

Nonlinear feature extraction of the stratosphere

Marijn van der Meer

Professor: Olivier Verscheure, Supervisors: Eniko Szekely, Raphaël de Fondeville

Swiss Data Science Center, EPFL Lausanne, Switzerland

Abstract—In this project, we perform nonlinear dimension reduction to explore the behaviour, patterns, and variability of the northern polar vortex. The northern polar vortex is a climate phenomenon associated with outbreaks of extremely cold temperatures in the Northern Hemisphere. However, expertise in seasonal forecasting for its effect on the surface weather in the Northern Hemisphere is currently lacking. A better understanding of the northern polar vortex variability and seasonality could allow for better predictability which, in the long term, could lead to improved weather predictions. To analyse the vortex dynamics, we perform nonlinear dimension reduction using an approach that approximates the eigenvalues of the Laplace-Beltrami operator over a manifold with a weighted Laplacian over an adjacency graph. The first Laplacian eigenmaps found using this approximation follow closely the patterns of Empirical Orthogonal Functions (EOF) on the data. When diverging from EOF patterns, Laplacian eigenmaps seem to capture organised and consistent nonlinear dynamics EOF struggle to detect. Furthermore, Laplacian eigenmaps retrieve an annual seasonality in our data. To capture further dynamics and search for other types of variability, extending Laplacian eigenmaps to nonlinear Laplacian spectral Analysis is suggested.

I. INTRODUCTION

The Earth's atmosphere is divided into four layers; the troposphere, closest to the Earth's surface; the stratosphere that houses the climate phenomenon we are interested in; the mesosphere characterised by colder low-density air and finally the thermosphere, a thin warm layer. In both hemispheres, the atmosphere over the polar regions is characterised during winter by a counter-clockwise rotating vortex of winds whose intensity peaks during mid-winter. In Northern Eurasia, this vortex encompasses an area of about 1'000 kilometres in diameter. During winter, the northern polar vortex expands, sending cold air southward. This happens fairly regularly and is often associated with outbreaks of bitterly cold temperatures in the Northern Hemisphere. Furthermore, approximately every two years, sudden stratospheric warmings (SSW) disturb

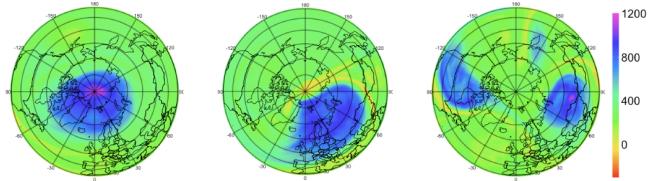


Fig. 1: The northern polar vortex in winter. Images show undisturbed polar vortex (left) and a disturbed polar vortex during two types of SSW: displacement (middle) and split (right). Colour scale shows the value of Potential Vorticity (PV) i.e., the absolute circulation of air enclosed between two surfaces, of the polar vortex. Image reproduced from [2].

the wind field and cause it to deform (Fig. 1). These distortions are responsible for long-term circulation anomalies in the troposphere which can strongly impact the weather at the surface up to three months following their appearance [1]. Gaining a thorough understanding of the northern polar vortex, its behaviour, variability, and seasonality is crucial in accurate weather forecasting on the surface in the Northern Hemisphere.

Unfortunately, as is explained by Sigmond M. *et al.* in [3], expertise in seasonal forecasting in extra-tropical regions is lacking. The current best seasonal forecasting expertise stems from the knowledge of El Niño Southern Oscillation's remote influences [3]. However, the further away we go from the tropical Pacific Ocean, the smaller El Niño's impact is on seasonal oscillations, particularly for Northern Eurasia [3]. So, the accuracy in forecasts gets poorer for these regions. This has encouraged the seasonal forecast community to look for further sources of predictability in extra-tropical regions [3]–[5]. As a solution, Sigmond M. *et al.* propose that insight into the state of the stratosphere can act as a source of increased seasonal predictability [3].

As aforementioned, circulation anomalies of the northern polar vortex are related to anomalies in the troposphere, with a strong effect on the weather that may last up to three months. In this project, we explore the north-

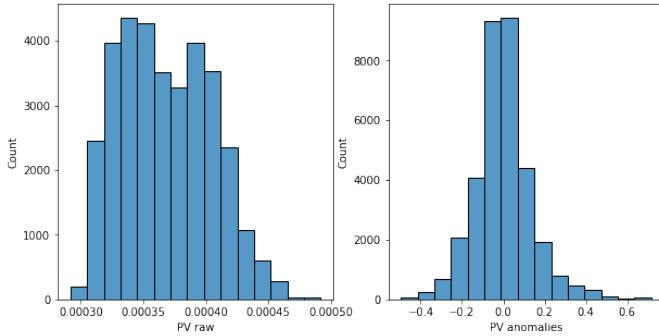


Fig. 2: Data distribution of the first dimension in the raw and anomalies data i.e., the first column $Y_{[:,1]} \in \mathbb{R}^N$ of Y^R and Y^A .

ern vortex' variability and anomalies with a long-term objective of producing improved long-range weather predictions [2]. To this end, we perform nonlinear feature extraction on Northern Hemisphere stratospheric measures to construct low-dimensional representations of the climate phenomenon and study its seasonality and variability.

II. STATE OF THE ART

In climate studies, the problem of constructing low-dimensional representations is classically addressed with Empirical Orthogonal Functions (EOF) analysis, known also as Principle Component Analysis (PCA) in statistics and physics. One way to use EOF is to examine the evolution in time of possible spatial modes of variability [6]. The eigenvalues of the spatially weighted anomaly covariance matrix of a field provide a measure of the percentage of variance explained by each pattern. Atmospheric processes are generally such that most of the variance is held within the first modes [6]. EOF techniques are usually accompanied by various pre-processing methods such as band pass filtering, running averaging, or seasonal partitioning [7]. To avoid any extensive need for pre-processing, several groups [8], [9] propose to perform dimension reduction using Laplacian eigenmaps. Laplacian Eigenmaps are a different, locality preserving, and nonlinear solution, inspired by the connections between the graph Laplacian and the Laplace Beltrami operator on a manifold, a topological space locally resembling an Euclidean space in the neighbourhood of each point.

In this project, we choose to apply this approach to stratosphere data of the northern polar vortex to perform nonlinear dimension reduction and study the variability and anomalies of the vortex.

III. MODELS AND METHODS

A. Data

To study the northern polar vortex, we use Potential Vorticity (PV), which quantifies the absolute circulation of a package of air enclosed between two isentropic i.e., constant-entropy, surfaces (Fig. 1). It is applied in atmospheric dynamics and meteorology to describe rotating fluids with vertical stratification. PV values in the troposphere are generally low, but increase rapidly from the troposphere to the stratosphere due to the substantial shift in static stability. The PV data for this project was extracted from an ERA-Interim reanalysis. Data was sourced every six hours every winter (November-April) from 1979 to 2018 for latitudes above 30° . This yielded $N = 33'960$ time snapshots of PV values. Furthermore, the data covered a regular longitude/latitude grid with 0.75° spacing, resulting in $D = 38'400$ grid cells. Overall, initial sampled data was of shape $X \in \mathbb{R}^{D \times N}$. We did not directly perform our analysis on X , but rather used three transformed datasets that originated from X :

- Raw data: N PV values measured four times per day during winter (November-April) from 1979 to 2018, $Y^R \in \mathbb{R}^{N \times d}$. The $d = 1001$ dimension corresponds to the coefficients of a functional basis whose elements are derived by solving the Laplacian on a spherical cap. The coefficients are then obtained by projecting the data on the elements of the basis.
- Anomalies data: for each grid cell, anomalies $Y^A \in \mathbb{R}^{N \times d}$ were calculated by computing a smooth estimate of the seasonality, that was then subtracted from data. The residuals were then normalised to have unit variance. An example of the difference between anomalies and raw values range can be seen in Fig. 2.
- Functional basis: $\Psi \in \mathbb{R}^{D \times d}$ with d basis columns $\psi_{[:,i]} \in \mathbb{R}^D \quad \forall i = 1, \dots, d$. Each row of a basis columns corresponds to a specific [longitude, latitude] grid coordinate.

To plot the geospatial representation of one time-step t of PV values $X_{[:,t]} \in \mathbb{R}^D$ as in Fig. 1, the columns $\psi_{[:,i]}$ of the basis $\Psi \in \mathbb{R}^{D \times d}$ are combined with the coordinates of one row of the data i.e., one time-step t , $Y_{[t,:]} \in \mathbb{R}^d$ with the basis as:

$$X_{[:,t]} = \sum_{i=1}^d Y_{[t,i]} \psi_{[:,i]} \quad (1)$$

Algorithm 1 Laplacian Eigenmaps decomposition [9]

- 1: **INPUT:** $Y \in \mathbb{R}^{N \times d}$, matrix of data samples. For simplicity, we write $y_i = Y_{[i,:]}$ for the i^{th} row of Y i.e., one data point.
- 2: **OUTPUT:** m Laplacian eigenmaps of Laplacian dimension reduction
- 3: Construct adjacency matrix $A \in \mathbb{R}^{N \times N}$: put edge between i and j if y_i and y_j are close i.e, either by
 - 1) ϵ -neighbourhood: e.g., if $\|y_i - y_j\|_2^2 < \epsilon$
 - 2) n -nearest neighbours of y_i
- 6: Construct weight matrix $W \in \mathbb{R}^{N \times N}$:
 - 1) Heat-kernel: $W_{ij} = e^{\frac{-\|y_i - y_j\|^2}{t}}$ for connected nodes and $W_{ij} = 0$ otherwise; t is the bandwidth of the kernel
 - 2) Simple kernel: $W_{ij} = 1$ for connected nodes and $W_{ij} = 0$ otherwise
- 9: Compute eigenvalues and eigenfunctions for the generalised eigendecomposition problem

$$L\phi = \lambda D\phi \quad (2)$$

where D is the diagonal degree matrix and its entries are column sums of W , $D_{ii} = \sum_j W_{ji}$ and $L = D - W$ is the positive semi-definite Laplacian matrix. Let $\phi_0, \dots, \phi_{m-1}$ be the eigenvector solutions to equation 2, ordered according to their eigenvalues ($L\phi_i = \lambda_i D\phi_i$). Then for $\lambda_0 = 0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_{m-1}$, we leave out λ_0 and take the next m eigenvectors for embedding in an m -dimensional Euclidean space: $y_i \rightarrow \Phi = (\phi_1(i), \dots, \phi_m(i))$

Where $Y_{[t,:]}$ is either a row of Y^R or Y^A depending on whether we take the raw or anomalies data. This vector is plotted with the help of the longitude and latitude coordinates.

For simplicity in the rest of this paper, we write $y_i = Y_{[i,:]}$ for the i^{th} row of Y i.e., one data point. We further use the superscript notation of R and A to indicate entities originating from respectively the raw or anomalies dataset.

B. Laplacian eigenmaps

The Laplacian minimizer aims to learn a low-dimensional representation of a dataset under the assumption that it lies on a low-dimensional manifold \mathcal{M} embedded in $\mathbb{R}^l, l > 0$. Under the assumption that there exists a map $\Phi : \mathcal{M} \rightarrow \mathbb{R}^m, l > m > 0$ preserving local distances i.e., points close on \mathcal{M} remain

close in \mathbb{R}^m , then the optimal map Φ is given by the m eigenfunctions ϕ_i corresponding to the m smallest (strictly positive) eigenvalues λ_i of the Laplace-Beltrami operator $L_{\mathcal{M}}$ [10]. $L_{\mathcal{M}}$ is a generalised Laplace operator on functions defined on manifolds \mathcal{M} in Euclidean space, and is defined as:

$$\Delta f = \nabla \cdot \nabla f \quad (3)$$

The eigenfunctions of $L_{\mathcal{M}}$ provide an optimal embedding for the manifold \mathcal{M} [9]. In practice the manifold is unknown and thus the eigenfunctions of the Laplace-Beltrami operator cannot be estimated. To address this, Belkin M. and Niyogi P. suggested an algorithm that builds a graph integrating neighbourhood information of the dataset [9]. An adjacency A matrix over the neighbourhood of data points serves to approximate the manifold \mathcal{M} . The Laplace-Beltrami operator $L_{\mathcal{M}}$ is approximated by a corresponding weighted graph Laplacian with an appropriately chosen weight function that constructs a weight matrix W from the adjacency matrix A (Algorithm 1).

1) *Eigenfunctions:* the embedding $\Phi = [\phi_0, \dots, \phi_{m-1}] \in \mathbb{R}^{N \times m}$ is given by the matrix of eigenvectors corresponding to the lowest eigenvalues λ_i . Those eigenvectors are solution to

$$L\phi_i = \lambda_i D\phi_i, \forall i \in \{1, \dots, m\} \quad (4)$$

and are interpreted as time series that can be considered as nonlinear analogues of EOF time series ζ_i [7]. Laplacian eigenfunctions ϕ_i form a basis for outlining relevant quantities on the manifold where the use of eigenvectors ϕ_i with smallest eigenvalues λ_i is analogous to sampling features that vary slowly on the nonlinear manifold. In this way, noise is reduced while avoiding overfitting [9], due to the fact that the eigenvalues λ_i do not, as in EOF, measure explained variance. Laplacian eigenvalues can be geometrically interpreted as an average gradient of the corresponding eigenvectors on the manifold [7], [11].

2) *Convergence:* The embedding map Φ given by the graph Laplacian eigenfunctions can be considered as a discrete approximation to continuous maps given by the eigenfunctions of the Laplace-Beltrami operator $L_{\mathcal{M}}$, intrinsically defined on the entire manifold \mathcal{M} . Under the assumptions that the map Φ is sufficiently regular, \mathcal{M} is uniformly sampled and that the function creating the weight matrix W is Gaussian, Belkin M. and Niyogi P. give a formal proof in [12] that the graph Laplacian

converges to the Laplace-Beltrami operator on \mathcal{M} [10].

3) Locality preserving map: The Belkin M. and Niyogi P. argue that the this embedding Φ preserves local distances i.e., points close on \mathcal{M} remain close in the embedding. The advantage of that is that the *locality-preserving character of the Laplacian eigenmap algorithm makes it relatively insensitive to outliers and noise. It is also not prone to short-circuiting, as only the local distances are used* [9]. By preserving local information in the embedding, their algorithm implicitly emphasises the natural clusters in the data and it exhibits stability with respect to the embedding [9]. Thus, as long as the embedding is isometric the representation does not change. Here, an isometric embedding is defined as a smooth embedding $\Phi : \mathcal{M} \rightarrow \mathcal{N}$ in Riemannian geometry where \mathcal{M} and \mathcal{N} are two manifolds and Φ preserves the length of curves.

4) Optimal weights: The problem of

$$L\phi_i = \lambda D\phi_i \quad (5)$$

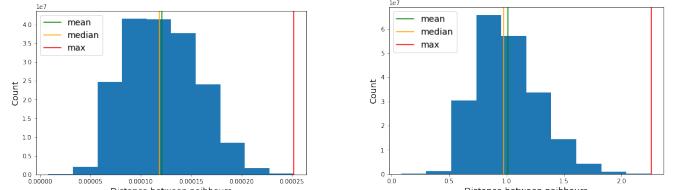
This problem can be reformulated as finding

$$\underset{\Phi^T D \Phi = 1}{\operatorname{argmin}} \Phi^T L \Phi \quad (6)$$

where the diagonal degree matrix D measures the importance of a vertex, as the bigger the values in matrix D_{ii} , the more neighbours a vertex y_i has [9]. The constraint $Y^T D Y = 1$ prevents collapse onto a subspace of dimension less than $m - 1$ [9]. The Laplace-Beltrami operator is intimately related to heat flow and in order to ensure that the approximation matrix is positive semi-definite, computing the graph Laplacian with a heat-kernel is optimal [9]. Thus, the weight matrix W is constructed such that $W_{ij} = e^{-\frac{\|y_i - y_j\|^2}{t}}$. Compared to a simple kernel, the heat-kernel heavily penalises neighbouring points that are mapped far apart in equation 6. It further approximates distances on the manifold when $\lim_{t \rightarrow 0}$ (in the limit of large data) as it then tends to Dirac's δ -function.

We adapted Algorithm 1 presented in [9], [12], the following way to construct our embedding map.

- Adjacency matrix: the adjacency graph on the raw and anomalies dataset was computed using the n -nearest neighbour approach, where n was chosen as 10% and 20% of the total number of nodes i.e., N , and the Euclidean distance as distance measure. The final adjacency matrix $A \in \mathbb{R}^{N \times N}$ was made



(a) Raw data

(b) anomalies data

Fig. 3: Histogram of non-zero distances across connected components created by step 1 of Algorithm 1 for raw (left) and anomalies (right) data. Different choices for the bandwidth of heat-kernel in red (max), yellow (median) and green (mean).

symmetric, meaning that if point y_i is a neighbour of y_j , then $A_{ji} = 1$ and $A_{ij} = 1$.

- Simple kernel: as a baseline, we started with the simple kernel as mentioned in Algorithm 1. The adjacency matrix A was used as weight matrix $W \in \mathbb{R}^{N \times N}$ where $W_{ij} = 1$ for connected nodes and 0 otherwise. In this case, all neighbouring vertices y_j of a vertex y_i carry the same weight regardless of their distance to y_i .
- Heat-kernel: from the simple binary kernel, we moved on to a heat-kernel, a pairwise measure of similarity that decays according to the squared Euclidean distance between nodes y_i and y_j .

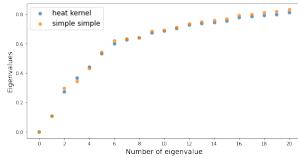
$$K(y_i, y_j) = e^{-\frac{(y_i - y_j)^2}{t}} \quad (7)$$

where t is the bandwidth. In our case, we chose to compute three different heat-kernels using the mean, median and maximum over all distances in the adjacency graph i.e., the distances in the connected components as the bandwidth t (Fig. 3).

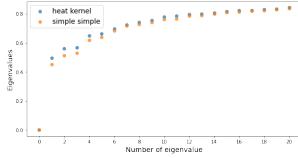
5) Geospatial patterns: To recover spatial patterns in the original data space i.e., geographical coordinates \mathbb{R}^D , we performed an operation inspired by the spatiotemporal recovery process in [7]

$$\hat{Y}_{[i,:]} = Y^T D \phi_i \in \mathbb{R}^d \quad (8)$$

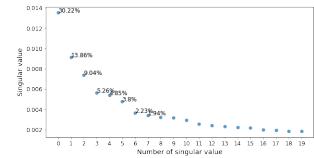
where $D \in \mathbb{R}^{N \times N}$ is the degree matrix as in Algorithm 1, $Y \in \mathbb{R}^{N \times d}$ the raw or anomalies PV values and $\phi_i \in \mathbb{R}^N$ a Laplacian eigenvector. Then, $\hat{Y}_{[i,:]} \in \mathbb{R}^d$ are transformed back to the original data space using equation 1 with the functional basis Ψ . This transformation does not maximise explained variance as EOF do, as it was proposed that this might not be the optimal way to recover dynamics, but following the manifold has a greater potential for recovering dynamically significant patterns [7], [13], [14].



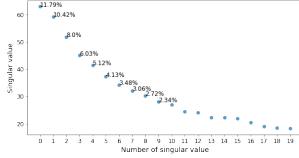
(a) Laplacian eigenvalues on raw data



(c) Laplacian eigenvalues on anomalies data



(b) EOF eigenvalues on raw data



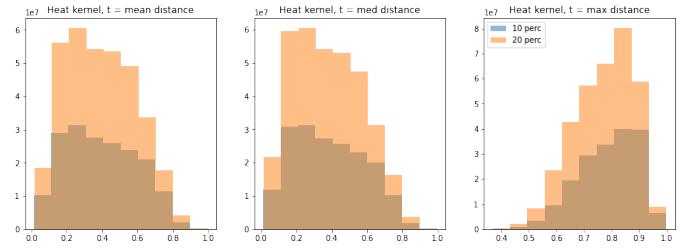
(d) EOF eigenvalues on anomalies data

Fig. 4: Leading $m = 20$ eigenvalues λ_i for raw (a) and anomalies (c) data ordered in increasing value. Eigenvalues were computed for the simple binary kernel (orange) and heat-kernel ($n = 10\%$ and $t = \text{mean}$) (blue) as in Fig. 5. Those are compared to the twenty first leading singular values found using with EOF, a linear dimension reduction. Here dimension reduction was done using singular value decomposition (SVD) on centered but not scaled data. Singular values are ordered according to their decreasing value for raw (b) and anomalies (d) data. For the leading singular values, their explained variance ratio is plotted above their points. For raw and anomalies data, the main two eigenvectors respectively explain approximately 30% and 14%, and 11% and 10% of the variance. See Fig. 9, 11 for the rest of the explained variance.

IV. RESULTS

A. Construction of Laplacian Eigenmaps

Several heat-kernels with different values as a bandwidth t were computed (Fig. 5) on an adjacency graph constructed with the nearest-neighbour method ($n = 10\%$ and 20% of total samples). These kernels were compared to a simple weight matrix with only binary values, 1 being the value attributed to connected nodes. We were looking to create a kernel that approximates distances on the manifold with a smooth transition from 0 to 1. A heat-kernel with $n = 10\%$ as the number of neighbours and the mean or median as a bandwidth provides this desired shape (Fig. 5). Using $n = 20\%$ did not seem to improve our results, while using the maximum as bandwidth tended to follow a left-tailed bell curve, more similar to the binary kernel. Following this result, we chose to use a heat-kernel with $t = \text{mean}$ and $n = 10\%$. With this kernel to construct the Laplacian and degree matrix as in Algorithm 1, we computed the eigenvectors and eigenvalues for both the raw and anomalies dataset. As can be seen in



(a) Raw data

1e7 Heat kernel, $t = \text{mean}$ distance

1e7 Heat kernel, $t = \text{med}$ distance

1e7 Heat kernel, $t = \text{max}$ distance

1e8 Heat kernel, $t = \text{max}$ distance

(b) anomalies data

1e7 Heat kernel, $t = \text{mean}$ distance

1e7 Heat kernel, $t = \text{med}$ distance

1e8 Heat kernel, $t = \text{max}$ distance

Fig. 5: Histogram of non-zero values of heat-kernels computed for Algorithm 1 according to equation 7 with different values of bandwidth t and different adjacency matrices. Adjacency matrix produced on raw (a) and anomalies (b) data using 10% (light orange) and 20% (darker orange) as the number of nearest neighbours. As bandwidth t : mean (left), median (middle) and maximum (right) of all distances over the adjacency graph.

Fig. 4a, 4c, the only non-zero eigenvalue is the first i.e., λ_0 , corresponding to eigenvector $\phi_0 = \mathbf{1}$. This indicates that the adjacency graph constructed in step 1 of Algorithm 1 is fully connected. As this trivial solution collapses all vertices of the adjacency graph on ϕ_0 [9] we ignored this eigenvector for the rest of the procedures.

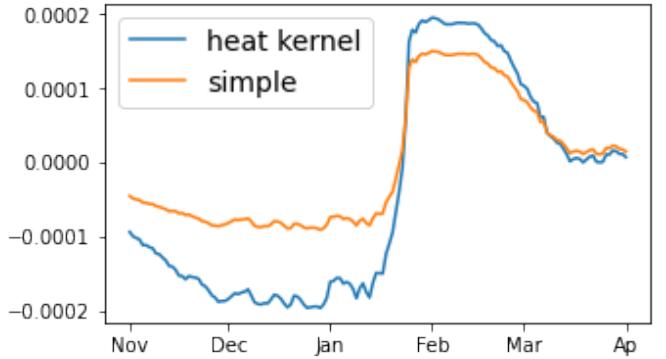
As previously outlined in Sect. III-B1, eigenfunctions ϕ_i , $\forall i \in \{1, \dots, m\}$, are interpreted as time-series that are nonlinear analogs of EOF time-series. Except revealing some highly oscillating patterns, those plots are hard to read but full time-series for the first twenty eigenvectors from 1979 to 2018 can be found in the Appendix (Fig. VI.1, VI.2). To improve readability, a single winter (from November 2008 to April 2009), was isolated and plotted for all twenty eigenvectors (Fig. 7, 8). This winter is known for having a strong split event in the polar vortex as illustrated in Fig. 1. For the raw data, the first eigenvector ϕ_1^R (Fig. 6a) is of particular interest as it seems to capture a seasonality, seeing how a peak appeared round the middle of

the winter and disappeared again afterwards. To see whether this pattern persisted over the whole forty years time-period, all winters were stacked over each other (Fig. 6b). From this, it seems that a peak arises for almost all years around mid to end winter. Nevertheless, the moment of appearance of the bell-curve shows very high variability. To continue the search for a regular pattern, the time-series' power-spectrum was computed. As our data spans only from November to April, time-series were zero-padded to expand to a full year. The power spectrum for ϕ_1^R (Fig. 6c) captures two main peaks, one for annual and one for biannual frequency. For smaller frequencies, there is again a lot of variability. The biannual frequency arose from the zero padding, but the second might be picking up on some seasonality in the rise of that bell-shaped peak throughout the year.

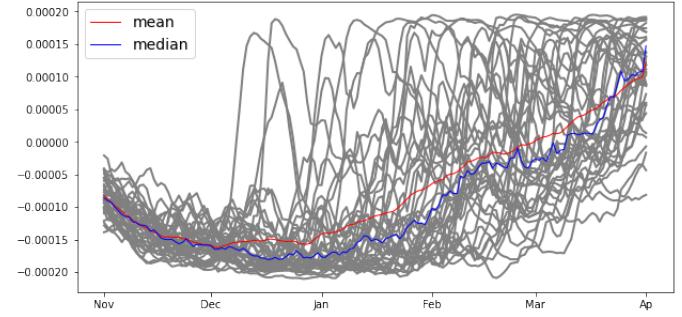
Looking at the other eigenvectors for raw and anomalies data, no flagrant pattern like the first non zero eigenvector on the raw data, ϕ_1^R , seem to arise. Their power-spectra were nevertheless computed to see whether they detected relevant signals (Fig. 7, 8). Interestingly, ϕ_2^R shows a high triennial frequency. Otherwise, except a seasonality of maximum once per year, the raw eigenvectors do not seem to pick up on frequencies above that e.g, no biennial frequencies which we might for example look for when searching for SSW events. The same applies for anomalies data.

B. Geospatial patterns of eigenmaps

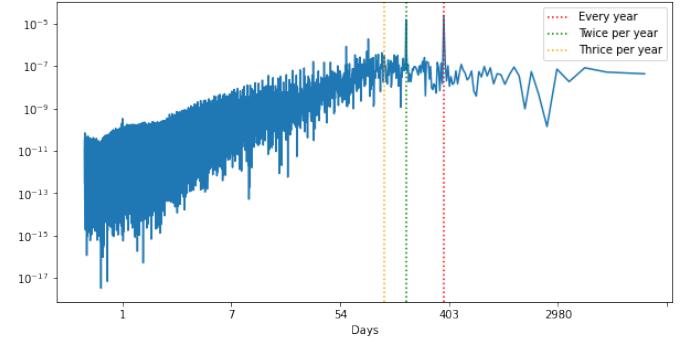
The results of nonlinear dimension reduction with Laplacian eigenmaps were compared to eigenvectors ζ_i found by EOF and calculated using singular value decomposition (SVD) (Fig. 9, 10, 11 and 12). For SVD, the data was centered but not scaled for each feature before applying the SVD. To compare Laplacian eigenmaps to EOF, equation 8 was used to recover spatiotemporal patterns from Laplacian eigenvectors $\Phi = [\phi_1, \dots, \phi_{m-1}]$ and equation 1 to transform those and the EOF ζ_i to the original data space. For raw data, until approximately Laplacian ϕ_{10}^R , eigenmaps ϕ_i^R show similar patterns as the EOF. Above that, EOF patterns become more fragmented and noisier (e.g. ζ_{17}^R Fig. 10), while Laplacian eigenmaps remain less noisy and clustered into larger forms. For anomalies data, EOF ζ_i^A also diverge from Laplacian eigenmaps around ϕ_{10}^A (Fig. 11).



(a) (RAW) First non-zero Laplacian eigenvector ϕ_1^R from the heat-kernel (blue) and Simple-kernel (orange) on raw data from November 2008 to April 2009.



(b) (RAW) First non-zero Laplacian eigenvector ϕ_1^R overlapping over 40 winters. In red the mean over the values for each time-step and in blue the median.



(c) (RAW) Power spectrum of the first non-zero Laplacian eigenvector ϕ_1 as in Fig. VI.1.

Fig. 6: First non-zero eigenvector ϕ_1^R computed on raw data using the heat-kernel ($n = 10\%$ and $t = \text{mean}$) plotted for one year of November 2008 to April 2009 (a), overlapping for forty winters from 1979 to 2018 (b), and its power-spectrum (c). The power-spectrum was computed as in IV-A.

C. Two-dimensional spatial patterns of eigenmaps

Due to their seemingly highly oscillating and difficult interpretable patterns, the leading five eigenvectors

ϕ_i , $\forall i \in \{1, \dots, 5\}$ were plotted against each other to explore two-dimensional patterns and nonlinear dependencies (Fig. 13, 14). Within these, three isolated winter trajectories (Nov-April) were highlighted: 2008-2009, known for being a year with a strong split event; 2010-2011, a winter with a strong displacement; 1998-1999, a winter starting with a strong displacement followed by a strong split (c.f. Fig. 1 for an example of a split and displacement event). For both datasets, the two-dimensional trajectories are highly nonlinear. For the raw data, interesting shapes show trajectories starting and ending in different regions. Those are U-shaped (e.g. ϕ_1^R vs ϕ_2^R) or butterfly-shaped (e.g. ϕ_1^R vs ϕ_5^R). Another relevant shape has trajectories starting and ending in approximately the same region; those are more circular (e.g. ϕ_3^R vs ϕ_4^R). Anomalies data show primarily circular patterns (e.g. ϕ_1^A vs ϕ_2^A) and some U-shaped (e.g. ϕ_1^A vs ϕ_4^A). Looking at the three isolated trajectories, they seem to start and end in the same region for almost all combinations of anomalies eigenvectors. We further notice that for both the raw and anomalies dataset, the three isolated trajectories follow distinct paths for almost all two-dimensional combinations.

Some combinations of eigenvectors i.e., ϕ_1^R vs ϕ_2^R and ϕ_5^R , ϕ_1^A vs ϕ_2^A and ϕ_4^A , were isolated for both datasets and dynamics were divided into regions of interest following approximately the three sampled yearly trajectories. The temporal ranges $\Delta t = [t_1, \dots, t_t]$ for those regions were extracted, and geospatial features on the original dataset Y^R or Y^A were averaged over these time-spans Δt . This provided reconstructed average spatiotemporal features in the original data space once equation 1 was applied using the functional basis Ψ (Fig. 15). Patterns extracted from the anomalies dataset show some kind of rotating double crescent-shape with a stronger circle of PV values appearing and disappearing above Northern-America. For the raw dataset, the idea of the bell-curve of eigenvector $\phi_{1,raw}$ in Fig. 7 seems to reappear here in Fig. 15a, 15b in the form of a disk-shaped polar region of contrasting PV values appearing in Northern-Europe and moving to Northern-America before disappearing.

V. DISCUSSION AND CONCLUSION

In this project we performed nonlinear dimension reduction on stratospheric northern polar vortex data using the approach proposed in [9]. The aim of this was to extract low-dimensional features that might help the study of the behaviour and variability of the northern

polar vortex in order to improve long-term and range weather forecasts in the Northern Hemisphere. The approach we used creates a low-dimensional embedding given by the leading eigenfunctions of the weighted Laplacian that approximates the Laplace-Beltrami operator over a manifold. The most relevant results were obtained using a nearest-neighbour method to construct an adjacency graph and an adapted appropriate heat-kernel as a weight function on this to create a weighted graph. Laplacian eigenmaps that were found, once projected back to the original data space, resemble closely the first ten eigenvectors ζ_i obtained by EOF analysis. When Laplacian eigenmaps ϕ_i start to diverge from EOF ζ_i , the latter become noisier, with a pattern diverging into four or more parts while eigenmaps are more organised and display a more consistent structure. EOF force orthogonality and search for linearity, while the physical system we study is probably nonlinear. The fact that Laplacian eigenmaps stay more concise when EOF get more dispersed, might indicate that the Laplacian method is picking up on dynamics in nonlinear patterns with low variance, which EOF struggle to find. The actual shapes of the projections of eigenmaps on the original data space are otherwise hard to interpret, especially given our lacking background in climate science, but it is reassuring that they capture the first same nodes as EOF and then clean, non-random, nonlinear, patterns when they finally diverge.

In addition, the first eigenmaps for raw data seemed to pick a seasonality, with average geospatial reconstructed patterns forming strong shapes of PV values above different regions of the poles (Northern-America or Europe) before disappearing. This seems reinforced by the power-spectra of both datasets, where the maximal important frequencies do not exceed an annual frequency. We were initially wondering whether this method would detect higher frequencies, such as biennial which might indicate SSW events, but these seem not to be present, or they might eclipsed by the more prominent biannual and annual signal peaks. In further work, other combinations of eigenvectors independent of ϕ_1 should be explored to see whether they capture something else than seasonality.

To search for additional dynamics that might not be detected by Laplacian eigenmaps and open our analysis, we extended our analysis to nonlinear Laplacian spectral analysis (NLSA) [15]–[17], a technique designed to better capture time varying dynamics. Compared to

Laplacian eigenmaps, NLSA adds an initial time-lagged embedding step (Takens embedding space) with a time-window τ before computing the Laplacian of the graph. This additional step transforms $Y \in \mathbb{R}^{N \times D}$ to a new matrix $\hat{Y} \in \mathbb{R}^{(N-\tau) \times (D\tau)}$ and is said to allow a better capture of dynamics [7]. Due to time constraints, we were only able to gather results for a temporal embedding window of $\tau = 7$ days. For such a short time-window, we collect results that resemble closely what we had without using a Takens embedding space (Fig. VI.3, VI.4, VI.5, VI.6). To explore whether dynamics are truly changed using NLSA, we should have used a longer time-window of 1 month or more, but this will be explored in further work.

In conclusion, using Laplacian eigenmaps we are able to obtain a relevant nonlinear low-dimensional embedding of the northern polar vortex behaviour. Those eigenmaps capture interesting different dynamics from EOF when the latter get noisy. Nevertheless, those patterns require the interpretation of a climate scientist to evaluate their relevance. Furthermore, Laplacian eigenmaps show informative seasonal patterns in the vortex's behaviour, suggesting that we capture annual or biannual dynamics. Nevertheless, to get relevant information about specific patterns such as SSW or the effects of the polar vortex on the weather at the surface, a more comprehensive understanding of the dynamics of the vortex is necessary. For this, further extensive research is required, notably the path of nonlinear Laplacian spectral analysis (NLSA) should more thoroughly be explored.

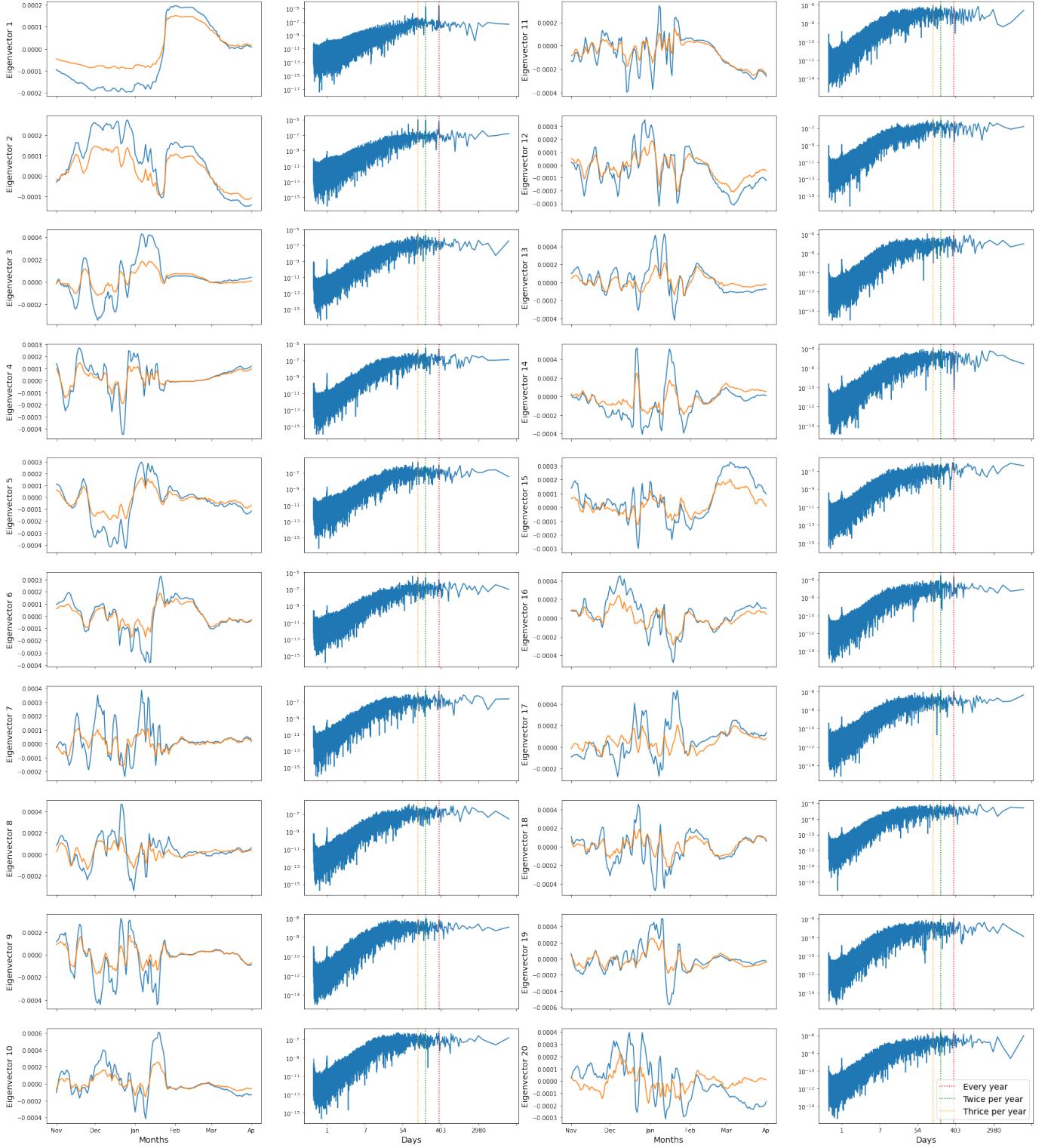


Fig. 7: (RAW) Leading twenty non-zero eigenfunctions ϕ_i^R for the winter time interval of November 2008 to April 2009, with their associated power-spectrum, an estimation of the spectral density of the signals. Power-spectrums were computed on eigenvectors as in Fig VI.1 but zero-padded in order to complete the data for the missing months. Eigenfunctions were computed on the raw data using the heat-kernel from Fig. 5 ($n = 10\%$ and $t = \text{mean}$) (blue) and the simple binary kernel (orange). Vertical lines in the power spectra indicate frequencies of once, twice and thrice per year respectively in yellow, green and red.

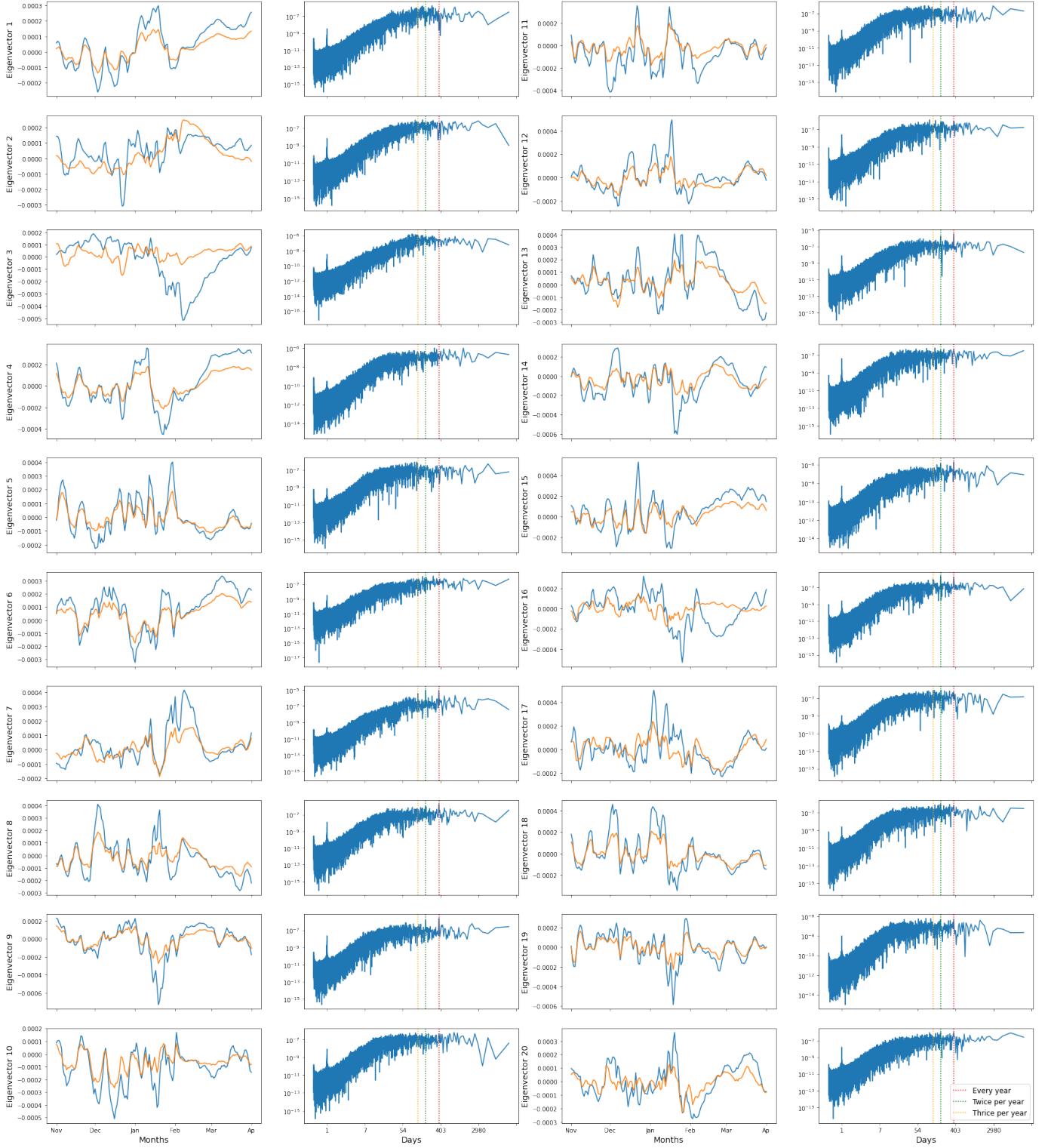


Fig. 8: (ANOMALIES) Leading twenty non-zero eigenfunctions ϕ_i^A for the winter time interval of November 2008 to April 2009, with their associated power-spectrum, an estimation of the spectral density of the signals. Power-spectra were computed on eigenvectors as in Fig VI.2 but zero-padded in order to complete the data for the missing months. Eigenfunctions were computed on the anomalies data using the heat-kernel from Fig. 5 ($n = 10\%$ and $t = \text{mean}$) (blue) and the simple binary kernel (orange). Vertical lines in the power spectra indicate frequencies of once, twice and thrice per year respectively in yellow, green and red.

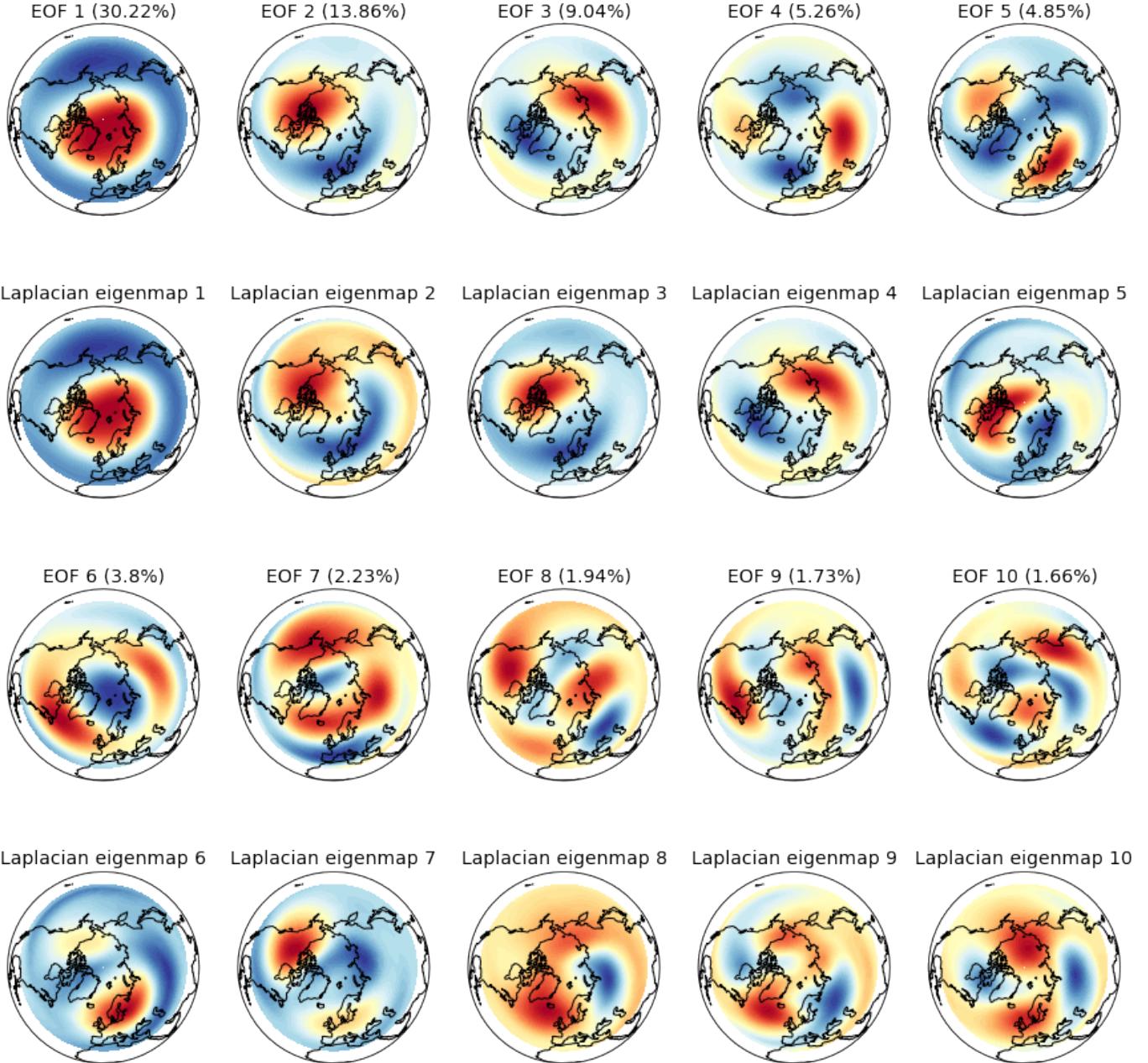


Fig. 9: (RAW) First ten leading EOF ζ_i^R to ζ_{10}^R obtained after linear dimension reduction using singular value decomposition on the raw dataset to project it to a lower dimensional space. The input data is centered but not scaled for each feature before applying SVD. Percentage of explained variance per EOF ζ_i^R between parentheses. Those eigenvectors are compared to the ten leading Laplacian eigenmaps ϕ_1^R to ϕ_{10}^R obtained using the heat-kernel as in Fig. 5 ($n = \%$, $t = \text{mean}$). Laplacian eigenmaps geospatial patterns are reconstructed according to equation 8 in section III-B5. Next ten eigenfunctions in Fig. 10.

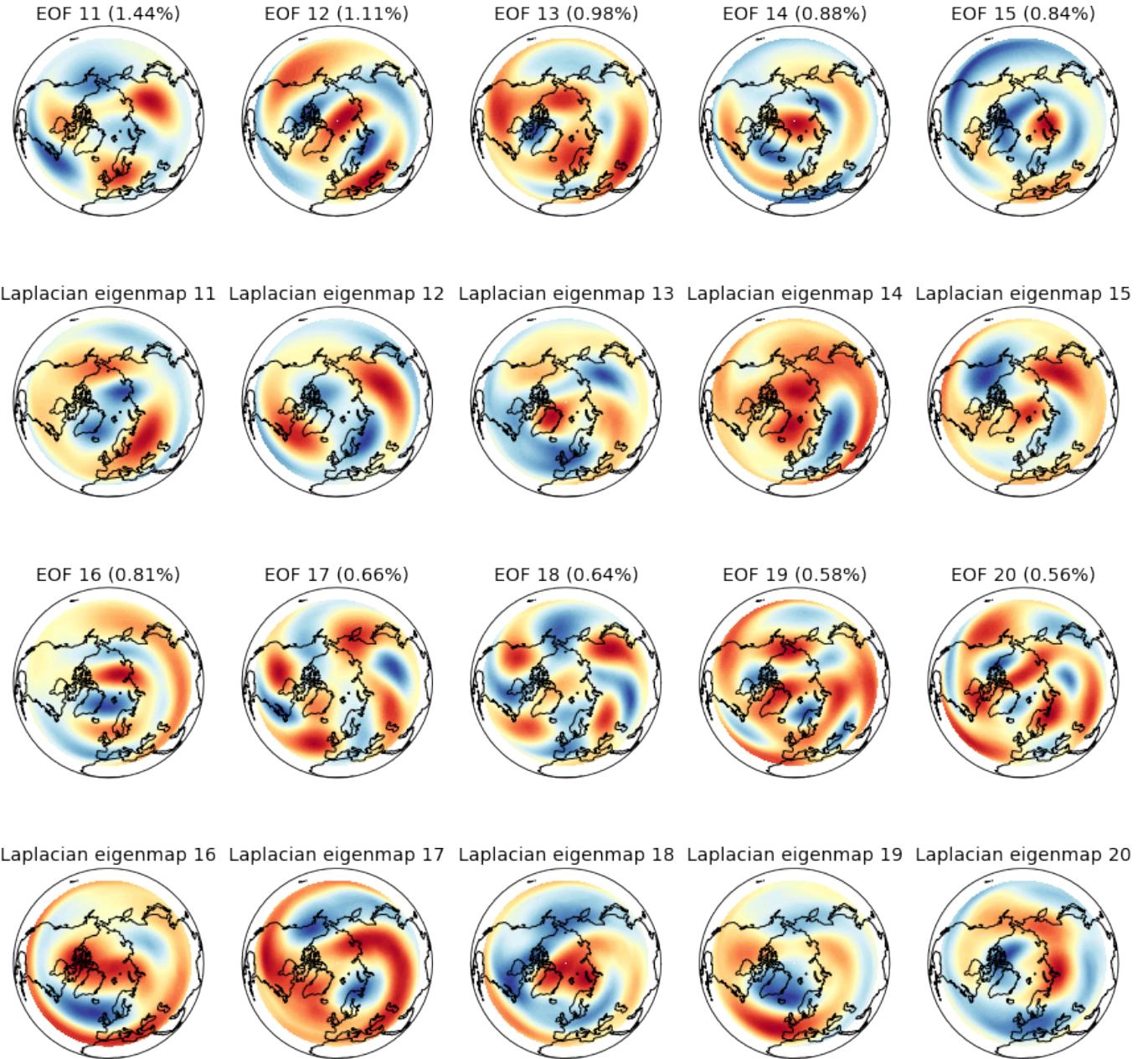


Fig. 10: (RAW) Leading EOF ten to twenty ζ_{11}^R to ζ_{20}^R obtained after linear dimension reduction using singular value decomposition on the raw dataset to project it to a lower dimensional space. The input data is centered but not scaled for each feature before applying SVD. Percentage of explained variance per EOF ζ_i between parentheses. Those eigenvectors are compared to leading Laplacian eigenmaps ϕ_{11}^R to ϕ_{20}^R obtained using the heat-kernel as in Fig. 5 ($n = \%$, $t = \text{mean}$). Laplacian eigenmaps geospatial patterns are reconstructed according to equation 8 in section III-B5.

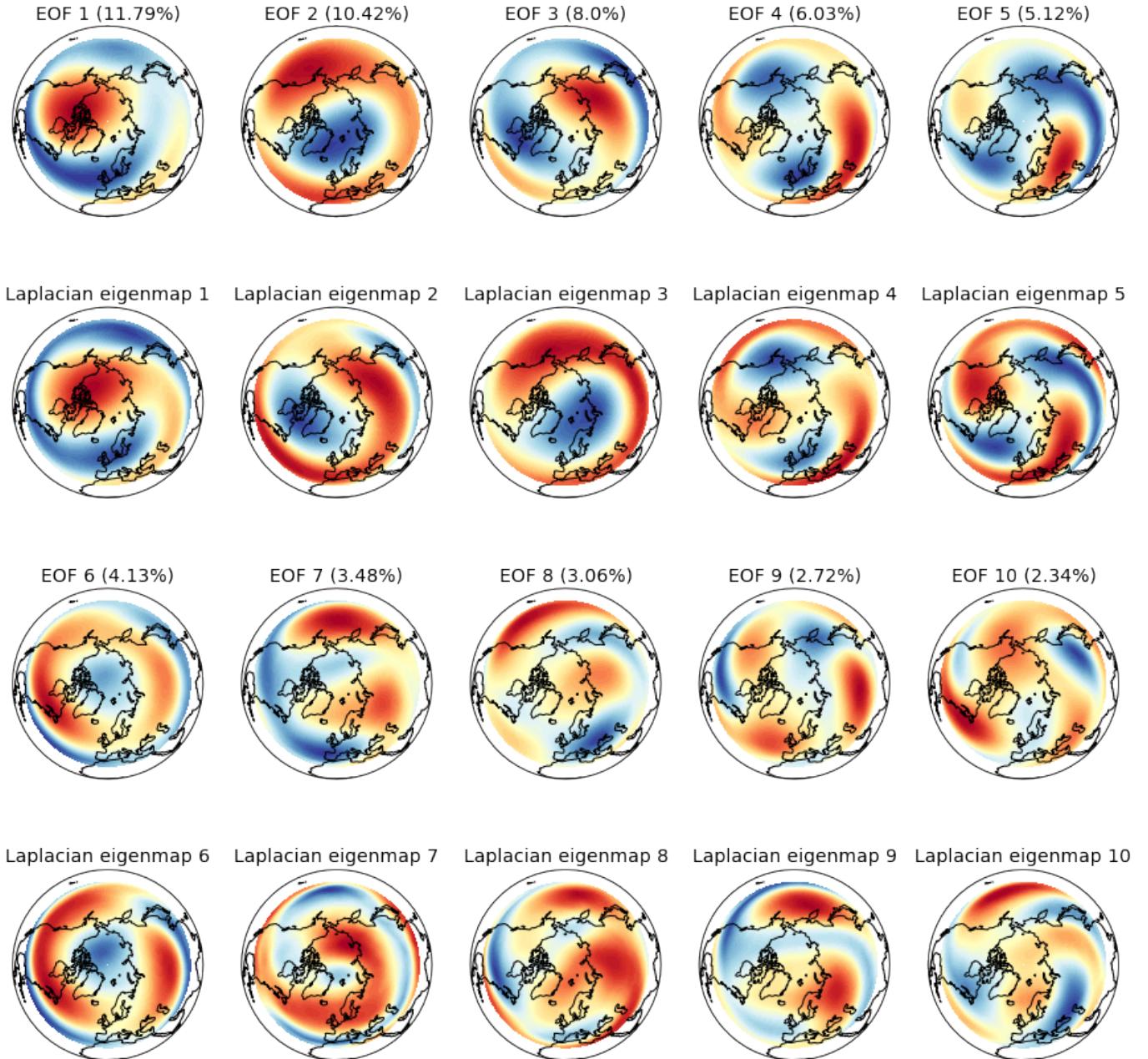


Fig. 11: (ANOMALIES) First ten leading EOF ζ_1^A to ζ_{10}^A obtained after linear dimension reduction using singular value decomposition on the anomalies dataset to project it to a lower dimensional space. The input data is centered but not scaled for each feature before applying SVD. Percentage of explained variance per EOF ζ_i^A between parentheses. Those eigenvectors are compared to the ten leading Laplacian eigenmaps ϕ_1^A to ϕ_{10}^A obtained using the heat-kernel as in Fig. 5 ($n = \%$, $t = \text{mean}$). Laplacian eigenmaps geospatial patterns are reconstructed according to equation 8 in section III-B5. Next ten eigenfunctions in Fig. 12.

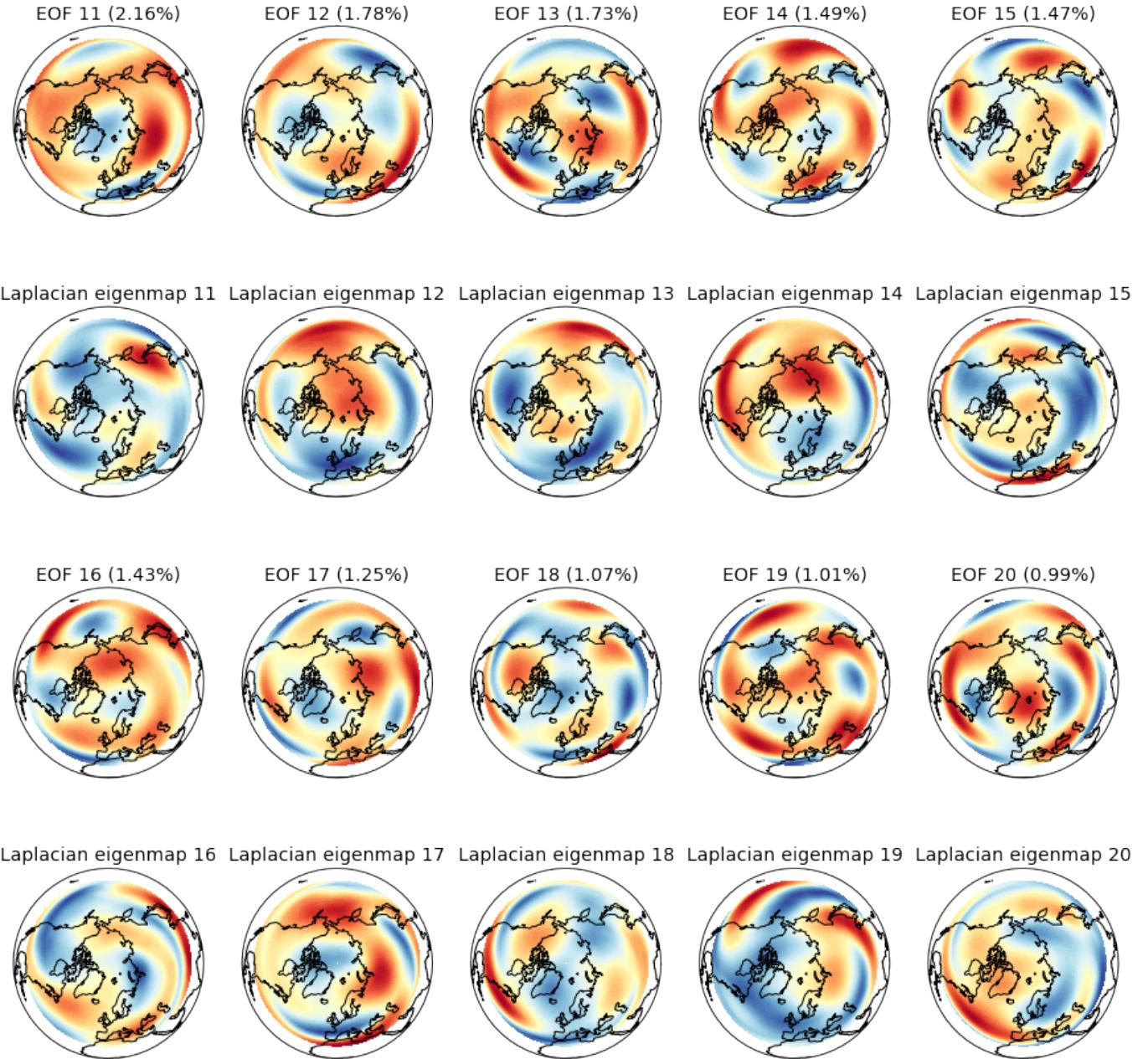


Fig. 12: (ANOMALIES) Leading EOF ten to twenty ζ_{11}^A to ζ_{20}^A obtained after linear dimension reduction using singular value decomposition on the anomalies dataset to project it to a lower dimensional space. The input data is centered but not scaled for each feature before applying SVD. Percentage of explained variance per EOF ζ_i^A between parentheses. Those eigenvectors are compared to leading Laplacian eigenmaps ϕ_{11}^A to ϕ_{20}^A obtained using the heat-kernel as in Fig. 5 ($n = \%$, $t = mean$). Laplacian eigenmaps geospatial patterns are reconstructed according to equation 8 in section III-B5.

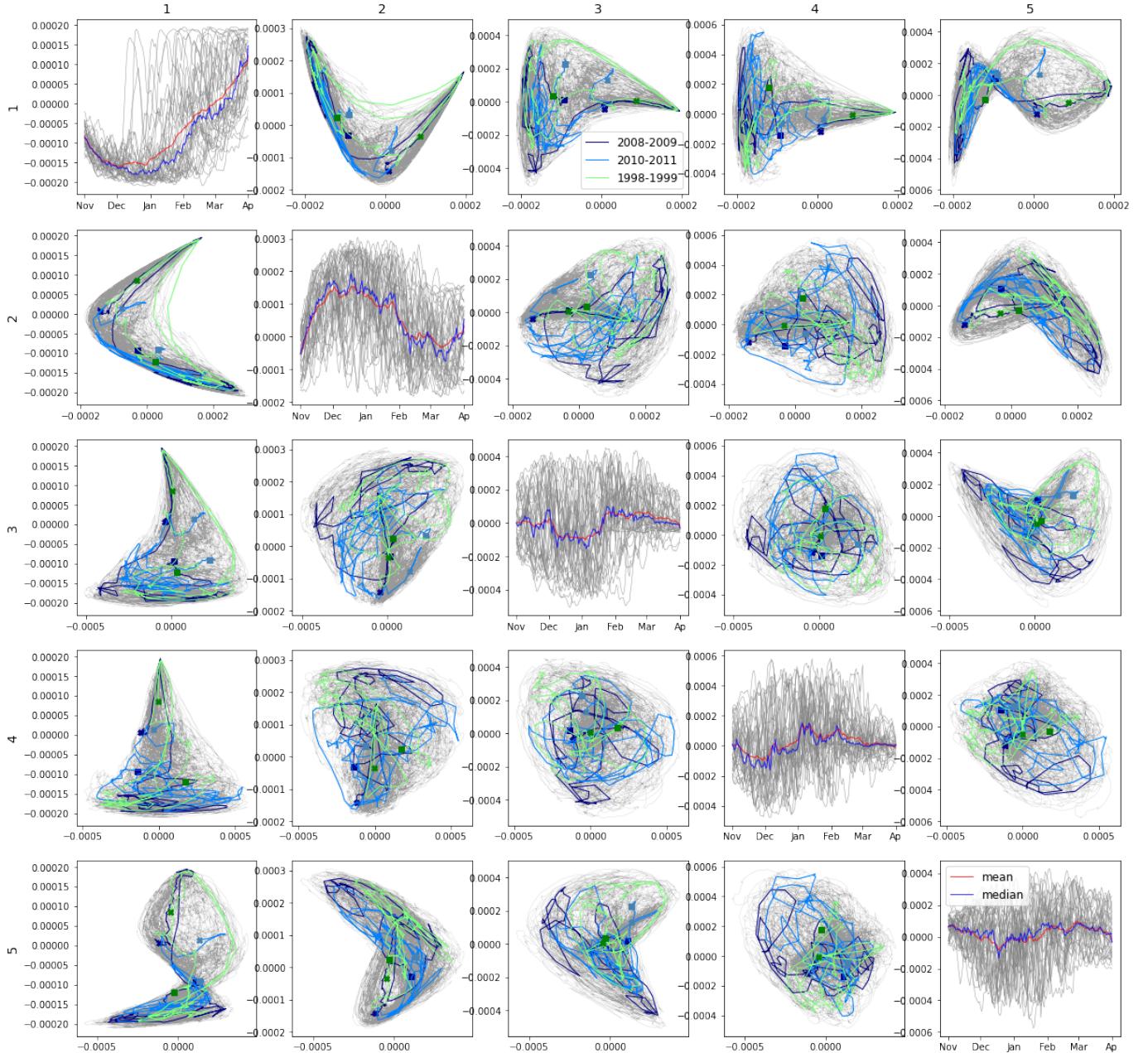


Fig. 13: (RAW) Five leading Laplacian eigenmaps ϕ_1^R to ϕ_5^R obtained using the heat-kernel as in Fig. 5 ($n = \%$, $t = \text{mean}$) on raw data. In the diagonal, eigenvectors for winters from 1979 to 2018 are plotted as overlapping from November to April. The red and blue curve respectively indicate the mean and median over the values for each time-step. Out of the diagonal, eigenvectors ϕ_1^R to ϕ_5^R are plotted against each-other. From left to right eigenvectors ϕ_1^R to ϕ_5^R as seen in Fig. VI.1 plotted in the x-axis. From top to bottom eigenvectors ϕ_1^R to ϕ_5^R plotted in the y-axis. The three highlighted curves indicate trajectories of the corresponding eigenvectors for a time-interval of November to April for 1998-1999 (green), 2008-2009 (dark blue) and 2010-2011 (light blue). The trajectories starts in November in at the square marks and end in April on the cross marks.

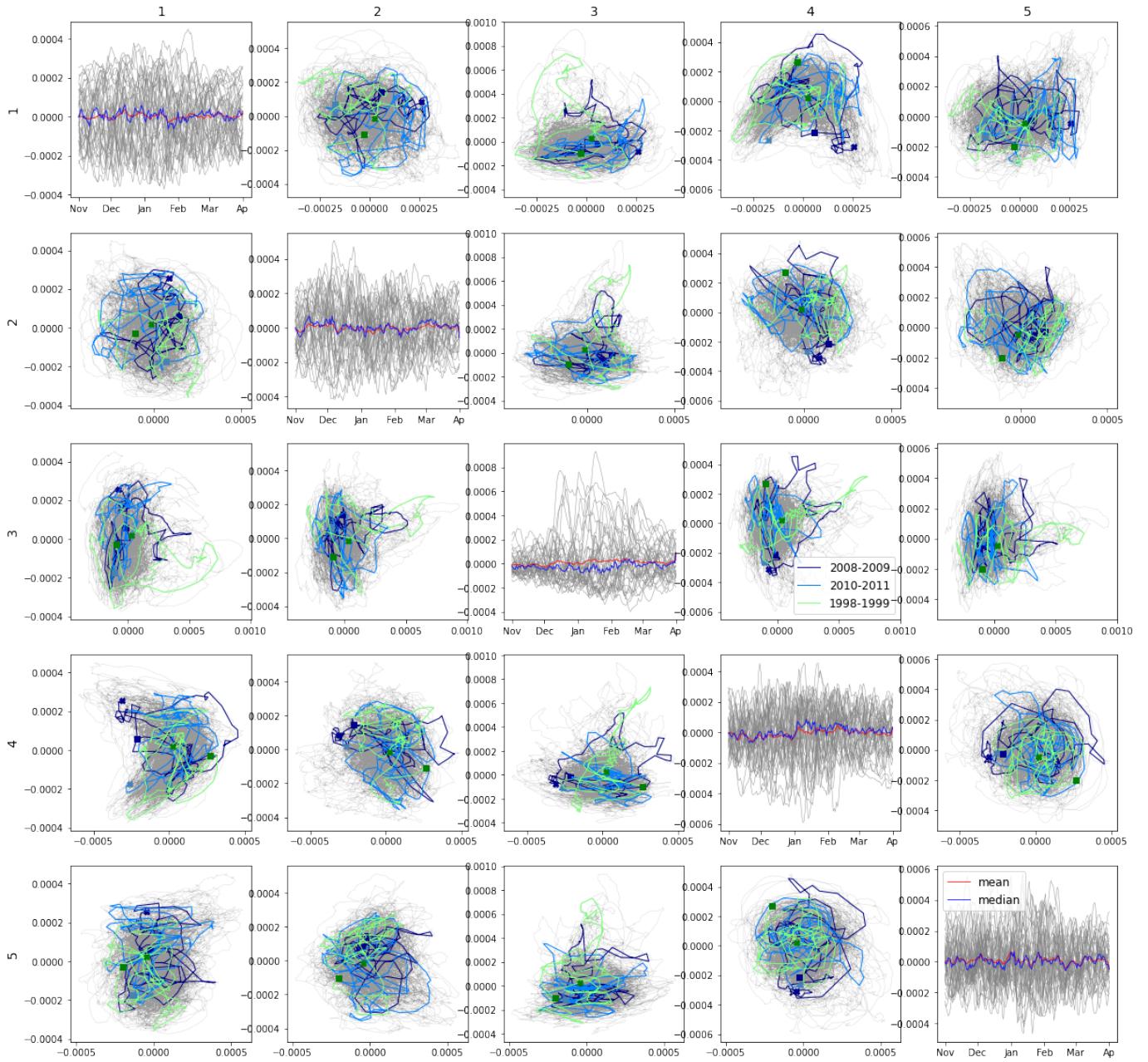
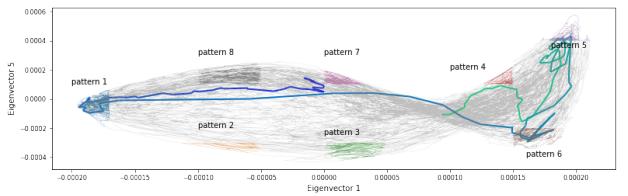
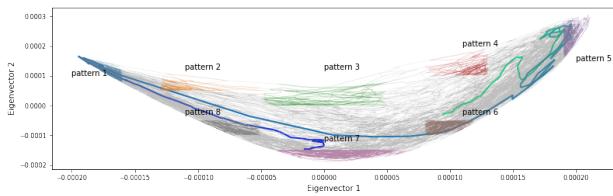
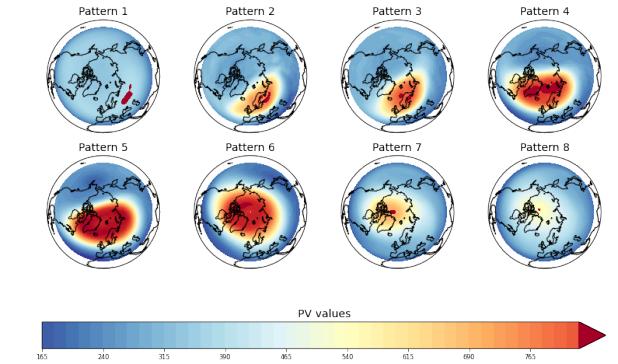
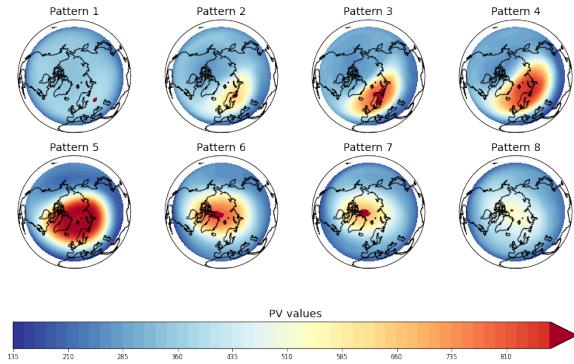
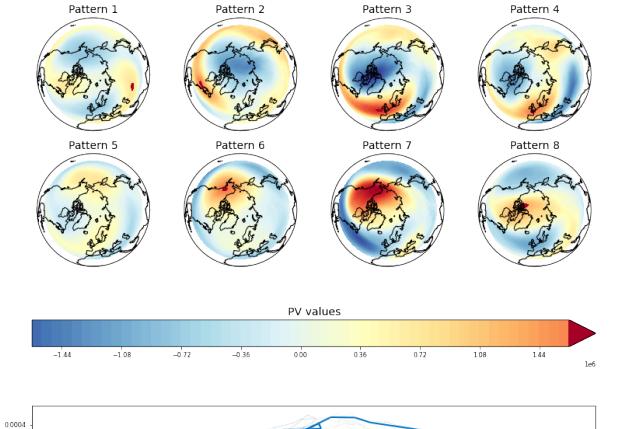
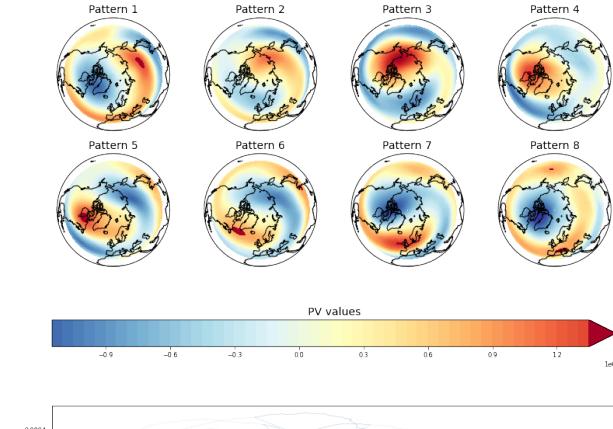


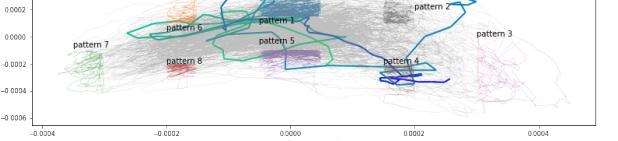
Fig. 14: (ANOMALIES) Five leading Laplacian eigenmaps ϕ_1^A to ϕ_5^A obtained using the heat-kernel as in Fig. 5 ($n = \%$, $t = \text{mean}$) on anomalies data. In the diagonal, eigenvectors for winters from 1979 to 2018 are plotted as overlapping from November to April. The red and blue curve respectively indicate the mean and median over the values for each time-step. Out of the diagonal, eigenvectors ϕ_1^A to ϕ_5^A are plotted against each-other. From left to right eigenvectors ϕ_1^A to ϕ_5^A as seen in Fig. VI.1 plotted in the x-axis. From top to bottom eigenvectors ϕ_1^A to ϕ_5^A plotted in the y-axis. The three highlighted curves indicate trajectories of the corresponding eigenvectors for a time-interval of November to April for 1998-1999 (green), 2008-2009 (dark blue) and 2010-2011 (light blue). The trajectories starts in November in at the square marks and end in April on the cross marks.



(a) (RAW) ϕ_1^R vs ϕ_2^R plotted against each-other (bottom) and average spatiotemporal patterns in the original data space (top) extracted for each region in the bottom plot.



(c) (ANOMALIES) ϕ_1^A vs ϕ_2^A plotted against each-other (bottom) and average spatiotemporal patterns in the original data space (top) extracted for each region in the bottom plot.



(d) (ANOMALIES) ϕ_1^A vs ϕ_4^A plotted against each-other (bottom) and average spatiotemporal patterns in the original data space (top) extracted for each region in the bottom plot.

Fig. 15: Eigenmaps from Fig. VI.1, VI.2 plotted against each-other for raw (a, b) and anomalies data (c, d). The curves in blue on the ϕ_i vs ϕ_j plots indicate a trajectory of the corresponding eigenvectors for a time-interval of November 2008 to April 2009. The trajectory starts in November 2008 in light green and ends in April 2009 in dark blue. For each of these plots, regions were chosen that follow the corresponding approximate yearly trajectories. From these, average spatiotemporal patterns were reconstructed in the original space, as explained in IV-C. Patterns were plotted using the basis Ψ as in equation 1.

REFERENCES

- [1] M. P. Baldwin and T. J. Dunkerton, "Stratospheric harbingers of anomalous weather regimes," *Science*, vol. 294, no. 5542, pp. 581–584, 2001. [Online]. Available: <https://science.sciencemag.org/content/294/5542/581>
- [2] R. de Fondeville, Z. Wu, E. Szekely, W. Ball, G. Obozinski, and D. Domeisen, "Toward enhanced seasonal forecasts: Characterization of the polar vortex variability," Presentation, 2020.
- [3] M. Sigmond, J. F. Scinocca, V. V. Kharin, and T. G. Shepherd, "Enhanced seasonal forecast skill following stratospheric sudden warmings," *Nature Geoscience*, vol. 6, no. 2, pp. 98–102, Feb 2013. [Online]. Available: <https://doi.org/10.1038/ngeo1698>
- [4] B. Kirtman and A. Pirani, "Wcrp position paper on seasonal prediction," *Report from the First WCRP Seasonal Prediction Workshop*, 2008.
- [5] N. R. Council, *Assessment of Intraseasonal to Interannual Climate Prediction and Predictability*. Washington, DC: The National Academies Press, 2010. [Online]. Available: <https://www.nap.edu/catalog/12878/assessment-of-intraseasonal-to-interannual-climate-prediction-and-predictability>
- [6] "Climate Data Guide empirical orthogonal function (eof) analysis and rotated eof analysis," <https://climatedataguide.ucar.edu/climate-data-tools-and-analysis/empirical-orthogonal-function-eof-analysis-and-rotated-eof-analysis>, accessed: 2020-11-06.
- [7] E. Székely, D. Giannakis, and A. Majda, "Extraction and predictability of coherent intraseasonal signals in infrared brightness temperature data," *Climate Dynamics*, vol. 46, 05 2015.
- [8] W.-w. Tung, D. Giannakis, and A. Majda, "Symmetric and antisymmetric convection signals in the madden–julian oscillation. part i: Basic modes in infrared brightness temperature," *Journal of the Atmospheric Sciences*, vol. 71, pp. 3302–3326, 09 2014.
- [9] M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural Computation*, vol. 15, no. 6, pp. 1373–1396, 2003.
- [10] R. de Fondeville, "Laplacian eigenmaps," Lecture, 2020.
- [11] K. Haven, A. Majda, and R. Abramov, "Quantifying predictability through information theory: Small sample estimation in a non-gaussian framework," *J. Comput. Phys.*, vol. 206, no. 1, p. 334–362, Jun. 2005. [Online]. Available: <https://doi.org/10.1016/j.jcp.2004.12.008>
- [12] M. Belkin and P. Niyogi, "Towards a theoretical foundation for laplacian-based manifold methods," *Journal of Computer and System Sciences*, vol. 74, no. 8, pp. 1289 – 1308, learning Theory 2005. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0022000007001274>
- [13] D. Giannakis, A. J. Majda, and I. Horenko, "Information theory, model error, and predictive skill of stochastic models for complex nonlinear systems," *Physica D: Nonlinear Phenomena*, vol. 241, no. 20, pp. 1735 – 1752, 2012. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0167278912001868>
- [14] N. Aubry, W. Lian, and E. Titi, *Preserving Symmetries in the Proper Orthogonal Decomposition*, ser. MSRI (Series). Mathematical Sciences Research Institute, 1991. [Online]. Available: <https://books.google.ch/books?id=5yr3NAAACAAJ>
- [15] D. Giannakis and A. J. Majda, "Limits of predictability in the north pacific sector of a comprehensive climate model," *Geophysical Research Letters*, vol. 39, no. 24, 2012. [Online]. Available: <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2012GL054273>
- [16] ——, "Nonlinear laplacian spectral analysis: capturing intermittent and low-frequency spatiotemporal patterns in high-dimensional data," *Statistical Analysis and Data Mining: The ASA Data Science Journal*, vol. 6, no. 3, pp. 180–194, 2013. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/sam.11171>
- [17] D. Giannakis and A. Majda, *Data-Driven Methods for Dynamical Systems: Quantifying Predictability and Extracting Spatiotemporal Patterns*. Wiley, May 2015, pp. 137–191.

VI. APPENDIX

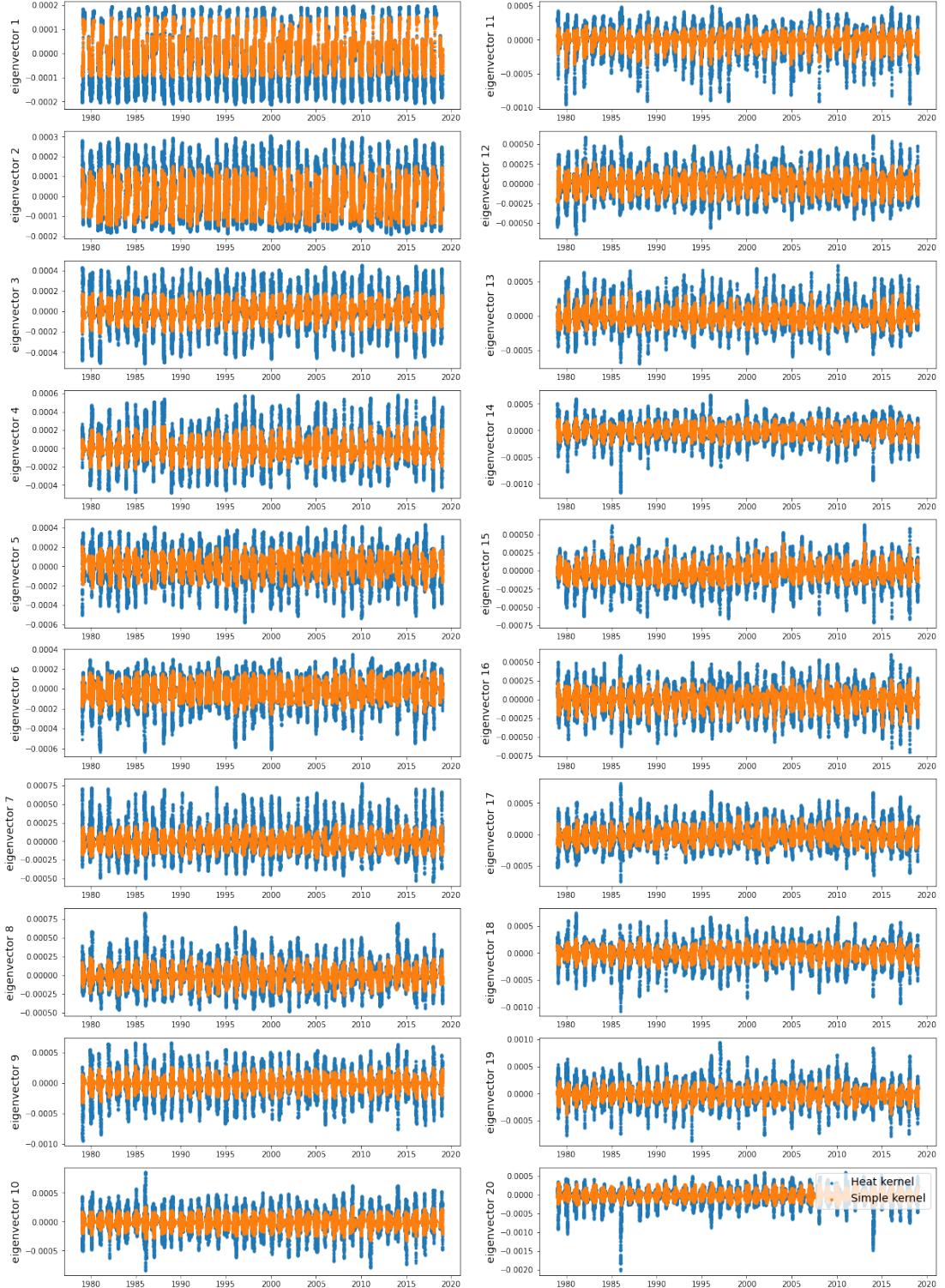


Fig. VI.1: (RAW) First twenty non zero eigenvectors on raw dataset using the heat-kernel from Fig. 5 ($n = 10\%$ and $t = \text{mean}$) (blue) versus first twenty non zero eigenvectors using the simple binary kernel (orange).

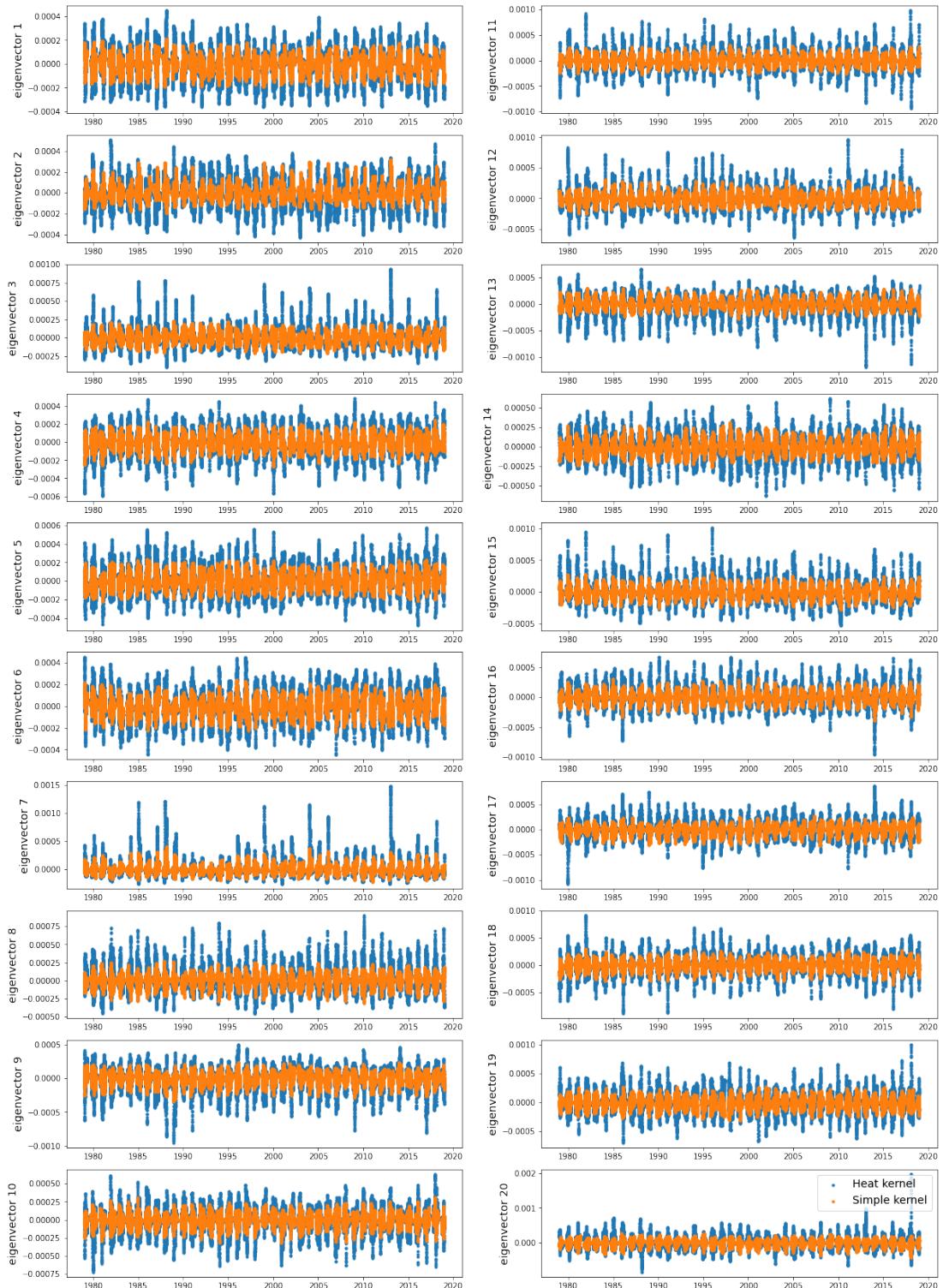


Fig. VI.2: (ANOMALIES) First twenty non zero eigenvectors on anomalies dataset using the heat-kernel from Fig. 5 ($n = 10\%$ and $t = \text{mean}$) (blue) versus first twenty non zero eigenvectors using the simple binary kernel (orange).

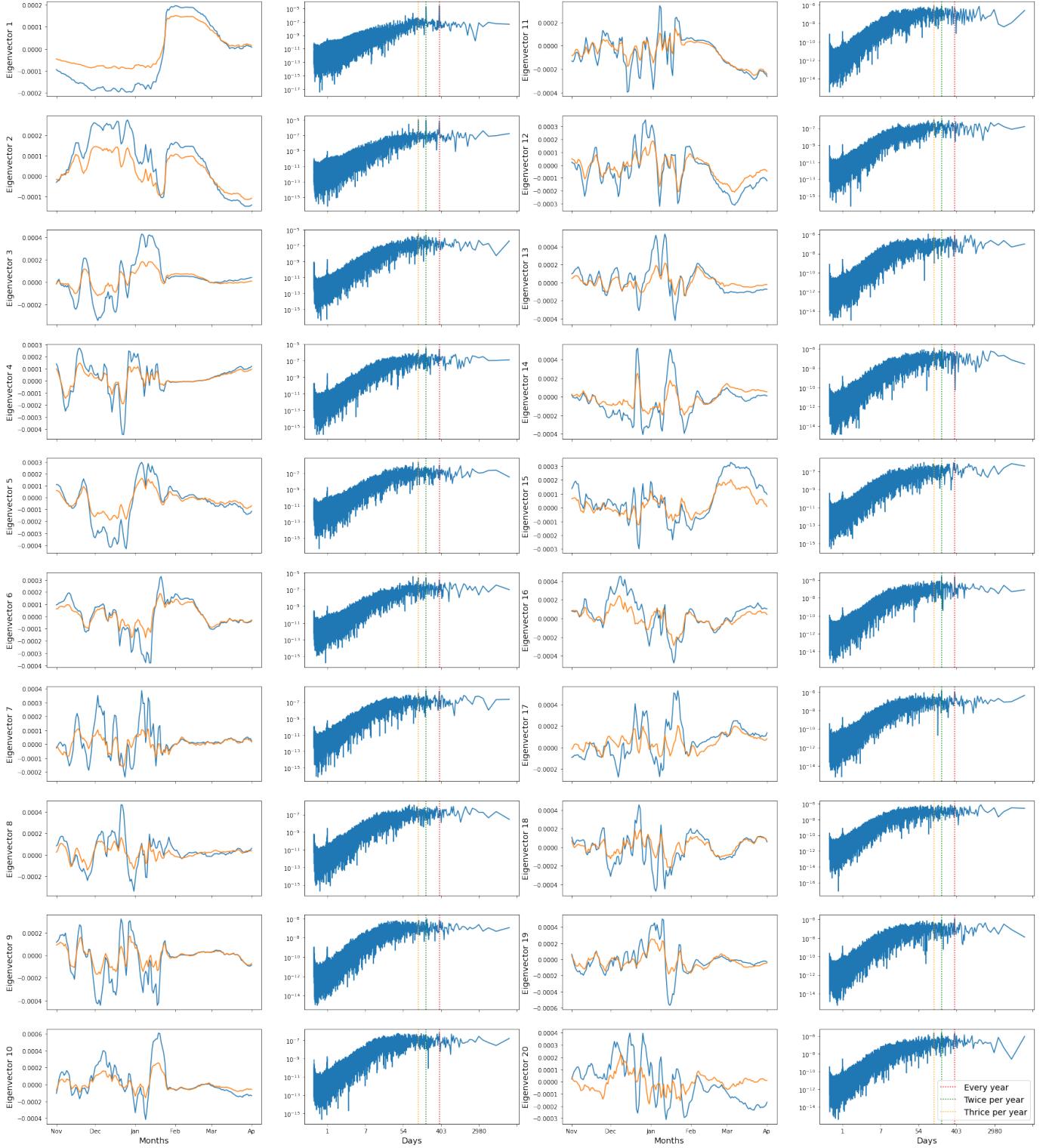


Fig. VI.3: (RAW, NLSA) Leading twenty non-zero NLSA eigenfunctions ϕ_i^R for the winter time interval of November 2008 to April 2009, with their associated power-spectrum, an estimation of the spectral density of the signals. NLSA was performed using a Takens embedding space with $\tau = 7$ days as a temporal extent of the embedding window. Power-spectrums were computed on eigenvectors as in Fig VI.1 but zero-padded in order to complete the data for the missing months. Eigenfunctions were computed on the raw data using the heat-kernel from Fig. 5 ($n = 10\%$ and $t = \text{mean}$) (blue) and the simple binary kernel (orange). Vertical lines in the power spectra indicate frequencies of once, twice and thrice per year respectively in yellow, green and red.

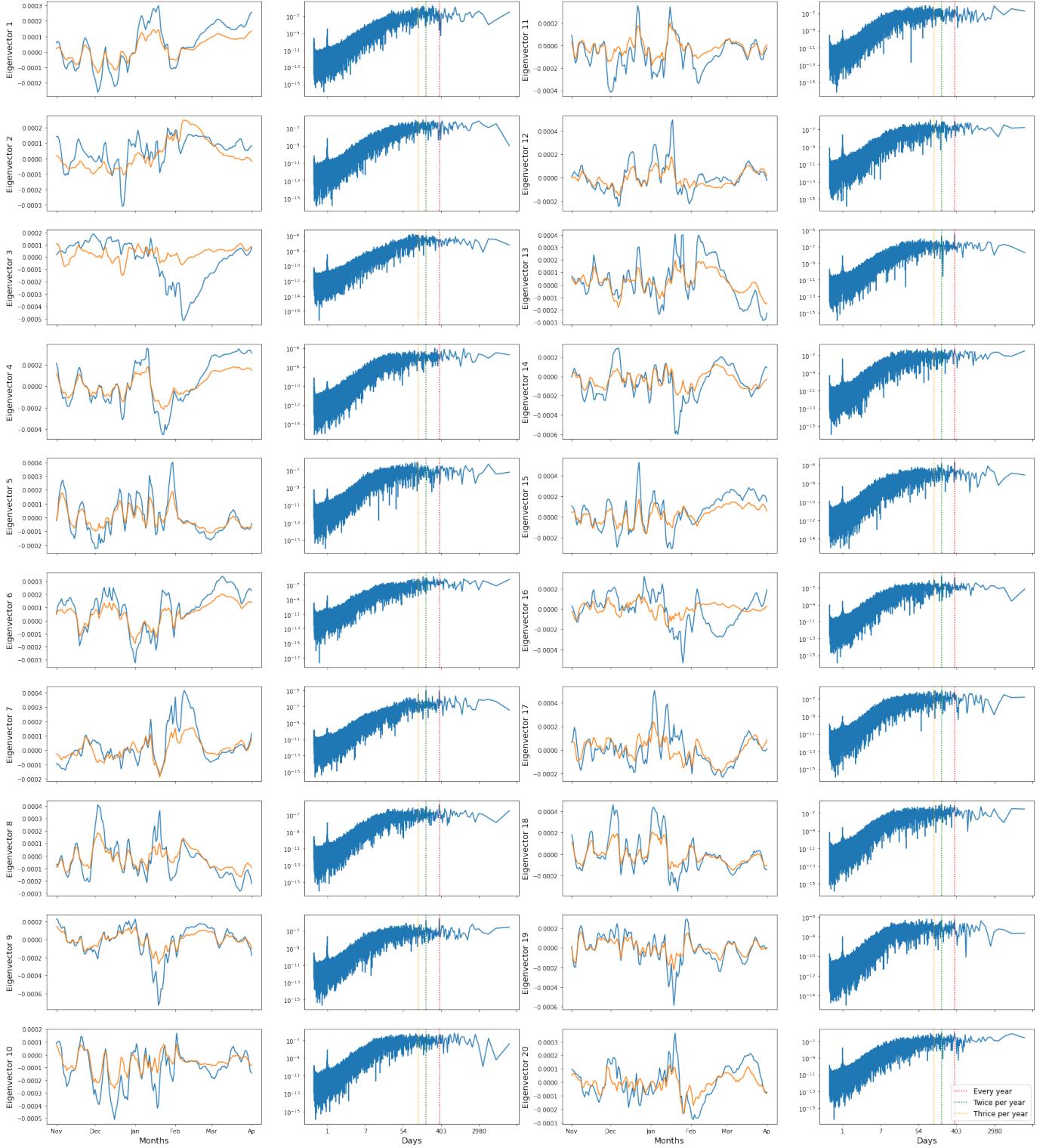


Fig. VI.4: (ANOMALIES, NLSA) Leading twenty non-zero NLSA eigenmaps ϕ_i^A for the winter time interval of November 2008 to April 2009, with their associated power-spectrum, an estimation of the spectral density of the signals. NLSA was performed using a Takens embedding space with $\tau = 7$ days as a temporal extent of the embedding window. Power-spectra were computed on eigenvectors as in Fig VI.2 but zero-padded in order to complete the data for the missing months. Eigenfunctions were computed on the anomalies data using the heat-kernel from Fig. 5 ($n = 10\%$ and $t = \text{mean}$) (blue) and the simple binary kernel (orange). Vertical lines in the power spectra indicate frequencies of once, twice and thrice per year respectively in yellow, green and red.

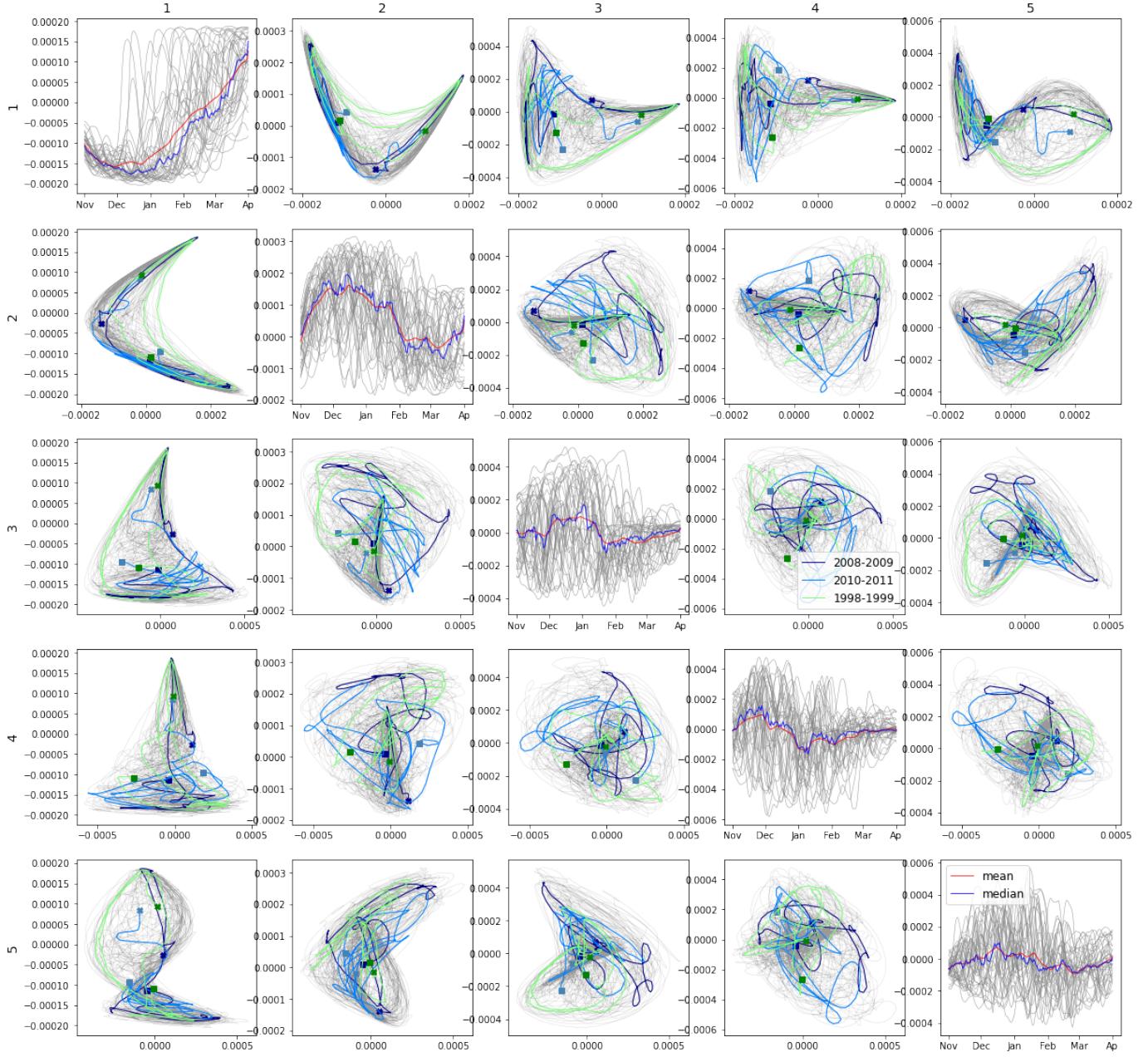


Fig. VI.5: (RAW, NLSA) Five leading NLSA eigenmaps ϕ_1^R to ϕ_5^R obtained using a Takens embedding space with $\tau = 7$ days as a temporal extent of the embedding window and the heat-kernel as in Fig. 5 ($n = \%$, $t = \text{mean}$) on raw data. In the diagonal, eigenvectors for winters from 1979 to 2018 are plotted as overlapping from November to April. The red and blue curve respectively indicate the mean and median over the values for each time-step. Out of the diagonal, eigenvectors ϕ_1^R to ϕ_5^R are plotted against each-other. From left to right eigenvectors ϕ_1^R to ϕ_5^R as seen in Fig. VI.1 plotted in the x-axis. From top to bottom eigenvectors ϕ_1^R to ϕ_5^R plotted in the y-axis. The three highlighted curves indicate trajectories of the corresponding eigenvectors for a time-interval of November to April for 1998-1999 (green), 2008-2009 (dark blue) and 2010-2011 (light blue). The trajectories starts in November in at the square marks and end in April on the cross marks.

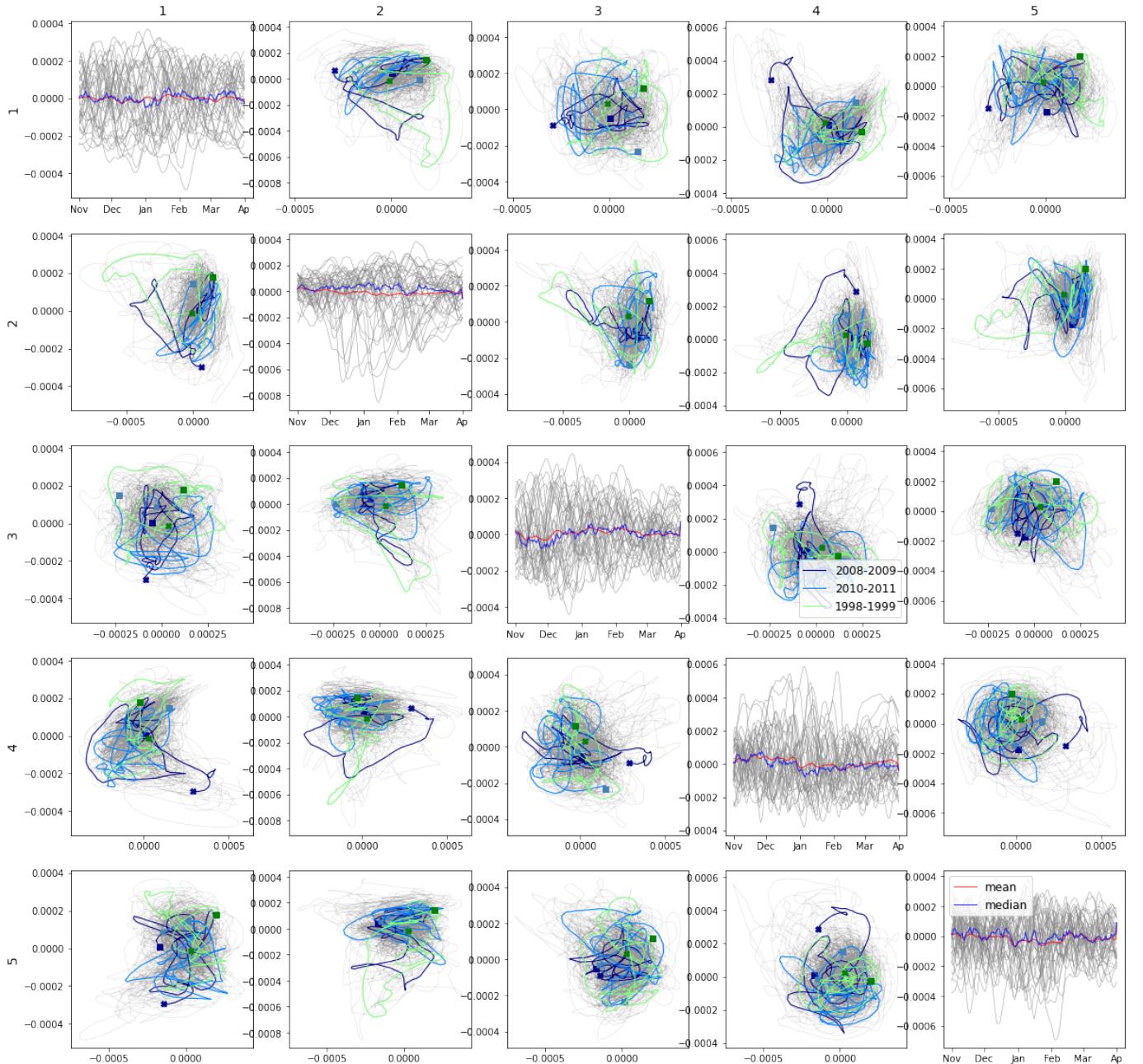


Fig. VI.6: (ANOMALIES, NLSA) Five leading NLSA eigenmaps ϕ_1^A to ϕ_5^A obtained using a Takens embedding space with $\tau = 7$ days as a temporal extent of the embedding window and the heat-kernel as in Fig. 5 ($n = \%$, $t = \text{mean}$) on anomalies data. In the diagonal, eigenvectors for winters from 1979 to 2018 are plotted as overlapping from November to April. The red and blue curve respectively indicate the mean and median over the values for each time-step. Out of the diagonal, eigenvectors ϕ_1^A to ϕ_5^A are plotted against each-other. From left to right eigenvectors ϕ_1^A to ϕ_5^A as seen in Fig. VI.1 plotted in the x-axis. From top to bottom eigenvectors ϕ_1^A to ϕ_5^A plotted in the y-axis. The three highlighted curves indicate trajectories of the corresponding eigenvectors for a time-interval of November to April for 1998-1999 (green), 2008-2009 (dark blue) and 2010-2011 (light blue). The trajectories starts in November in at the square marks and end in April on the cross marks.