^{(r-1)})$ is negative. The rationale for this choice is that $\text{Min}(S^{(r)})$ is typically increasing with $r$ because it is a sum of non-negative values. If there is a drop from $r - 1$ to $r$ it means that despite adding an extra term to the sum, the sum decreases, and thus overall the new $m_j$ values (based on $r$ input vectors) are smaller than the old $m_j$ (based on $r - 1$ input vectors), indicating that the new sparse basis $s_1, \ldots, s_r$ has better separated the $r$ features. Note that such a heuristic implicitly assumes that *all* $r$ sparse vectors are reliable, and no sparse vector will be rejected from the collection. In Section 5.2 we will combine this heuristic with the reliability heuristic of Section 4.1 to suggest a way to simultaneously select both $r$ (the number of input vectors), and $k$ (the number of reliable sparse vectors), with $k < r$ after some sparse vectors have been rejected.

## 4.3 Extraction of a (sub)-partition from a sparse basis output

Recall that we wish to separate the features encoded in the eigenvectors $v_1, \ldots, v_r$ resulting from a spectral clustering. Algorithm 3.1 has already done most of the hard work by forming a sparse approximate basis of the space spanned by the eigenvectors. Indeed one would often be very satisfied with this separation procedure. A simple and effective way of summarising the separated features is to create a "superposition vector".

**Superposition vector:** $\mathfrak{s} := \min\{1, \sum_{j=1}^{r} \max\{S_{ij}, 0\}\}$. This vector combines membership likelihoods across features, and may be interpreted as the likelihood that the $i^{\text{th}}$ index belongs to *some* feature. In MATLAB, this is realised as `super=min(1,sum(max(S,0),2));` superposition vectors $\mathfrak{s}$ are illustrated in Figures 9, 13, 14, and 22.

If a hard separation is required, then a final step is to threshold the sparse vectors to fix which indices are in or out of a feature. If there is a two-dimensional phase space associated

to the sparse vectors, then one may apply by inspection a manual threshold $\tau$ so that feature $j$ consists of grid cells or nodes with index $i$ satisfying $S_{ij} > \tau$.

For situations where visualisation is not possible, the main thing we desire is to ensure that distinct features are disjoint. We consider three algorithms. Algorithm 4.1 applies the minimum threshold $\tau^{pu}$ to ensure that after thresholding, the columns $S_j$ form a sub-partition of unity: that is $\sum_{j=1}^{r} S_{ij} \leq 1$ for each $i = 1, \ldots, p$. Algorithm 4.2 selects the minimum threshold $\tau^{dp}$ to ensure that after thresholding, the $s_j$ have disjoint support. In each case, the principle of maximum likelihood is then applied to produce a feature vector $a \in \{0, \ldots, r\}^p$ with $a_i = j$ if element $i$ belongs to feature $j$ for $i = 1 \ldots, p$ and $j = 1, \ldots, r$, and $a_i = 0$ if feature $i$ is unassigned. A third algorithm, described at the end of this section, simply applies the above maximum likelihood procedure without thresholding. Define a thresholding transformation $H_\mu : \mathbb{R} \to \mathbb{R}$ by $H_\mu(z) = z$ if $|z| > \mu$ and $H_\mu(z) = 0$ otherwise; $H_\mu$ is applied to vectors elementwise.

**Algorithm 4.1** (Input: vectors $\{s_1, \ldots, s_r\} \subset \mathbb{R}^p$ produced by Algorithm 3.1; output: thresholded $\{s_1, \ldots, s_r\}$ vectors and feature vector $a \in \{0, \ldots, r\}^p$).

1. *For $j = 1, \ldots, r$, set $s_j \to \max\{s_j, 0\}$ to make the vectors nonnegative.*

2. *For each row $i = 1, \ldots, p$, let $\tilde{s}_{i1}, \ldots, \tilde{s}_{ir}$ be the values of $s_{i1}, \ldots, s_{ir}$ in decreasing order.*

3. *Set the threshold $\tau^{pu} := \max_{1 \leq i \leq p, 1 \leq j \leq r}\{\tilde{s}_{ij} : \sum_{l=1}^{j} \tilde{s}_{il} > 1\}$. For each $j = 1, \ldots, r$, hard threshold $s_j \to H_{\tau^{pu}}(s_j)$.*

4. *For each $i = 1, \ldots, p$ let $j^* = \arg\max_{1 \leq j \leq r} s_{ij}$ (resolving ties arbitrarily). If $s_{ij^*} > 0$, set $a_i = j^*$, otherwise set $a_i = 0$.*

**Algorithm 4.2** (Input: vectors $\{s_1, \ldots, s_r\} \subset \mathbb{R}^p$ produced by Algorithm 3.1; output: thresholded $\{s_1, \ldots, s_r\}$ vectors and feature vector $a \in \{0, \ldots, r\}^p$).

1. *For $j = 1, \ldots, r$, set $s_j \to \max\{s_j, 0\}$ to make the vectors nonnegative.*

2. *For each row $i = 1, \ldots, p$, let $\tilde{s}_{i1}, \ldots, \tilde{s}_{ir}$ be the values of $s_{i1}, \ldots, s_{ir}$ in decreasing order.*

3. *Set the threshold $\tau^{dp} := \max_{1 \leq i \leq p} \tilde{s}_{i2}$, and for each $j = 1, \ldots, r$ hard threshold $s_j \to H_{\tau^{dp}}(s_j)$.*

4. *For each $i = 1, \ldots, p$, $j = 1, \ldots, r$, if any $s_{ij} > 0$ for $j = 1, \ldots, r$, set $a_i = j$, otherwise, set $a_i = 0$.*

Because a disjoint partition is a stronger requirement than a partition of unity, one has $\tau^{dp} \geq \tau^{pu}$. If $\tau^{dp} \geq 0.5$, then $\tau^{dp} = \tau^{pu}$. We illustrate both algorithms in Figure 5, where three synthetic sparse vectors $s_1, s_2, s_3 \in \mathbb{R}^{5000}$ are plotted as dotted curves. These vectors have overlapping supports, and $s_1 + s_2 > 1$ on part of the domain. The post-thresholding vectors are shown as solid curves in Figure 5.

One is free to apply other thresholdings to achieve specific outcomes; the above two algorithms are of generic use for feature separation. Thresholds may be applied vectorwise

Figure 5: Upper left: $\tau^{pu}$ from Algorithm 4.1 is indicated by the horizontal dotted line. $H_{\tau^{pu}}(s_1), H_{\tau^{pu}}(s_2), H_{\tau^{pu}}(s_3)$ (solid lines) are equal to $s_1, s_2, s_3$ wherever those vectors are greater than $\tau^{pu}$, and zero elsewhere. Note that $H_{\tau^{pu}}(s_1) + H_{\tau^{pu}}(s_2) + H_{\tau^{pu}}(s_3) \leq 1$ as required for a partition of unity. Lower left: The resulting partition $a$, obtained by performing maximum likelihood on $H_{\tau^{pu}}(s_1), H_{\tau^{pu}}(s_2), H_{\tau^{pu}}(s_3)$. Upper right: $\tau^{dp}$ from Algorithm 4.2 is indicated by the horizontal dotted line. $H_{\tau^{dp}}(s_1), H_{\tau^{dp}}(s_2), H_{\tau^{dp}}(s_3)$ are plotted as solid lines. Lower right: The resulting partition $a$ obtained by performing maximum likelihood on $H_{\tau^{dp}}(s_1), H_{\tau^{dp}}(s_2), H_{\tau^{dp}}(s_3)$.

to individual vectors, or they may be computed from individual vectors and applied uniformly across the vectors as in the two algorithms above. In Section 5.2 we illustrate another option specifically related to coherent sets in nonlinear dynamical systems.

If no thresholding method is effective for a particular system (for example, when feature separations are not clear), maximum likelihood can be applied directly to the output of Algorithm 3.1; that is, one can apply the last step of Algorithm 4.1 to obtain a feature vector $a$. This can be realised in MATLAB by `[m,a]=max(S,[],2); a(m<=0)=0`. The support of $a$ coincides with the positive support of $\max_j s_{ij}$.

## 4.4   Weighting entries of the data vectors

In certain situations, one may wish to apply a weight vector $0 < \nu \in \mathbb{R}^p$ to the calculation of the objective in (1). The weight vector $\nu$ allows one to emphasise (or deemphasise) specific indices $i = 1, \ldots, p$. One situation where this may be appropriate is if one is constructing a transfer operator from a grid with differently sized grid cells; larger cells can be weighted according to their area or volume.

We define a weighted inner product $\langle v, w \rangle_\nu := \sum_{i=1}^p \nu_i v_i w_i$, denote the corresponding

17

norm by $\| \cdot \|_{\ell_{2,\nu}}$, and denote a weighted Frobenius norm by $\|A\|_{F,\nu} := \sqrt{\sum_{i=1}^{p} \sum_{j=1}^{r} \nu_i A_{ij}^2}$. Let $\mathfrak{U}_\nu^{p,r} = \{A \in \mathbb{R}^{p \times r} : \text{ each column of } A \text{ has } \ell_{2,\nu} \text{ norm } 1\}$. We define a weighted $\ell_{1,1}$ matrix norm by $\|A\|_{1,1,\nu} := \sum_{i=1}^{p} \sum_{j=1}^{r} \nu_i |A_{ij}|$. Let $V \in \mathbb{R}^{p \times r}$ be a matrix whose columns are pairwise orthonormal in the inner product $\langle \cdot, \cdot \rangle_\nu$. Problem (1) can be generalised to

$$\underset{\substack{S \in \mathfrak{U}_\nu^{p,r} \\ R \in \mathfrak{S}^r}}{\arg\min} \frac{1}{2} \|V - SR\|_{F,\nu}^2 + \mu \|S\|_{1,1,\nu}. \tag{2}$$

Denoting by $D_\nu$ the diagonal matrix with $\nu$ on the diagonal, we have that $\|A\|_{1,1,\nu} = \|D_\nu A\|_{1,1}$ and $\|A\|_{F,\nu} = \|D_\nu^{1/2} A\|_F$. Substituting $V' = D_\nu^{1/2} V$ and $S' = D_\nu^{1/2} S$, problem (2) is equivalent to

$$\underset{\substack{S' \in \mathfrak{U}^{p,r} \\ R \in \mathfrak{S}^r}}{\arg\min} \frac{1}{2} \|V' - S'R\|_F^2 + \mu \left\| D_\nu^{\frac{1}{2}} S' \right\|_{1,1}. \tag{3}$$

We generalise the soft threshold function $C_\mu$ to $C'_{\mu,\nu} : \mathbb{R}^p \to \mathbb{R}^p$, defined by $[C'_{\mu,\nu}(v)]_i = C_{\mu \nu_i^{1/2}}(v_i)$. We prove in Appendix B that Problem (3) can be solved using Algorithm 3.1 with $V$ replaced with $V'$, $\mu = 0.99/\sqrt{p}$ replaced with $\mu = 0.99/\sqrt{\|\nu\|_1}$ in Step 1, $C_\mu$ replaced with $C'_{\mu,\nu}$ in Step 2, and Step 5 replaced with $S_j \to \text{sign}(\nu^\top S_j) S_j$. After running Algorithm 3.1, set $S = D_\nu^{-1/2} S'$ to obtain a solution for Problem (2). The parameter $\mu$ should be chosen less than $1/\sqrt{\|\nu\|_1}$ to ensure $C'_{\mu,\nu}(D_\nu^{1/2} v)$ is nonzero for every $v \in \mathbb{R}^p$ with $\|v\|_{2,\nu} = 1$; taking $\mu \geq 1/\sqrt{\|\nu\|_1}$ will result in $C'_{\mu,\nu}(D_\nu^{1/2} v)$ being the zero vector when $v$ is the constant $\ell_{2,\nu}$-unit vector.

# 5    Examples

We perform numerical experiments on three systems: the Bickley jet, a turbulence simulation obtained from the Navier-Stokes equations, and North Atlantic ocean currents derived from satellite altimetry. We use the transfer operator (Construction 5 in Table 1) and the dynamic Laplacian (Construction 6 in Table 1) to detect coherent sets.

We can quantify the goodness of fit and sparsity of our solutions using the corresponding parts of the objective (1) being minimised.

- $\ell_2$ **subspace error**: For subspace approximation, note that for $p \times r$ orthogonal $V$, one has $\|V\|_F^2 = r$. Thus, reporting $\frac{1}{r} \|V - SR\|_F^2$ quantifies the subspace approximation error relative to the "zero" estimate $S = 0$; this value should lie between 0 and 1 (lower value means better approximation).

- **Absolute sparsity**: We report the absolute sparsity as $\|S\|_{0,1}/(pr)$, where $\|S\|_{0,1} := \sum_{i=1}^{p} \|S_i\|_0$. The absolute sparsity is the proportion of elements of $S$ that are nonzero, and therefore lies between zero and one (lower value means sparser).

- **Relative sparsity**: We report the relative ($\ell_1$) sparsity $\|S\|_{1,1}/\|V\|_{1,1}$. The $\|\cdot\|_{1,1}$ norm is used as the sparsity penalty term in Algorithm 3.1. The relative sparsity should lie between 0 and 1 (lower value means sparser).

In Sections 5.1 and 5.3 all time integrations are carried out using an explicit Runge-Kutta method of order (4,5) with adaptive step size, as implemented by MATLAB's `ode45` integrator, using the absolute and relative error tolerance $10^{-3}$. The methodologies listed in Table 1 apply to any finite dimension, we restrict ourselves to 2-dimensional examples to present our results with greater transparency.

## 5.1 Bickley jet

The Bickley jet is an idealised system introduced as a model of "banded chaos" by [7], which has become a popular test case for coherent set detection methods [44, 52, 31, 46, 43, 30, 21]. The Bickley jet is a Hamiltonian system modeling a meandering jet with vortices on either side. Its stream function is

$$\psi(x, y, t) = -U_0 L_0 \tanh(y/L_0) + \sum_{i=1}^{3} A_i U_0 L_0 \operatorname{sech}^2(y/L) cos(k_i(x - c_i t)). \tag{4}$$

We use parameter values $U_0 = 62.66$ m/s, $L_0 = 1770$ km, $A_1 = 0.0075$, $A_2 = 0.15$, $A_3 = 0.3$, $c_1 = 0.1446U_0$, $c_2 = 0.205U_0$, $c_3 = 0.461U_0$, $k_1 = 2/r_e$; $k_2 = 4/r_e$, $k_3 = 6/r_e$, $r_e = 6371$ km as in [30]. We consider the associated flow on the initial domain $\mathcal{M} = [0, 20] \times [-3, 3]$, periodic in the $x$-direction, over the time interval $[0, 40]$ days. We calculate the dynamic Laplacian using the "adaptive transfer operator method" of [21] on a Delaunay triangulation (alpha complex with $\alpha = 0.06$) of $\mathcal{M}$ with nodes on an initial grid of $600 \times 180$ points, on the time set $\mathcal{T} = \{0, 40\}$. There are clear eigengaps after the second and eighth eigenvalues, in both the unnormalised and Weyl-normalised plots (Figure 4). These two eigengaps characterise two natural spatiotemporal scales for the Bickley dynamics.

We applied Algorithm 3.1 with $r = 2$ and $r = 8$ to the eigenvectors of the dynamic Laplacian. With $r = 2$, the two eigenvectors are the leading (constant) eigenvector and a second non-constant eigenvector. Algorithm 3.1 produces two sparse vectors whose supports separate the domain into upper and lower pieces; see Figure 6a. The six vortices are also highlighted within the supports of the two sparse vectors. The span of these two sparse vectors is a reasonable approximation of the span of the two leading eigenvectors (subspace error of 11.8%). The sparse vectors have an absolute sparsity of 45.2% and a relative sparsity of 65.0%.

We apply Algorithm 4.1 and 4.2 to obtain a partition of the domain $A_1 \sqcup A_2 \subset M$. Since the sparse vectors in Figure 6a have disjoint support, both algorithms return the same sub-partition, formed by the supports of the two sparse vectors; see Figure 6b. Note that the central jet between the upper and lower halves is also visible (zeroed as dark blue), using only information from the first nontrivial dynamic Laplacian eigenvector. k-means clustering creates a partition of the entire phase space; see Figure 6c.

19

(b) Hard thresholding result from Algorithms 4.1 and 4.2.



(a) Sparse vectors approximating the subspace spanned by the two leading eigenvectors of the dynamic Laplacian ($r = 2$).

(c) Result of k-means (requesting two clusters) applied to the two leading eigenvectors of the dynamic Laplacian.

Figure 6: Bickley jet, with $r = k = 2$.

With $r = 8$, Algorithm 3.1 produces the sparse vectors in Figure 3 (right column). The subspace error is low at 6.4%, as are the absolute sparsity and relative sparsity at 14.1% and 44.6%, respectively. The low absolute sparsity indicates that we have effectively separated the domain into disjoint features, with little overlap between the features.

We can again use either Algorithm 4.1 or Algorithm 4.2 to obtain a partition of the domain; see Figures 7a and 7b respectively. Algorithm 4.1 produces a sub-partition of the domain consisting of almost the entire domain, while Algorithm 4.2 produces a sub-partition of the domain with a larger gap between the upper and lower regions corresponding to a central jet. Both of these results are similar to the k-means clustering result with $k = 8$ clusters, as shown in Figure 7c.

**Remark 5.1.** The k-means algorithm can be viewed as a coarse type of sparse eigenbasis approximation in the following way. Apply k-means to the vectors $v_1, \ldots, v_r$ arising from a spectral clustering, requesting $r$ clusters. For $j = 1, \ldots, r$ define vectors $\hat{s}_j \in \mathbb{R}^p$ as $\hat{s}_{ij} = 1$ if index $i$ is contained in cluster $j$, and $\hat{s}_{ij} = 0$ otherwise. In the "ideal" situation where data points are very clearly separated into distinct clusters, the eigenvectors $v_1, \ldots, v_r$ will be approximately constant, and the (sparse) vectors $\hat{s}_1, \ldots, \hat{s}_r$ will approximately span the same space. By solving a Procrustes problem, one can also infer an orthogonal rotation implicitly

20

(a) Algorithm 4.1      (b) Algorithm 4.2      (c) k-means

Figure 7: State space partition of Bickley dynamics using dynamic Laplacian eigenvectors and Algorithm 3.1 ($r = 8$), followed by: (a) Algorithm 4.1; (b) Algorithm 4.2; (c) k-means with 8 clusters.

performed by k-means to approximately rotate the vectors $v_1, \ldots, v_r$ to the vectors $\hat{s}_1, \ldots, \hat{s}_r$. This relationship with sparse eigenbasis approximation can help explain why k-means can be very effective in some (closer to "ideal") cases where one expects to classify *every* data point (i.e. one *partitions* the dataset into features that exhaust the dataset).

## 5.2   Turbulent flow

We consider a numerical solution to the Navier-Stokes equations in two dimensions with random-in-phase forcing, introduced in [13]. Specifically, we consider the velocity field $u :$ $\mathbb{T}^2 \times \mathbb{R} \to \mathbb{R}^2$ which solves the equations

$$\partial_t u + u \cdot \nabla u = -\nabla p + \nu \Delta u + f$$
$$\nabla \cdot u = 0 \tag{5}$$

where the toral domain $\mathbb{T}^2 = [0, 2\pi] \times [0, 2\pi]$; see [13] for details. We use the Ulam approximation $P$ of the transfer operator $\mathcal{P}$ constructed on a $256 \times 256$ grid, for a flow time of 50 time units, exactly as in [30]. Using the method of [28, 15], we compute singular vectors of the matrix $L$ (in the notation of [28, 15], which is a slight reweighting of the stochastic matrix $P$). The second left and right singular vectors of $L$ are shown in Figure 8. The left



Figure 8: Second left singular vector (left); second right singular vector (right).

21

and right images in Figure 8 may be compared directly with [30, Figure 9(g) and Figure 11(a)], respectively, where in the latter figures, only a binary plot of positive and negative values was displayed. In Figure 8 there are three highlighted coherent regions: two with extreme positive values (red), and one with extreme negative values (dark blue). When applying Algorithm 3.1 to the first two left singular vectors (the leading left singular vector is constant and the second left singular vector is shown in Figure 8 (left)), these three features are separated and highlighted in red; see Figure 9. The sparse basis vectors shown in Figure



Figure 9: Superposition vector $\mathfrak{s}$ (see Section 4.3 for the definition) of the two sparse vectors generated for the turbulence flow using two singular vectors ($r = 2$) (left). Push forward of the left image $\mathfrak{s}$ by $L$ (right).

9 are considerably more informative than either (i) a simple partition based on the sign of the singular vectors or (ii) the result of a hard partitioning method such as k-means. The sparsity quantifiers in the case of $r = 2$ are 31.6% (subspace error), 38.4% (absolute sparsity) and 54.6% (relative sparsity).

   Turbulence is often cited as a canonical example of dynamics operating over a broad range of spatial scales. This fact is backed up by the lack of a clear eigengap for the transfer operator singular value spectrum; see Figure 10 (left). The Weyl rescaling in Figure 10 (right)



Figure 10: Plot of $\log \lambda_r$ vs $r$ (left) and $\log \lambda_r / (r - 1)$ vs $r$ (right) for the turbulence flow.

shows some minor drops, which may provide weak eigengap information. In situations where

22

there is no clear eigengap, we propose to use information from the sparse vectors output by Algorithm 3.1 to decide how many sparse vectors contain "reliable" coherent features. The upper envelope (dotted black) of Figure 11 shows $\text{Min}(S^{(r)})$ vs. $r$. There are several "drops" in this envelope, with the three largest drops at $r = 5, 30$, and $36$. In this turbulence example, we illustrate a refinement of this choice of $r$, simultaneously taking into account the "reliability" of the sparse vectors.

### 5.2.1 A sparse vector heuristic for simultaneously selecting the number of input vectors $r$ and number of features $k$

For some maximum number of columns $r_{\max}$ and each $r = 2, \ldots, r_{\max}$, we apply Algorithm 3.1 to the first $r$ columns of the eigenvector matrix $V$, and denote the resulting arrays $S^{(r)}$. For each $r$ and $k$, let $\text{Min}(S^{(r)}, k) := \sum_{j=1}^{k} - \min_{1 \leq i \leq p} S_{ij}^{(r)}$ be the sum of the minimum values of the first $k$ sparse vectors. In Figure 11, for each $k \leq r$, we plot $\text{Min}(S^{(r)}, k)$ vs. $r$ as a coloured line, for each $2 \leq r \leq r_{max}$. The $k^{\text{th}}$ coloured line shows the "performance"



Figure 11: "Stacked" plot of sparse vector minimum value curves.

(as determined by the minimum value heuristic) of Algorithm 3.1 at separating $k$ features from various $r$ input vectors. The best performance for $k$ features is achieved when a single coloured line attains a minimum of $\text{Min}(S^{(r)}, k)$ across the various values of $r$. For a given $k$, these minimal $r$ are denoted $r_{\min}(k) := \arg\min_{k \leq r \leq r_{\max}} \text{Min}(S^{(r)}, k)$ (when multiple $r$ values give the same $\text{Min}(S^{(r)}, k)$ value, we take $r_{\min}(k)$ as the smallest of these). These "optimal"

$(k, r)$ combinations are shown as black diamonds on Figure 11. Figure 12 shows a plot of $r_{\min}(k)$ vs. $k$.



Figure 12: Turbulence example: $r_{\min}(k)$ vs number of coherent features $k$.

Notice that there are $r_{\min}(k)$ takes on only certain $r$ values. These common $r$ values (e.g. $r = 30, 38, \ldots$ in Figure 11) could be interpreted as representing *natural spatial scales* of the dynamics because one expects the eigenvectors (or singular vectors) input to Algorithm 3.1 to become more oscillatory as $r$ is increased, highlighting increasingly smaller spatial features. One may select any of these optimal $(k, r)$ combinations. To be conservative, in our experiments we select $k$ at the lower end of each "vertical strip" of black diamonds. If we increase $k$ we often simply obtain more coherent sets. For this example, we select $(k, r) = (7, 30)$ and $(21, 38)$ for further analysis; the corresponding sparse superposition vectors and their forward-time images are shown in Figures 13 and 14. From these two



Figure 13: Superposition vector $\mathsf{s}$ of the 7 most reliable features at initial time (left) according to the minimum value reliability criterion ($r = 30$) and its forward-time image (right).

figures one already sees 7 (resp. 21) clearly highlighted coherent sets in red that remain coherent after the flow time of 40 time units. The span of the collection of sparse vectors

24

Figure 14: Superposition vector $\mathfrak{s}$ of the 21 most reliable features at initial time (left) according to the minimum value reliability criterion ($r = 38$) and its forward-time image (right).

in both the case of $r = 30$ and that of $r = 38$ produce reasonable approximations of the span of the corresponding leading eigenvectors with a subspace error of 10.4% (resp. 9.8%). The sparse vectors have an absolute sparsity of 7.7% (resp. 7.1%) and a relative sparsity of 26.8% (resp. 25.2%).

### 5.2.2 Extraction of a subpartition from sparse vector output

One could apply a single global threshold to our selected sparse vectors according to Algorithms 4.1 and 4.2. We found that because Figure 13 is already a collection of vectors with mostly disjoint supports, both of these algorithms produced a very low threshold with little change to the sparse vectors. On the other hand, the sparse vectors that are combined in Figure 14 have "overlaps" of large values in small parts of phase space, causing Algorithms 4.1 and 4.2 to produce very high thresholds close to 1. In place of these thresholding algorithms, one could easily apply a global threshold manually by inspecting Figure 13 or 14; for example a threshold of around 0.7–0.75 would produce coherent sets with small boundary lengths relative to enclosed area at the initial and final times.

To end this section we apply thresholding according to the Cheeger ratio. Such a thresholding has been applied to specific single vectors [16, 20, 23, 30], but because Algorithm 3.1 has automatically separated *all* features, we may apply this thresholding to either the superposition vector $\mathfrak{s}$ or to each individual vector sparse vector $s_j$, $j = 1, \ldots, k$. Given a vector (corresponding to a discrete representation of a real-valued function on a manifold $M$), we select a threshold by choosing the level set that minimises a quantity we call the *scale-invariant Cheeger ratio*. Let $M$ be a $d$-dimensional manifold and let $T : M \to M$ be a volume-preserving transformation representing our dynamics. For simplicity of presentation, we assume that one application of $T$ represents the nonlinear dynamics over the full finite-time interval and that we compute the coherence via the Cheeger ratio at only the initial and final time; see [16, 23] for more detail and generalisations. Let $\Gamma \subset M$ be a $d-1$ dimensional

25

submanifold disconnecting $M$ into two pieces $M_1, M_2$. The scale-invariant Cheeger ratio is

$$h^D(\Gamma) := \frac{\ell_{d-1}(\Gamma) + \ell_{d-1}(T(\Gamma))}{2(\min\{\ell_d(M_1), \ell_d(M_2)\})^{\frac{d-1}{d}}}, \tag{6}$$

where $\ell_{d-1}$ and $\ell_d$ are $(d-1)$- and $d$-dimensional volume measure, respectively. In this example we treat each sparse vector $s_j$ separately, and each $\Gamma$ will be selected as a level set of $s_j$, $j = 1, \ldots, k$. In particular we choose the level set value so that (6) is minimised. In terms of numerics, having selected a $(k,r)$ combination, for each $1 \leq j \leq k$, we linearly interpolate $s_j$ on a fine regular grid on $M$ to obtain a function $u_j : M \to \mathbb{R}$. For threshold values $\tau$ ranging from $\min_x u_j(x)$ to $\max_x u_j(x)$ we calculate $h^D(\Gamma_\tau)$, where $\Gamma_\tau$ is the level set $\{x \in M : u_j(x) = \tau\}$, and select the $\tau$ yielding the minimum value. We use MATLAB's `contourc` command to calculate level sets; `isosurface` is the corresponding command for three dimensions. This takes approximately two minutes for 21 interpolated sparse vectors in $\mathbb{R}^{65536}$ on an Intel Core i7-6700 processor with 16GB RAM. Note that in contrast to Algorithm 4.1 or 4.2, we calculate a *separate threshold* for each vector before applying a hard thresholding.

We apply this thresholding procedure to the $(k,r)$ combination (7,30) whose superposition of sparse vectors is shown in Figure 13); the result is the partition in Figure 15. All features obtained are highly coherent, though two features are perhaps overly small. The



Figure 15: Partition $(k,r) = (7,30)$ of phase space showing 7 coherent sets at initial- (left) and forward- time (right) obtained by thresholding.

scale-invariant dynamic Cheeger ratio for these two small sets has a very flat minimum over a large range of thresholds and so these two sets could be further enlarged without changing appreciably the Cheeger ratio (and therefore the finite-time coherence).

We also apply this thresholding procedure to the $(k,r)$ combination (21,38), whose superposition vector $\mathfrak{s}$ is shown in Figure 14. The resulting thresholded sets are displayed in Figure 16, coloured according to their reliability (a reliability index of 1 is the most reliable (coherent)). We retain all the coherent sets identified in Figure 15, and gain another 15 coherent sets. In this case, while some of the additional coherent sets are relatively small, many of these sets cannot be significantly increased without introducing filamentation; that is, they do not have flat minima for their Cheeger ratios.

Figure 16: Partition $(k, r) = (21, 38)$ of phase space showing 21 coherent sets at initial- (left) and forward- time (right) obtained by thresholding.

### 5.2.3 Extraction of a partition with k-means

The use of sparse vectors via Algorithm 3.1 clearly outperforms k-means clustering of the singular vectors. As an example, we apply k-means to the $(k, r)$ combination (21,38), namely using 38 eigenvectors and requesting $k = 21$ clusters; see Figure 17. We add one extra cluster to give k-means a chance to find the complement of 21 coherent sets as a 22nd coherent set The k-means clustering correctly detects some coherent sets, notably many of those in Figure 15, but also produces many more highly filamented sets, and misses several very coherent features visible in Figure 16. To a large extent, k-means fails in difficult examples such as turbulence because k-means *imposes* a partition of the phase space into coherent sets, when a partition is not reasonable. The superposition of the sparse vectors in Figure 14 highlights the coherent sets in a much more precise way than k-means can, and we can further refine Figure 14 by thresholding, as in Figure 16. Increasing $k$ causes k-means to add more filamented features, while decreasing $k$ causes k-means to miss some of the coherent features in Figures 14 and 16. We also experimented with the use of the k-means sum of squared errors (the objective that k-means attempts to minimise) to determine the number of clusters, but we failed to find any "kinks" to identify any values for $k$, as described for example in [49].

## 5.3 North Atlantic surface currents

As our final example, we consider a real-world dataset, namely geostrophic AVISO [1] velocity fields. We consider the 90-day period 15 January 2015 to 15 April 2015; the velocity fields are updated every 2 hours. We consider a synthetic initial $291 \times 224$ Mercator grid of tracers on the domain $\mathcal{M} = [-65°, -35°] \times [30°, 50°]$. After removing trajectories initialised on land, we numerically integrate the remaining 56205 trajectories, linearly interpolating the velocity field between the AVISO-provided velocity fields spaced one day apart. Using trajectory positions[7] every 3 days, we compute the dynamic Laplacian using the "adaptive transfer

---

[7]In particular, derivatives of the velocity field or flow map need not be estimated.

Figure 17: Initial and forward-time plot of the k-means partition of state space into 22 clusters ($r = 38$).

operator approach" of [21] on alpha complexes[8] with zero Neumann boundary conditions; see [21] for details on implementation.

The second eigenfunction highlights the Gulf of St Lawrence in the extreme northwest of the 15 January 2015 domain as a coherent set evolving distinctly from the rest of the domain (not shown), and the third eigenfunction separates the domain at the Gulf Stream; see the sharp red/green interface in Figure 18. One may directly compare Figure 18 (lower) with [2, Figure 9]. The latter figure shows a satellite-based map of sea surface temperature at 5 April 2015. The Gulf Stream is seen as a very strong temperature gradient emanating from the east coast of the USA. There is a very close match with the strong red/green interface (and also the red/orange interface between -50 and -40 longitude) in the eigenvector in Figure 18 (lower). This demonstrates that the strong temperature gradient of the surface ocean is highly correlated with coherent transport of water parcels on the surface.

The spectrum for this dataset, shown in Figure 19 (left) reveals little separation of scales; the largest drops in the Weyl rescaling (Figure 19; right) are at $r = 3$ and $r = 30$. Therefore, we again select values for $r$ and $k$ using the heuristic described in Section 5.2.1. We construct the cumulative minimum value plot (Figure 20) and plot $r_{\min}(k)$ in Figure 21. Our initial choice of $r_{\max} = 50$ as our maximum number of eigenvectors was possibly too small to identify many of the coherent sets, because $r_{\min}(k) = 50$ for almost all values of $2 \leq k \leq 50$. Increasing $r_{\max}$ to 100, the coloured curves for $k = 1, \ldots, 40$ all start to plateau before $r_{\max} = 100$, and so we consider this $r_{\max}$ sufficiently large. We only plot the coloured curves for $k = 1, \ldots, 40$, to simplify the figure.

---

[8]With the lowest alpha such that the alpha complexes are connected.

Figure 18: Third eigenfunction at 15 January 2015 (upper) and its forward-time image at 15 April 2015 (lower).

As discussed in Subsection 5.2.1, we want to consider the first $k$ value in a run of $k$ values with the same $r_{\min}(k)$. We take $r = 96, k = 36$. We repeated the experiment with various $k$ and $r_{\min}(k)$ pairs with $r_{\min}(k)$ around 77–81, but the resulting regions were similar to a subset of those for $r = 96, k = 36$. The span of these 96 sparse vectors produced an approximation of the original subspace with an error of 3.9%. The absolute sparsity is 3% and the relative sparsity is 19.1%. We plot the sum of the resulting sparse vectors in Figure 22, together with their forward time image. We obtain many coherent features, as well as some partial coherent features at the boundary. A zoom of two of the sparse vectors highlighting two different eddies at 28 February 2015 is shown in Figure 23. The black dots in Figure 23 indicate trajectory locations used in the construction of the dynamic Laplacian. Note that we resolve high-resolution estimates of the eddy shapes using the finite-element approach of [21] from lower-resolution trajectory data.

The eddy field in Figure 22 can be compared directly with a backward finite-time Lya-

Figure 19: Spectrum of the dynamic Laplacian for the North Atlantic dataset. Plot of $\lambda_r$ vs $r$ (left) and $\lambda_r/(r-1)$ vs $r$ (right).

punov exponent field in [2, Figure 10 (upper right)]. The FTLE field suggests a Gulf Stream feature as in Figure 18 (lower) and some eddies of similar size as Figure 22; in both cases, the FTLE field is not cleanly identifying these figures. Figure 4 [45] shows fields of encounter volume, exact diffusivity, long-time diffusivity, and diffusive timescale at 11 January 2015 on a domain approximately -75 to -45 degrees longitude by 35 to 42.5 degrees latitude. These four (spatially smaller) fields can be directly compared with the corresponding part of the larger domain in Figure 22 (upper). The three eddies highlighted in Figure 22 (upper) in the subdomain of [45, Figure 4] appear in the latter figure to varying extents, but not as sharply as in Figure 22. Other features one might find in [45, Figure 4] do not appear in Figure 22 (upper), possibly because they are less coherent than the many other features highlighted. The partition produced by k-means again proved ineffective; see Figure 24.

Figure 20: North Atlantic negative proportion stacked bar chart



Figure 21: $r_{\min}$ vs number of coherent features.

Figure 22: Superposition vector $\mathfrak{s}$ of the 36 most reliable features for the North Atlantic dataset using Algorithm 3.1 with $r = 96$ at 15 January 2015 (upper) and its forward-time image at 15 April 2015 plotting only final trajectory locations (lower).

Figure 23: Zoom of two sparse vectors produced by Algorithm 3.1 with $(k, r) = (36, 96)$, plotted at 28 February 2015. The black dots are locations of trajectory data used to compute the dynamic Laplace operator.



Figure 24: Initial (left) and forward-time (right) plot of the k-means partition of state space for the North Atlantic dataset into 37 clusters using $r = 96$ eigenvectors.

# 6 Discussion

We proposed a simple, powerful post-processing of the output of spectral clustering to separate individual features or clusters. Spectral clustering outputs a truncated eigenspectrum of some Markov- or Laplace-type matrix or operator, with associated eigenbasis. Features or clusters are encoded in the vectors of this eigenbasis, but may be mixed together within several eigenvectors. We adapted sparse PCA methodologies to identify an approximate sparse basis for the given eigenbasis. The use of sparsity (i) efficiently and reliably separates individual features and (ii) produces a natural cutoff to identify those data points that do not belong to any cluster.

Spectral clustering often relies on the existence of an eigengap to determine (i) the number of clusters or features and (ii) where to truncate the eigenspectrum. In many situations, there is no clear eigengap. We proposed a refinement of the standard eigengap procedure in Section 4.2.1 that takes into account the spectral structure of Markov and Laplace-type operators, to correctly scale the eigenspectrum to allow fair comparison of gaps along the spectrum. We further introduced vector-based heuristics in Section 4.2.2 based on the output of our sparse eigenbasis approximation (SEBA) algorithm to suggest a good choice of the size of eigenbasis $r$. As described in Section 4.2, these vector-based heuristics determine natural *spatial scales*, while the eigenvalues at eigengaps determine natural *temporal scales*. A rejection criterion to help determine suitable numbers of features or clusters $k$ was proposed in Section 4.1. We described a method for simultaneously selecting $k$ and $r$ in Section 5.2.1.

We believe that our techniques are a useful toolkit for general spectral clustering problems. One specific motivation was the identification of almost-invariant and coherent sets in time-dependent nonlinear dynamical systems; in particular, simply and reliably extracting many such sets and determining corresponding natural spatial and temporal scales. We applied our techniques to three example nonlinear systems.

Firstly, the Bickley jet, which possesses a modest number of strongly coherent sets that were easily separated and extracted using eigenfunctions of the dynamic Laplace operator. In this example, other post-processing techniques such as k-means clustering also perform very well. Secondly, turbulent, randomly forced Navier-Stokes flow on a 2-torus, with coherent sets at multiple spatial scales and with no clear temporal scale separation. We identified good combinations for the "number of coherent sets" and the "eigenbasis size" and produced well-separated coherent sets through the superposition vector $\mathfrak{s}$. In contrast, post-processing methods such as k-means performed poorly, as expected for the situation where much of the phase space should not be classified as coherent. Thirdly, we analysed trajectory data in the North Atlantic Ocean derived from satellite altimetry. We found a large separation of natural spatial scales with a jump from scales commensurate with the entire ocean basin, to features around one degree of longitude across. Leading eigenvectors identified the Gulf of St Lawrence as a coherent set strongly dynamically disconnected from the rest of the domain, and the Gulf Stream as a persistent large-scale coherent transport barrier. No intermediate scales were identified; the next natural spatial scale occurred at the scale of mesoscale eddies, at which we identified numerous eddies.

# 7    Acknowledgements

# A    Code

## A.1    MATLAB code for Algorithm 3.1

```matlab
function S=SEBA(V)

%V is pxr matrix (r vectors of length p as columns, assumed orthonormal)
%S is pxr matrix with columns approximately spanning the column space of V

[V,~]=qr(V,0); %Enforce orthonormality
[p,r]=size(V);
mu=0.99/sqrt(p);
S=zeros(size(V));
Rnew=speye(r);   %Initialise rotation
R=0;

while norm(Rnew-R)>1e-14 %tolerance may be decreased
    R=Rnew;
    Z=V*R';
    %Thresholding to solve sparse approximation problem
    for i=1:r
        S(:,i)=sign(Z(:,i)).*max(abs(Z(:,i))-mu,0);
        S(:,i)=S(:,i)/norm(S(:,i));
    end
    %Polar decomposition to solve Procrustes problem
    [P,~,Q]=svd(S'*V,0);
    Rnew=P*Q';
end

%Choose correct parity of vectors and scale so largest value is 1.
for i=1:r
    S(:,i)=S(:,i)*sign(sum(S(:,i)));
    S(:,i)=S(:,i)/max(S(:,i));
end

%Sort so that most reliable vectors appear first (starting with column 1
    in S).
[~, I] = sort(min(S), 'descend');
S = S(:, I);
```

## A.2 MATLAB code for Algorithm 4.1

```matlab
function [S,A,taupu] = subpartition_unity(S)
%SUBPARTITION_UNITY Apply hard thresholding to obtain a sub-partition of
% unity from the columns of S. Also create a maxmimum-likelihood
% sub-partition.

S=max(S,0);                             % Take non-negative part
S_descend=sort(S,2,'descend');          % Sort each row in descending order
S_sum=cumsum(S_descend,2);              % Must be <=1 for partition of unity
taupu=max([S_descend(S_sum>1);0]);      % Largest element where row sum > 1
S(S<=taupu)=0;                          % Apply hard thresholding

[M,A]=max(S,[],2);  % Find column with largest element, by row
A(M==0)=0;          % Suppress zero rows
```

## A.3 MATLAB code for Algorithm 4.2

```matlab
function [S,A,taudp] = disjoint_support(S)
%DISJOINT_SUPPORT Apply hard thresholding to obtain disjoint supports for
% the columns of S. Also create a maximum-likelihood sub-partition.

S=max(S,0);                         % Take non-negative part of S
S_descend=sort(S,2,'descend');      % Sort to find largest element in each row
taudp=max(S_descend(:, 2));         % Take largest non-leading element
S(S<=taudp)=0;                      % Apply hard thresholding

[M,A]=max(S,[],2);  % Find column with largest element, by row
A(M==0)=0;          % Suppress zero rows
```

## A.4 MATLAB code for producing a "natural spatial scales" plot of $(k, r)$ combinations

```matlab
function SpatialScalesPlot(V)
%V is pxR matrix (r vectors of length p as columns, assumed orthonormal)

R=size(V,2);
column_minima=nan(R-1,R);
for i=1:R-1
    S=SEBA(V(:,1:i+1));    % Column minima are in descending order
    column_minima(i,1:i+1)=-min(S);
end
minimum_sums=cumsum(column_minima,2);

[~,best_r]=min(minimum_sums);
cla reset;
plot(best_r(1:R)+1,(1:R),['o','-'],'Color','k','MarkerSize',4,...
```

```
                'MarkerFaceColor','k');% Plot best r for each k=1,..,r
```

## A.5   MATLAB code for producing a "stacked" sparse vector minimum value plot

Note that the first seven lines coincide with the code in Section A.4, so these lines need not be recomputed if one is creating both plots.

```matlab
function MinValStackedPlot(V,kmax)
%V is pxR matrix (r vectors of length p as columns, assumed orthonormal)
%kmax is number of lines to output, assumed <= R

R=size(V,2);
column_minima=nan(R-1,R);
for i=1:R-1
    S=SEBA(V(:,1:i+1));    % Column minima are in descending order
    column_minima(i,1:i+1)=-min(S);
end
minimum_sums=cumsum(column_minima,2);

[best_sum,best_r]=min(minimum_sums);
cla reset; hold on
for j=1:kmax
    iStart=max(j-1,1);
    plot(iStart+1:R,minimum_sums(iStart:R-1,j),'o-',...
        'MarkerSize',2,'Linewidth',1.5);
end
plot(best_r(1:kmax)+1,best_sum(1:kmax),'kd','MarkerSize',4,...
    'MarkerFaceColor','k');% Plot best r for each k
plot([2,2:kmax],minimum_sums([1,R:R:R*(kmax-1)]),'k:',...
    'MarkerSize',2.5) % Plot upper envelope
hold off
```

# B   Proof that the expression $(2)$ can be minimised as outlined in Section 4.4

We wish to solve problem (3), which is equivalent to problem (2). We can apply the same alternating optimisation scheme as in Section 3. With fixed $S'$, problem (3) is the same as problem (1) with $V$ replaced with $V'$ and $S$ replaced with $S'$. This can be solved using Step 2 of Algorithm 3.1.

With fixed $R$, we use the property of the (weighted) Frobenius norm that $\|A\|_{F,\nu} = \|AR\|_{F,\nu}$ for every matrix $A \in \mathbb{R}^{p \times r}$ and orthogonal $R \in \mathbb{R}^{r \times r}$. To see this, note that $\|A\|_{F,\nu}^2 = \text{Trace}(A^\top D_\nu A)$, so

$$\|A\|_{F,\nu}^2 = \text{Trace}(A^\top D_\nu A(RR^\top)) = \text{Trace}(R^\top A^\top D_\nu AR) = \|AR\|_{F,\nu}^2.$$

This lets us reformulate problem (3), for fixed $R$, as

$$\underset{S' \in \mathfrak{U}^{p,r}}{\arg\min} \frac{1}{2}\|V'R^\top - S'\|_F^2 + \mu\|D_\nu^{\frac{1}{2}} S'\|_{1,1} = \underset{S' \in \mathfrak{U}^{p,r}}{\arg\min} \sum_{j=1}^{r} \frac{1}{2}\|Z_j - S_j'\|_2^2 + \mu\|D_\nu^{\frac{1}{2}} S_j'\|_1, \quad (7)$$

where $Z_j$ is the $j$th column of $V'R^\top$ and $S_j'$ is the $j$th column of $S'$. Since $\|Z_j\|_2 = \|S_j'\|_2 = 1$, we have $\|Z_j - S_j'\|_2 = \|Z_j\|_2 + \|S_j'\|_2 - 2Z_j^\top S_j' = 2 - 2Z_j^\top S_j'$, so problem (7) is equivalent to $\arg\max_{S' \in \mathfrak{U}^{p,1}} Z_j^\top S_j' - \mu\|D_\nu^{1/2} S_j'\|_1$. By a slight modification of equations (5)-(7) in [37], taking $p = 1$, $A = Z_j^\top$ and $x = 1$ and replacing $\gamma\|z\|_1$ with $\sum_{i=1}^{n} \mu\nu_i^{1/2}|z_i|$, it can be shown that our problem (7) is solved for each $j$ by setting $S_j' = C'_{\mu,\nu}((V'R^\top)_j)/\|C'_{\mu,\nu}((V'R^\top)_j)\|_2$, exactly as in Step 3 of Algorithm 3.1 with $V$ replaced with $V'$ and $C_\mu$ replaced with $C'_{\mu,\nu}$.

Lastly, we prove that $C'_{\mu,\nu}(D_\nu^{1/2}z)$ is nonzero for every $z \in \mathbb{R}^p$ with $\|z\|_{2,\nu} = 1$, if and only if $\mu < 1/\sqrt{\|\nu\|_1}$. We start by showing that if $\mu < 1/\sqrt{\|\nu\|_1}$ then $C'_{\mu,\nu}(D_\nu^{1/2}z)$ is nonzero for every $z \in \mathbb{R}^p$ with $\|z\|_{2,\nu} = 1$. To see this, choose any such $z$, then

$$1 = \|z\|_{2,\nu} = \sqrt{\sum_{i=1}^{p} \nu_i z_i^2} \le \sqrt{\sum_{i=1}^{p} \nu_i \max_{1 \le j \le p} z_j} = \sqrt{\|\nu\|_1} \max_{1 \le j \le p} z_j.$$

That is, there is some index $j$ with $z_j \ge 1/\sqrt{\|\nu\|_1}$. Recalling that $\nu_i > 0$ for every $i = 1, \ldots, p$, $\mu < 1/\sqrt{\|\nu\|_1}$ implies $|\nu_j^{1/2} z_j| - \mu\nu_j^{1/2} > 0$, i.e. the $j$th element of $C'_{\mu,\nu}(D_\nu^{1/2}z)$ is nonzero as required. To prove the reverse direction, note that the vector $c \in \mathbb{R}^p$ with every element equal to $1/\sqrt{\|\nu\|_1}$ has $\|c\|_{2,\nu} = 1$. But for $\mu \ge 1/\sqrt{\|\nu\|_1}$, $|\nu_i^{1/2} c_i| - \mu\nu_i^{1/2} \le 0$ for every $i = 1, \ldots, p$, so $C'_{\mu,\nu}(D_\nu^{1/2}c) = 0$, and the proof is complete.

# References

[1] https://www.aviso.altimetry.fr/en/data.html and http://marine.copernicus.eu/.

[2] S. Balasuriya, N. T. Ouellette, and I. I. Rypina. Generalized Lagrangian coherent structures. *Physica D*, 372:31–51, 2018.

[3] R. Banisch and P. Koltai. Understanding the geometry of transport: Diffusion maps for Lagrangian trajectory data unravel coherent sets. *Chaos*, 27(3):035804, 2017.

[4] T. Berry and T. Sauer. Spectral clustering from a geometrical viewpoint. Technical report, 2015.

[5] A. d'Aspremont, F. Bach, and L. El Ghaoui. Optimal solutions for sparse principal component analysis. *J. Mach. Learn. Res.*, 9:1269–1294, 2008.

[6] A. d'Aspremont, L. El Ghaoui, M. Jordan, and G. Lanckriet. A direct formulation for sparse PCA using semidefinite programming. *SIAM Review*, 49(3):434–448, 2007.

[7] D. del-Castillo-Negrete and P. J. Morrison. Chaotic transport by Rossby waves in shear flow. *Physics of Fluids A: Fluid Dynamics*, 5(4):948–965, 1993.

[8] M. Dellnitz and O. Junge. On the approximation of complicated dynamical behavior. *SIAM Journal on Numerical Analysis*, 36(2):491–515, 1999.

[9] A. Denner, O. Junge, and D. Matthes. Computing coherent sets using the Fokker-Planck equation. *Journal of Computational Dynamics*, 3(2):163–177, 2016.

[10] P. Deuflhard and C. Schütte. Molecular conformation dynamics and computational drug design. In *Applied Mathemetics Entering the 21st Century. Proceedings ICIAM*, pages 91–119, 2004.

[11] P. Deuflhard and M. Weber. Robust Perron cluster analysis in conformation dynamics. *Linear Algebra and its Applications*, 398:161–184, 2005.

[12] D. L. Donoho. For most large underdetermined systems of linear equations the minimal 1-norm solution is also the sparsest solution. *Communications on Pure and Applied Mathematics*, 59(6):797–829, 2006.

[13] M. Farazmand and G. Haller. Attracting and repelling Lagrangian coherent structures from a single computation. *Chaos*, 23(2):023101, 2013.

[14] G. Froyland. Statistically optimal almost-invariant sets. *Physica D*, 200(3):205–219, 2005.

[15] G. Froyland. An analytic framework for identifying finite-time coherent sets in time-dependent dynamical systems. *Physica D*, 250:1–19, 2013.

[16] G. Froyland. Dynamic isoperimetry and the geometry of Lagrangian coherent structures. *Nonlinearity*, 28(10):3587–3622, 2015.

[17] G. Froyland and M. Dellnitz. Detecting and locating near-optimal almost-invariant sets and cycles. *SIAM Journal on Scientific Computing*, 24(6):1839–1863, 2003.

[18] G. Froyland, C. Horenkamp, V. Rossi, N. Santitissadeekorn, and A. Sen Gupta. Three-dimensional characterization and tracking of an Agulhas ring. *Ocean Modelling*, 52-53:69–75, 2012.

[19] G. Froyland, C. Horenkamp, V. Rossi, and E. van Sebille. Studying an Agulhas ring's long-term pathway and decay with finite-time coherent sets. *Chaos*, 25(8):083119, 2015.

[20] G. Froyland and O. Junge. On fast computation of finite-time coherent sets using radial basis functions. *Chaos*, 25(8):087409, 2015.

[21] G. Froyland and O. Junge. Robust FEM-based extraction of finite-time coherent sets using scattered, sparse, and incomplete trajectories. *SIAM Journal on Applied Dynamical Systems*, 17(2):1891–1924, 2018.

[22] G. Froyland, O. Junge, and P. Koltai. Estimating long-term behavior of flows without trajectory integration: The infinitesimal generator approach. *SIAM Journal on Numerical Analysis*, 51(1):223–247, 2013.

[23] G. Froyland and P. Koltai. Estimating long-term behavior of periodically driven flows without trajectory integration. *Nonlinearity*, 30(5):1948–1986, 2017.

[24] G. Froyland and E. Kwok. A dynamic Laplacian for identifying Lagrangian coherent structures on weighted Riemannian manifolds. *Journal of Nonlinear Science*, Jun 2017.

[25] G. Froyland, S. Lloyd, and N. Santitissadeekorn. Coherent sets for nonautonomous dynamical systems. *Physica D*, 239(16):1527–1541, 2010.

[26] G. Froyland, K. Padberg, M. H. England, and A. M. Treguier. Detection of coherent oceanic structures via transfer operators. *Phys. Rev. Lett.*, 98:224503, 2007.

[27] G. Froyland and K. Padberg-Gehle. Almost-invariant and finite-time coherent sets: directionality, duration, and diffusion. In W. Bahsoun, C. Bose, and G. Froyland, editors, *Ergodic Theory, Open Dynamics, and Coherent Structures*, pages 171–216. Springer, 2014.

[28] G. Froyland, N. Santitissadeekorn, and A. Monahan. Transport in time-dependent dynamical systems: Finite-time coherent sets. *Chaos*, 20(4):043116, 2010.

[29] G. Froyland, R. M. Stuart, and E. van Sebille. How well-connected is the surface of the global ocean? *Chaos*, 24(3):033126, 2014.

[30] A. Hadjighasem, M. Farazmand, D. Blazevski, G. Froyland, and G. Haller. A critical comparison of Lagrangian methods for coherent structure detection. *Chaos*, 27(5):053104, 2017.

[31] A. Hadjighasem, D. Karrasch, H. Teramoto, and G. Haller. Spectral-clustering approach to Lagrangian vortex detection. *Phys. Rev. E*, 93:063107, 2016.

[32] J. Han, K. Xiong, and F. Nie. Orthogonal and nonnegative graph reconstruction for large scale clustering. In *Twenty-Sixth International Joint Conference on Artificial Intelligence*, pages 1809–1815, 2017.

[33] Z. Hu, G. Pan, Y. Wang, and Z. Wu. Sparse principal component analysis via rotation and truncation. *IEEE Transactions on Neural Networks and Learning Systems*, 27(4):875–890, 2016.

[34] A. K. Jain. Data clustering: 50 years beyond k-means. *Pattern Recognition Letters*, 31(8):651–666, 2010.

[35] A. K. Jain and R. C. Dubes. *Algorithms for Clustering Data*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1988.

[36] I. T. Jolliffe, N. T. Trendafilov, and M. Uddin. A modified principal component technique based on the LASSO. *Journal of Computational and Graphical Statistics*, 12(3):531–547, 2003.

[37] M. Journée, Y. Nesterov, P. Richtárik, and R. Sepulchre. Generalized power method for sparse principal component analysis. *J. Mach. Learn. Res.*, 11:517–553, 2010.

[38] D. Karrasch and J. Keller. A geometric heat-flow theory of Lagrangian coherent structures. *arXiv e-print*, 2016.

[39] L. Kaufman and P. J. Rousseeuw. Clustering by means of medoids. In *Reports of the Faculty of Technical Mathematics and Informatics*. Delft University of Technology, 1987.

[40] P. Koltai. A stochastic approach for computing the domain of attraction without trajectory simulation. *Disc. Cont. Dynam. Sys., Supplement*, page 854–863, 2011.

[41] M. Lee, J. Lee, H. Lee, and N. Kwak. Membership representation for detecting block-diagonal structure in low-rank or sparse subspace clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1648–1656, 2015.

[42] C. Lu, S. Yan, and Z. Lin. Convex sparse spectral clustering: Single-view to multi-view. *IEEE Transactions on Image Processing*, 25(6):2833–2843, 2016.

[43] K. Padberg-Gehle and C. Schneide. Network-based study of Lagrangian transport and mixing. *Nonlinear Processes in Geophysics*, 24(4):661–671, 2017.

[44] I. I. Rypina, M. G. Brown, F. J. Beron-Vera, H. Koçak, M. J. Olascoaga, and I. A. Udovydchenkov. On the Lagrangian dynamics of atmospheric zonal jets and the permeability of the stratospheric polar vortex. *Journal of the Atmospheric Sciences*, 64(10):3595–3610, 2007.

[45] I. I. Rypina, S. G. Llewellyn Smith, and L. J. Pratt. Connection between encounter volume and diffusivity in geophysical flows. *Nonlinear Processes in Geophysics*, 25(2):267–278, 2018.

[46] K. L. Schlueter-Kuck and J. O. Dabiri. Coherent structure colouring: identification of coherent structures from sparse data using graph theory. *Journal of Fluid Mechanics*, 811:468–486, 2017.

[47] B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput.*, 10(5):1299–1319, 1998.

[48] M. A. Stremler, S. D. Ross, P. Grover, and P. Kumar. Topological chaos and periodic braiding of almost-cyclic sets. *Physical Review Letters*, 106(11):114101, 2011.

[49] C. A. Sugar, L. A. Lenert, and R. A. Olshen. An application of cluster analysis to health services research: empirically defined health states for depression from the SF-12. Technical Report 203, Stanford University, 1999.

[50] U. von Luxburg. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, 2007.

[51] H. Weyl. Das asymptotische Verteilungsgesetz der Eigenwerte linearer partieller Differentialgleichungen (mit einer Anwendung auf die Theorie der Hohlraumstrahlung). *Mathematische Annalen*, 71(4):441–479, 1912.

[52] M. O. Williams, I. I. Rypina, and C. W. Rowley. Identifying finite-time coherent sets from limited quantities of Lagrangian data. *Chaos*, 25(8):087408, 2015.

[53] S. X. Yu and J. Shi. Multiclass spectral clustering. In *Proceedings of the Ninth IEEE International Conference on Computer Vision-Volume 2*, page 313, 2003.

[54] H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15(2):265–286, 2006.