

Time Series Clustering using Trend, Seasonal and Autoregressive Components: Patterns of Change of Maximum Temperature in Iberian Peninsula

Arnobio PALACIOS GUTIERREZ^{1,2,✉} and Jose Luis VALENCIA DELFA¹

¹Complutense University of Madrid, Faculty of Statistical Studies,
Madrid, Community of Madrid, Spain

²Technological University of Chocó, Group Valoración y Aprovechamiento de la Biodiversidad,
Quibdó, Chocó-Colombia

✉ arnobiop@ucm.es

Abstract

Time series clustering is an important field of data mining and can be used to identify interesting patterns. This study introduces a new way to obtain clusters of time series by representing them with feature vectors that define the trend, seasonality and noise components of each series, in order to identify areas of the Iberian Peninsula that follow the same pattern of change in their maximum temperature during 1931–2009. Singular spectrum analysis decomposition in a sequential manner is used for dimensionality reduction, which allows the extraction of the trend, seasonality and residual components of each time series corresponding to an area of the Iberian region; then, the feature vectors of the time series are obtained by modelling the extracted components and estimating the parameters. Finally, the series are clustered using a clustering algorithm, and the clusters are defined according to the centroids. The results identified three differentiated zones, allowing to describe how the maximum temperature varied: in the north and central zones, an increase in temperature was noted over time, and in the south, a slight decrease, moreover different seasonal variations were noted according to zones.

Keywords: clustering, maximum temperature time series, singular spectrum analysis, feature vectors of time series, Iberian Peninsula.

1. INTRODUCTION

Climate change is a global problem that has a significant impact on society and ecosystems and is increasingly noticeable. In consequence, studies on climate change have become increasingly more important, especially those relating to temperatures, which have been increasing as stated in the report of the Sixth Assessment of the Intergovernmental Panel on Climate Change (IPCC) that reports that global mean surface temperature (GMST) has increased by 1.1 °C between the 2001–2021 and 1850–1900 periods, after accelerating its rate after the 1970s (IPCC 2021). Moreover, in the very near future, 2020–2050, GMST can warm up as much as 0.25 °C per decade, according to some climate model predictions (Samset et al. 2020; Tebaldi et al. 2021).

Although these numbers may seem low, the changes and effects are really remarkable, as manifested by global warming, prolonged droughts, heat waves and forest fires. In Europe and the Iberian Peninsula (IP) have been experiencing these conditions in the last years (Kuglitsch et al. 2010; Russo et al. 2015; Molina et al. 2020; Calheiros et al. 2021) and according to some climate predictions, is expected these conditions continue for the foreseeable future (IPCC 2014; King and Karoly 2017; Dosio et al. 2018; Vicedo-Cabrera et al. 2018).

Analysis of the temperature changes experienced by the IP, as well as the projections that have arisen around this issue will be a more manageable process if studies that define these changes by zones or sub-regions are include and if analysis that take into account temperature extremes are considered, which tell us about unusual changes (Gebremichael et al. 2022).

A way to define the extreme temperature changes experienced by a geographical area for sub-region is by obtaining time series (TS) clusters of these temperatures defined in points or areas distributed over the geographical area, since; TS clustering is used to identify interesting patterns in TS data sets. There are mainly three categories or approaches to TS clustering (Warren Liao 2005; Rani and Sikka 2012; Aghabozorgi et al. 2015; Ergüner Özkoç 2021), depending on whether they work directly with raw data, indirectly with features or characteristics extracted from the raw data, or indirectly with models built from raw data.

This study is framed within the clustering of TS based on the approach of extracting features from data and proposes a procedure to cluster TS by their trend, seasonality, and main autocorrelations, so that patterns of change in maximum temperature (TMAX) can be identified for zone in the IP during the period 1931–2009. The novelty of our methodology is the use the decomposition of TS using singular spectrum analysis (SSA). In this decomposition process, three components associated with the trend, seasonality and residual of the initial TS are reconstructed, allowing the extraction of the parameters that describe these components. Secondly, the representation of each TS is obtained from a feature vector generated on the basis of the calculated parameters, which allows clustering the TS using unsupervised learning algorithms, such as k-means (Hartigan and Wong 1979), C-medoids (Park and Jun 2009), hierarchical agglomerative (HA) (Lukasová 1979), and Kohonen self-organising maps (SOM) (Kohonen et al. 1996), which are known and representative conventional algorithms that use the Euclidean distance. Finally, in our experiment on a climatic database, after comparing the clusters obtained with the different methods, a hybrid approach that combines HA and k-means, called hkmeans (Lee et al. 2010; Kassambara 2017), is selected as a clustering algorithm to define TS that are similar and follow a pattern. The results made it possible to identify three differentiated zones according to their TMAX level and trend. In addition, was observed that the identified zones show different seasonal variations.

The remainder of this document is organised as follows. Section 2 describes the TS decomposition method using SSA in a sequential manner. Section 3 proposes the new method for defining the trend, seasonality and autoregression patterns of TS. Section 4 presents the results of the method. Section 5 presents the main conclusions of the paper.

2. SEQUENTIAL SSA DECOMPOSITION METHOD

This technique is based on the singular value decomposition (SVD) of a specific matrix obtained from a TS and aims to decompose an original TS into a sum of a small number of interpretable components, such as the trend that is smooth and slowly varying, oscillatory components that are periodic or pure quasiperiodic or amplitude-modulated, and noise without any pattern or structure (Golyandina et al. 2001; Golyandina and Korobeynikov 2014).

In the following, the SSA method is presented formally.

Input: $\mathbb{T} = (\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_N)$ the initial TS, which is one-dimensional N -order TS.

Result: A decomposition of \mathbb{T} into a sum of identifiable components $\mathbb{T} = \tilde{\mathbb{T}}_1 + \tilde{\mathbb{T}}_2 + \dots + \tilde{\mathbb{T}}_m$.

Step 1: Embedding. The so-called “trajectory matrix” is obtained as $\mathbf{X} = \mathcal{J}(\mathbb{T})$, where \mathcal{J} is a linear map that transforms the TS \mathbb{T} into a matrix of order $L \times K$, where L is an integer that is called the “window length”, $1 < L < N$, and $K = N - L + 1$.

The set of all possible path matrices can be denoted as $\mathcal{M}_{L,K}^{(\mathcal{H})}$. \mathcal{H} / denotes Hankel matrices, where all elements along the diagonal are equal. If N and L are fixed, then there is a biunivocal correspondence between the path matrices and the TS.

The trajectory matrix \mathbf{X} constructed from lagged vectors generated from the TS \mathbb{T} can be represented as:

$$\mathcal{J}(\mathbb{T}) = \begin{pmatrix} \mathbf{t}_1 & \mathbf{t}_2 & \mathbf{t}_3 & \dots & \mathbf{t}_K \\ \mathbf{t}_2 & \mathbf{t}_3 & \mathbf{t}_4 & \dots & \mathbf{t}_{K+1} \\ \mathbf{t}_3 & \mathbf{t}_4 & \mathbf{t}_5 & \dots & \mathbf{t}_{K+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{t}_L & \mathbf{t}_{L+1} & \mathbf{t}_{L+2} & \dots & \mathbf{t}_N \end{pmatrix}. \quad (1)$$

Step 2: Decomposition of \mathbf{X} into a sum of the matrices of rank 1. The result obtained in this step is the decomposition:

$$\mathbf{X} = \sum_i \mathbf{X}_i, \quad \mathbf{X}_i = \sigma_i U_i V_i^T, \quad (2)$$

where $U_i \in R^L$ and $V_i \in R^K$ are vectors such that $\|U_i\| = 1$ and $\|V_i\| = 1$ for all i and σ_i denotes nonnegative numbers.

If such a decomposition is performed by conventional SVD, the corresponding SSA method is “Basic SSA”, and the singular value decomposition of the matrix \mathbf{X} is calculated via the eigenvalues and eigenvectors of the matrix $S = \mathbf{X}\mathbf{X}^T$ of size $L \times L$. Here, $\lambda_1, \dots, \lambda_L$ denotes the eigenvalues of the matrix S taken in decreasing order of magnitude ($\lambda_1 \geq \dots \geq \lambda_L \geq 0$) and U_1, \dots, U_L denotes the orthonormal system of the eigenvectors of the matrix S corresponding to these eigenvalues. If $d = \text{rank}(\mathbf{X}) = \max\{i, \text{ such that } \lambda_i \geq 0\}$ and $V_i = \mathbf{X}^T U_i / \sqrt{\lambda_i}$, ($i = 1, \dots, d$) are factor vectors, then $\mathbf{X}_i = \sqrt{\lambda_i} U_i V_i^T$ are matrices of rank 1, so they are elementary matrices. Thus, the SVD of the trajectory matrix can be written as:

$$\mathbf{X} = \mathbf{X}_1 + \dots + \mathbf{X}_d. \quad (3)$$

The collection $(\sqrt{\lambda_i}, U_i, V_i^T)$ is called an SVD eigenvector of order i and consists of the singular value $= \sqrt{\lambda_i}$, an eigenvector U_i (the left singular vector) and a factor vector V_i (the right singular vector).

Step 3: Grouping. The input of this step is expansion (2) and a specification of how to cluster the components of Eq. (2). The index set $I = \{1, 2, \dots, d\}$ must be segmented into m disjoint

subsets. I_1, I_2, \dots, I_m . Let $I = \{i_1, i_2, \dots, i_p\} \subset \{1, 2, \dots, d\}$ be a subset of indices; then, the resulting matrix \mathbf{X}_I corresponding to the group I is defined as:

$$\mathbf{X}_I = \mathbf{X}_{i_1} + \mathbf{X}_{i_2} + \dots + \mathbf{X}_{i_p} . \quad (4)$$

Thus, if a partition is specified in m disjoint subsets of the index set $\{1, 2, \dots, d\}$, then, by expansion (2), the result of the grouping step leads to the following decomposition:

$$\mathbf{X} = \mathbf{X}_{I_1} + \mathbf{X}_{I_2} + \dots + \mathbf{X}_{I_m} . \quad (5)$$

The above procedure for choosing the sets I_1, I_2, \dots, I_m is called the ‘‘eigentriple grouping’’ procedure. The grouping of expansion (2), where $I_k = \{k\}$, is called ‘‘elementary’’.

Step 4: Reconstruction. In this step, each matrix \mathbf{X}_{I_k} from lumped decomposition (5) is transferred into the form of the input object \mathbb{T} , which is a TS of length N . Such a transformation is optimally performed as follows: Let $\mathbf{Y} \in R^{L \times K}$ be a matrix with elements y_{ij} , $1 \leq i \leq L, 1 \leq j \leq K$; we look for an object $\tilde{\mathbb{Y}} \in \mathcal{M}$ that provides the minimum of $\|\mathbf{Y} - \mathcal{J}(\tilde{\mathbb{Y}})\|_F$, where $\|\mathbf{Z}\|_F = (\sum_{ij} |z_{ij}|^2)^{1/2}$ is the Frobenius norm of $\mathbf{Z} = [z_{ij}] \in R^{L \times K}$.

Let $\Pi_{\mathcal{H}} : R^{L \times K} \rightarrow \mathcal{M}_{L,K}^{(\mathcal{H})}$ be the orthogonal projection of $R^{L \times K}$ onto $\mathcal{M}_{L,K}^{(\mathcal{H})}$ in the Frobenius norm. Then, $\tilde{\mathbb{Y}} = \mathcal{J}^{-1} \circ \Pi_{\mathcal{H}}(\mathbf{Y})$. The projection $\Pi_{\mathcal{H}}$ is simply the average of the entries corresponding to a given element of an object (Golyandina et al. 2018; Section 1.1.2.6). In Basic SSA, the composite mapping $\mathcal{J}^{-1} \circ \Pi_{\mathcal{H}}$ uses the long average of antidiagonals so that $\tilde{y}_k = \sum_{(i,j) \in \mathcal{A}_k} (\mathbf{Y}_{ij}) / |\mathcal{A}_k|$, where $\mathcal{A}_k = \{(i, j) : i + j = k + 1, 1 \leq i \leq L, 1 \leq j \leq K\}$.

If $\tilde{\mathbf{X}}_k = \mathbf{X}_{I_k}$ are the reconstructed matrices, $\tilde{\mathbf{X}}_k = \Pi_{\mathcal{H}} \tilde{\mathbf{X}}_k$ are their corresponding path matrices, and $\tilde{\mathbb{T}}_k = \mathcal{J}^{-1}(\tilde{\mathbf{X}}_k)$ are the reconstructed objects. Then, the resulting decomposition of the input object \mathbb{T} is:

$$\mathbb{T} = \tilde{\mathbb{T}}_1 + \tilde{\mathbb{T}}_2 + \dots + \tilde{\mathbb{T}}_m . \quad (6)$$

If the grouping is elementary, the reconstructed objects $\tilde{\mathbb{T}}_k$ in Eq. (6) are called ‘‘elementary components’’.

The SSA parameters, i.e., length of the window L and the way in which \mathbf{X}_{I_k} matrices are grouped, are very important for the outcome of the decomposition and depend on properties of the initial TS and the objective of the analysis, check Golyandina et al. (2001) for more details.

SSA can be performed sequentially, which is recommended when the TS structure is complex (Golyandina et al. 2012). **Sequential SSA** consists of two stages: the first stage performs the extraction of the TS trend with a small L , and in the second stage, the periodic components of the residue are detected and extracted with $L \sim N / 2$.

3. TREND, SEASONALITY AND AUTOREGRESSION SSA-BASED TS PATTERN IDENTIFICATION

The algorithm for the identification of trend, seasonality and autoregression patterns in TS proposed in this study can be summarised in the following steps:

- 1) Perform sequential SSA to extract the 3 components of the initial TS associated with the trend, seasonality and residual part.
- 2) Model the extracted TS in such a way that their associated characteristics can be extracted:
 - a. Trend component: from $\mathbb{T}_{\text{trend}} = \mu - \beta t$, estimate μ and β ;

- b. Seasonal component: from $\mathbb{T}_{\text{seasonal}} = c_1 \sin(2\pi t/T) + c_2 \cos(2\pi t/T)$, where T is the period, estimate c_1 and c_2 ;
- c. Residual part: from $\mathbb{T}_{\text{residual}}$, obtain an $AR(T)$ and calculate the autocorrelations $\varphi_1, \varphi_2, \varphi_3$, and φ_s , where $s = T$, the period.

A feature vector is constructed for each initial TS by considering the estimated parameters.

- 3) Use a conventional algorithm to obtain a similar TS.
- 4) By averaging the initial TS of each group, the representative patterns of each group are obtained based on the defined characteristics.

4. RESULTS AND ANALYSIS

A set of 1776 points from a grid of $25 \times 25 \text{ km}^2$ elaborated through spatial interpolation kriging by the “Servicio de Desarrollos Climatológicos” of the Meteorological Spanish State Agency was used. This grid includes points distributed in Spain, Portugal and the closest areas of the Atlantic Ocean and the Mediterranean Sea, for each point a monthly TS of TMAX from January 1931 to 2009 is considered.

After comparing the four clustering algorithms, k-means, k-medoids, HA and SOM, on the dataset, given the superiority of HA and K-Means, a hybrid method, called hierarchical k-means (hkmeans), was selected.

Here, hkmeans is applied to the set formed by the feature vectors of the TS. The results are shown in Fig. 1.

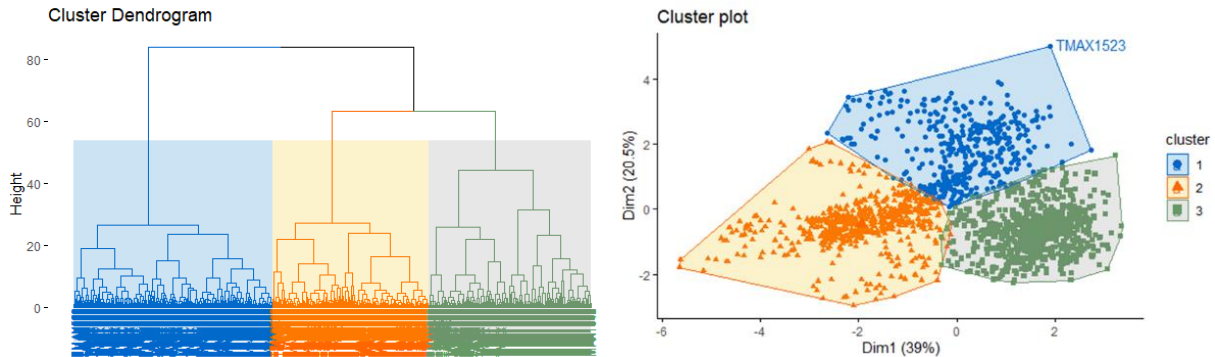


Fig. 1. Result of the hkmeans cluster.

Figure 2 illustrates the distribution (the grid is composed of longitudes and latitudes in UTM coordinates) of the points in Spain according to the clusters obtained.

Clearly, three clusters can be observed; zone 1 situated in the north of the IP, where the areas with the lowest TMAX and with a higher proportion of increase compared to the other areas are found, zone 2, more to the south, where the areas with the highest TMAX are located, but showing a slight decline over the period, and zone 3, with areas more towards the centre and with intermediate TMAX, which also show an increase over time. In addition, was observed that the identified zones show different seasonal variations.

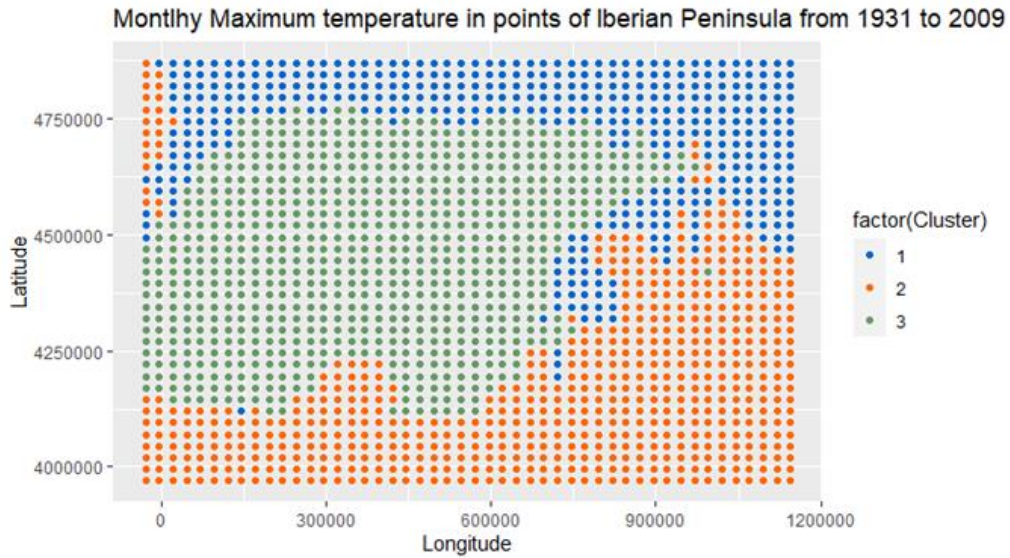


Fig. 2. Distribution of points in Spain according to geographical location and clusters.

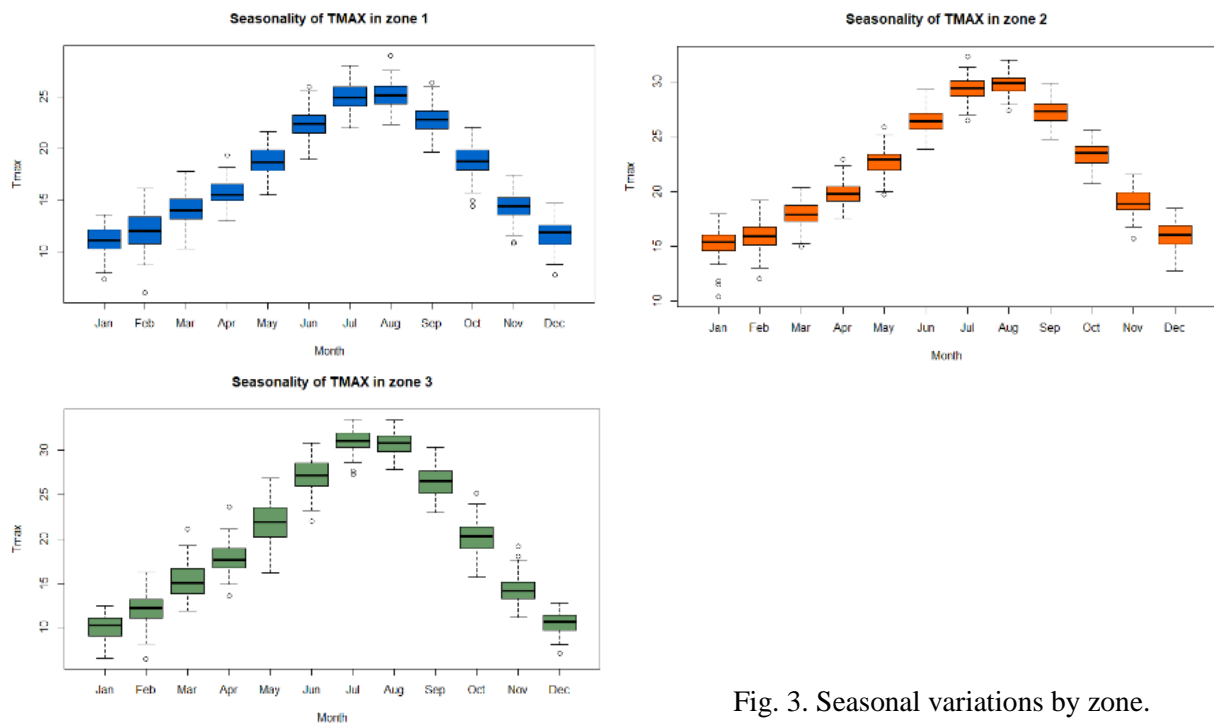


Fig. 3. Seasonal variations by zone.

Figure 3 shows the seasonal variations by zone. It can be noted that, zone 1 shows its largest variations in winter months and, zone 3, in spring and autumn months. Zone 2 does not show marked differences in its monthly variation.

5. CONCLUSIONS

In this paper, we present a new method for clustering TS by taking into account their trend, seasonality, and residual components. The procedure allowed to describe how the maximum temperature varied in the Iberian Peninsula during 1931–2009 through three zones defined according to their trend and monthly variation. The north of the Iberian Peninsula, where the

areas with the lowest maximum temperatures are found, experienced a 0.2034°C increase in its maximum temperature per decade between 1931 and 2009, the south, where the areas with the highest maximum temperatures are located, only showed a slight decline, and the central zone, showed an increase of 0.135°C per decade.

Acknowledgments. We express our gratitude to “Servicio de Desarrollos Climatológicos” of the Meteorological Spanish State Agency for providing the data used in the study, and the Colombian Ministry of Science for the support in the doctoral formation of Arnobio Palacios.

References

- Aghabozorgi, S., A. Seyed Shirkorshidi, and T. Ying Wah (2015), Time-series clustering – A decade review, *Inform. Syst.* **53**, 16–38, DOI: 10.1016/j.is.2015.04.007.
- Calheiros, T., M.G. Pereira, and J.P. Nunes (2021), Assessing impacts of future climate change on extreme fire weather and pyro-regions in Iberian Peninsula, *Sci. Total Environ.* **754**, 142233, DOI: 10.1016/j.scitotenv.2020.142233.
- Dosio, A., L. Mentaschi, E.M. Fischer, and K. Wyser (2018), Extreme heat waves under 1.5°C and 2°C global warming, *Environ. Res. Lett.* **13**, 5, 054006, DOI: 10.1088/1748-9326/aab827.
- Ergüner Özkoç, E. (2021), Clustering of time-series data. **In:** D. Birant (ed.), *Data Mining – Methods, Applications and Systems*, IntechOpen, London, 87–106, DOI: 10.5772/intechopen.84490.
- Gebremichael, H.B., G.A. Raba, K.T. Beketie, G.L. Feyisa, and T. Siyoum (2022), Changes in daily rainfall and temperature extremes of upper Awash Basin, Ethiopia, *Sci. Afr.* **16**, e01173, DOI: 10.1016/j.sciaf.2022.e01173.
- Golyandina, N., and A. Korobeynikov (2014), Basic Singular Spectrum Analysis and forecasting with R, *Comput. Stat. Data An.* **71**, 934–954, DOI: 10.1016/j.csda.2013.04.009.
- Golyandina, N., A. Pepelyshev, and A. Steland (2012), New approaches to nonparametric density estimation and selection of smoothing parameters, *Comput. Stat. Data An.* **56**, 7, 2206–2218, DOI: 10.1016/j.csda.2011.12.019.
- Golyandina, N., A. Korobeynikov, and A. Zhigljavsky (2018), *Singular Spectrum Analysis with R*, Springer-Verlag, Berlin, DOI: 10.1007/978-3-662-57380-8.
- Golyandina, N., V. Nekrutkin, and A.A. Zhigljavsky (2001), *Analysis of Time Series Structure: SSA and Related Techniques*, Monographs on Statistics and Applied Probability, Vol. 90, Chapman & Hall/CRC, Boca Raton, 320 pp.
- Hartigan, J.A., and M.A. Wong (1979), Algorithm AS 136: A K-Means Clustering Algorithm, *J. Roy. Stat. Soc. C: Appl. Statist.* **28**, 1, 100–108, DOI: 10.2307/2346830.
- IPCC (2014), *Climate Change 2014: Synthesis Report. Contribution of Working Groups I, II and III to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, IPCC, Geneva, Switzerland, 151 pp.
- IPCC (2021), *Climate Change 2021: The Physical Science Basis. Contribution of Working Group I to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change* (V. Masson-Delmotte et al. (eds.)), Intergovernmental Panel on Climate Change, available from: https://www.ipcc.ch/report/ar6/wg1/downloads/report/IPCC_AR6_WGI_FrontMatter.pdf www.ipcc.ch.
- Kassambara, A. (2017), *Practical Guide To Cluster Analysis in R. Unsupervised Machine Learning, Multivariate Analysis*, Vol. 1, STHDA, 187 pp.
- King, A.D., and D.J. Karoly (2017), Climate extremes in Europe at 1.5 and 2 degrees of global warming, *Environ. Res. Lett.* **12**, 11, 114031, DOI: 10.1088/1748-9326/aa8e2c.

- Kohonen, T., E. Oja, O. Simula, A. Visa, and J. Kangas (1996), Engineering applications of the self-organizing map, *Proc. IEEE* **84**, 10, 1358–1384, DOI: 10.1109/5.537105.
- Kuglitsch, F.G., A. Toreti, E. Xoplaki, P.M. Della-Marta, C.S. Zerefos, M. Türkeş, and J. Luterbacher (2010), Heat wave changes in the eastern Mediterranean since 1960, *Geophys. Res. Lett.* **37**, 4, L04802, DOI: 10.1029/2009GL041841.
- Lee, A.J.T., M.C. Lin, R.T. Kao, and K.T. Chen (2010), An effective clustering approach to stock market prediction, *PACIS 2010 Proceedings* **54**, 345–354, available from: <https://aisel.aisnet.org/pacis2010/54>.
- Lukasová, A. (1979), Hierarchical agglomerative clustering procedure, *Pattern Recogn.* **11**, 5–6, 365–381, DOI: 10.1016/0031-3203(79)90049-9.
- Molina, M.O., E. Sánchez, and C. Gutiérrez (2020), Future heat waves over the Mediterranean from an Euro-CORDEX regional climate model ensemble, *Sci. Rep.* **10**, 1, 8801, DOI: 10.1038/s41598-020-65663-0.
- Park, H.S., and C.H. Jun (2009), A simple and fast algorithm for K-medoids clustering, *Expert Syst. Appl.* **36**, 2, part 2, 3336–3341, DOI: 10.1016/J.ESWA.2008.01.039.
- Rani, S., and G. Sikka (2012), Recent techniques of clustering of time series data: A survey, *Int. J. Comput. Appl.* **52**, 15, 1–9, DOI: 10.5120/8282-1278.
- Russo, S., J. Sillmann, and E.M. Fischer (2015), Top ten European heatwaves since 1950 and their occurrence in the coming decades, *Environ. Res. Lett.* **10**, 12, 124003. DOI: 10.1088/1748-9326/10/12/124003.
- Samset, B.H., J.S. Fuglestedt, and M.T. Lund (2020), Delayed emergence of a global temperature response after emission mitigation, *Nat. Commun.* **11**, 1, 3261, DOI: 10.1038/s41467-020-17001-1.
- Tebaldi, C., K. Debeire, V. Eyring, E. Fischer, J. Fyfe, P. Friedlingstein, R. Knutti, J. Lowe, B. O'Neill, B. Sanderson, D. van Vuuren, K. Riahi, M. Meinshausen, Z. Nicholls, K.B. Tokarska, G. Hurtt, E. Kriegler, J.-F. Lamarque, G. Meehl, et al. (2021), Climate model projections from the Scenario Model Intercomparison Project (ScenarioMIP) of CMIP6, *Earth Syst. Dynam.* **12**, 1, 253–293, DOI: 10.5194/esd-12-253-2021.
- Vicedo-Cabrera, A.M., Y. Guo, F. Sera, V. Huber, C.-F. Schleussner, D. Mitchell, S. Tong, M. de S.Z.S. Coelho, P.H.N. Saldiva, E. Lavigne, P.M. Correa, N.V. Ortega, H. Kan, S. Osorio, J. Kyselý, A. Urban, J.J.K. Jaakkola, N.R.I. Rytí, M. Pascal, et al. (2018), Temperature-related mortality impacts under and beyond Paris Agreement climate change scenarios, *Climatic Change* **150**, 3–4, 391–402, DOI: 10.1007/s10584-018-2274-3.
- Warren Liao, T. (2005), Clustering of time series data – A survey, *Pattern Recogn.* **38**, 11, 1857–1874, DOI: 10.1016/j.patcog.2005.01.025.

Received 17 November 2022

Accepted 20 December 2022