

CGC: a scalable Python package for co- and tri-clustering of geodata cubes

Francesco Nattino¹, Ou Ku¹, Meiert W. Grootes¹, Emma Izquierdo-Verdiguier², Serkan Girgin³, and Raul Zurita-Milla³

¹ Netherlands eScience Center, Science Park 140, 1098 XG Amsterdam, The Netherlands ² Institute of Geomatics, University of Natural Resources and Life Science (BOKU), 1190, Vienna, Austria ³ Faculty of Geo-Information Science and Earth Observation (ITC), University of Twente, PO Box 217, 7500 AE, Enschede, the Netherlands

DOI: [10.21105/joss.04032](https://doi.org/10.21105/joss.04032)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Hugo Ledoux](#) ↗

Reviewers:

- [@Subho07](#)
- [@Narayana-Rao](#)

Submitted: 17 December 2021

Published: 15 February 2022

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Multidimensional data cubes are increasingly ubiquitous, in particular in the geosciences. Clustering techniques encompassing their full dimensionality are necessary to identify patterns “hidden” within these cubes. Clustering Geodata Cubes (CGC) is a Python package designed for partitional clustering, which identifies groups of similar data across two (e.g., spatial and temporal) or three (e.g., spatial, temporal, and thematic) dimensions. CGC provides efficient and scalable co- and tri-clustering functionality appropriate to analyze both small and large datasets as well as a cluster refinement functionality that supports users in their quest to make sense of complex datasets.

Introduction

Faced with the increasing ubiquity of large datasets, data mining techniques have become essential to extracting patterns and generating insights. In this regard, clustering techniques, which aim to identify groups or subgroups with similar properties within a larger dataset, are becoming ever more popular.

Traditional clustering techniques focus on a single dimension and may therefore obfuscate relevant groups ([Cheng & Church, 2000](#); [Hartigan, 1972](#)). Hence, clustering techniques capable of simultaneously grouping data along multiple dimensions are needed. These techniques - referred to as co- or bi-clustering and tri-clustering in the case of two and three dimensions, respectively - have seen significant development and adoption in fields ranging from bioinformatics ([Cheng & Church, 2000](#)) to finance ([Shi et al., 2018](#)) and natural language processing ([Dhillon, 2001](#)).

The exponential growth of multidimensional data referring to geographical features (e.g., time series of satellite images) has resulted in a wide variety of geodata cubes, which can benefit from co- and tri-clustering. Indeed, following the development of a general information-theoretical approach ([Dhillon et al., 2003](#)) to partitional co-clustering ([Banerjee et al., 2007](#)), Wu et al. presented an application of co-clustering to geodata ([Wu et al., 2015](#)), as well as its extension to three dimensions ([Wu et al., 2018](#)).

In light of the eminent employability of co- and tri-clustering approaches to geospatial disciplines like geo-information science and Earth observation, and the transferability to other (geo)scientific domains, this paper presents the Clustering Geodata Cubes (CGC) package. CGC provides efficient and scalable co- and tri-clustering methods to analyze both small and

large datasets on single or multiple computing node systems. It also features a cluster refinement functionality that allows users to make sense of complex datasets more easily. An example of its application is the ongoing work on the analysis of spring onset datasets at high spatial resolution and continental scale, a preview of which is presented in the CGC tutorial. Although the package aims to meet the needs of geospatial data scientists, the algorithms implemented remain widely applicable and can easily be applied in other domains, which require performing partitional clustering of positive data matrices.

Statement of need

The CGC package focuses on the needs of geospatial data scientists who require tools to make sense of multi-dimensional data cubes by providing the following features and functionalities:

- **Partitional co- and tri-clustering methods suitable for spatiotemporal (multi-dimensional) data.** CGC includes clustering methods designed to simultaneously group elements into disjoint clusters along two or three dimensions. These methods are advantageous over one-dimensional clustering in that they provide a strategy to identify patterns that unfold across multiple dimensions, e.g. space and time. In addition, CGC provides functionality to refine the identified co- and tri-clusters. This post-processing step reduces the number of clusters to facilitate the identification and visualization of patterns.
- **Scalable clustering of small and big datasets.** CGC offers functionality to efficiently utilize available computational resources (ranging from single machines to computing clusters) and to tackle a wide range of dataset sizes. For single machine execution the package offers optimized support of multi-core CPUs and/or limited system memory. For large datasets CGC supports the use of distributed data and computation on computing clusters.
- **Easy integration into geospatial analysis workflows.** CGC is written in Python, which is widely used for geospatial scripting and applications, and employs Numpy ([Harris et al., 2020](#)) and Dask ([Dask Development Team, 2016](#)) arrays as input and output data types, guaranteeing seamless integration to the Python ecosystem and interoperability with the libraries prevalent in the field of big (geo)data. This furthermore ensures the interoperability of CGC with the Xarray package ([Hoyer & Hamman, 2017](#)), so that this versatile and popular tool can be used for data loading and manipulation before and after analyses with CGC.
- **Ease of use and reproducibility.** To facilitate community use and adoption, documentation and tutorials illustrating domain-science examples, applications, and use cases are available via the [publicly accessible repository](#), where development takes place, and which provides a platform for issue tracking. CGC is distributed via the [Python Package Index \(PyPI\)](#), and code-release snapshots archived on [Zenodo](#) facilitate reproducible analysis.

Algorithms

Co-clustering

CGC implements the Bregman block average co-clustering (BBAC) algorithm from [Banerjee et al. \(2007\)](#) as inspired by the MATLAB code of ([Merugu & Banerjee, 2004](#)). Briefly, the BBAC algorithm iteratively optimizes the clustering of rows and columns of a data matrix

83 starting from a random initial assignment until convergence. The information loss from the
84 original matrix to the clustered one, which is constructed as the matrix of the co-cluster
85 means, is minimized using a loss function that is based on the I-divergence (Dhillon et al.,
86 2003; Kullback & Leibler, 1951). CGC also supports a user-defined convergence threshold.
87 To limit the influence of the initial conditions on the final clustering and to avoid local minima
88 several runs are carried out, with the cluster assignments from the lowest loss function value
89 ultimately being selected. Number of runs can be altered by the user.

90 Note that in the CGC implementation of the algorithm, the update in the row- and column-
91 cluster assignments is computed only from the previous iteration's row and column. Contrarily
92 to the original MATLAB implementation (Merugu & Banerjee, 2004), this makes the algorithm
93 independent of the order in which the dimensions are considered, while still leading to an
94 optimal clustering solution.

95 Tri-clustering

96 For tri-clustering CGC implements the so-called Bregman cube average tri-clustering (BCAT)
97 algorithm, which is a generalization of the BBAC algorithm to three dimensions (Wu et
98 al., 2018). The algorithmic implementation with respect to the loss function and the update
99 schema are analogous to that described above for the co-clustering. Similarly, the tri-clustering
100 outcome is independent of the order in which the three dimensions of the input array are
101 provided.

102 Cluster refinement

103 The CGC package implements an optional, secondary cluster refinement step based on the k-
104 means method (Wu et al., 2016) and optimized using the Silhouette metric (Rousseeuw, 1987)
105 as implemented in the scikit-learn package (Pedregosa et al., 2011). This secondary grouping
106 is based on statistical properties of the co- or tri-clusters (see the [package documentation](#))
107 and helps to better capture patterns in the data by combining clusters and going beyond strict
108 checkerboard structures.

109 Related Software

110 A number of co-/bi-clustering implementations based on different algorithms exist. While
111 some of the available packages focus on specific applications, like gene expression data (Barkow
112 et al., 2006; Eren et al., 2012), generic co-clustering tools include biclust (R) and CoClust
113 (Python) (Role et al., 2019), as well as two implementations from the scikit-learn Python
114 library (Pedregosa et al., 2011). Most of the available algorithms target tabular data with a
115 hidden blocked-diagonal structure, where each row and column of the input matrix is assigned
116 to only one co-cluster. In contrast, in spatio-temporal data a set of spatial elements often
117 exhibits the same behaviour within multiple time windows and vice versa. This type of data
118 requires algorithms where a subset taken along one dimension can be associated to multiple
119 subsets taken along the other dimension. Partitioning algorithms of this type, such as the
120 BBAC and BCAT algorithms implemented in CGC, can discover checkerboard-like patterns
121 in the input data matrix. The package additionally enables the user to recover more of the
122 intrinsic structure of the data in a second step, going beyond the limits on structure imposed

123 Within the Python ecosystem, prominent implementations targeting checkerboard-like struc-
124 tures akin to CGC's initial step are the scikit-learn SpectralBiclustering algorithm, which imple-
125 ments the method from (Kluger et al., 2003), and the CoclustInfo algorithm from CoClust, also
126 based on information-theoretic co-clustering (Dhillon et al., 2003; Govaert & Nadif, 2018).

127 However, both approaches differ from CGC in the field of applicability. Spectral methods
128 like the former are fast, but their accuracy has been shown to be limited (Ailem et al., 2015;
129 Role et al., 2019). The latter focuses on datasets representing joint-probability distributions,
130 preserving statistics other than the co-cluster average (the row and column averages) in the
131 search for the optimal clustering solution (Banerjee et al., 2007). This requirement does not
132 generally apply to the geospatial datasets which CGC is targeting. Neither, however, offers
133 the ability to “break out” of the algorithmically-imposed checkerboard pattern.

134 Finally, to the best of our knowledge, CGC is unique in being designed from the outset for
135 use with big data (e.g. by including an implementation that supports multi-node distributed
136 computing), being able to analyze more complex data cubes via tri-clustering, and being able
137 to perform secondary cluster refinement.

138 Software package overview

139 The CGC software is structured in the following main modules, details of which are described
140 in the [online package documentation](#):

- 141 ▪ [coclustering](#), containing the following implementations of the co-clustering algorithm:
 - 142 – The [Numpy-based, vectorized single machine implementation with threading sup-](#)
143 [port for optimal usage of multi-core CPUs](#).
 - 144 – The [Numpy-based single machine implementation with a reduced memory foot-](#)
145 [print](#). This implementation trades performance for low memory usage, but uses
146 Numba’s just-in-time compilation (Lam et al., 2015) to mitigate performance loss.
 - 147 – The [Dask-based implementation](#). This implementation provides support for clus-
 - 148 tering large, out-of-core datasets by distributing the computation across multiple
149 nodes using Dask arrays.
- 150 ▪ [triclustering](#), containing the following tri-clustering implementations:
 - 151 – A [Numpy-based implementation](#) analogous to the co-clustering one described
152 above (note that the low-memory version is currently not available).
 - 153 – A [Dask-based implementation](#), also analogous to the corresponding co-clustering
154 version described above.
- 155 ▪ [kmeans](#), which implements the k-means cluster refinement step for both co- and tri-
156 clustering.
- 157 ▪ [utils](#), which includes a collection of utility functions e.g., memory consumption esti-
158 mation and cluster averaging.

159 Performance comparisons between the various [co-clustering](#) and [tri-clustering](#) implementations
160 are also briefly discussed in the package documentation.

161 Tutorial

162 The software package is complemented by an [online tutorial](#) (in the form of Jupyter notebooks
163 (Kluyver et al., 2016)) that illustrates how to perform cluster analysis of geospatial datasets
164 using CGC. Notebooks are made directly available via a [dedicated GitHub repository](#), and
165 are also published as [static web pages](#) for reference and linked to the [CGC documentation](#)
166 via a ‘Tutorials’ tab. The tutorials are designed to run on systems with limited CPU/memory

capacity, which, together with environment requirement specifications in a standardized format (conda YAML file) and binder badges, give users the possibility to easily run the notebooks live on `mybinder.org` (Jupyter Project et al., 2018).

Tutorials cover the following topics:

- Co- and tri-cluster analysis.
- K-means-based cluster refinement.
- Choice of the suitable implementation for a given problem size/infrastructure available.
- Loading of geospatial data, common data-manipulation tasks in a workflow involving CGC, visualization of the output.

Note that while the tutorial is aimed at geospatial uses cases, it illustrates some real-case applications that are likely to make it easier for users to carry out cluster analysis using CGC in other fields as well.

Acknowledgements

The authors would like to thank Dr. Yifat Dzigan for the helpful discussions and support and Dr. Romulo Goncalves for the preliminary work that led to the development of the software package presented here. We also would like to thank SURF for providing computational resources to test the first versions of the CGC package via the e-infra190130 grant.

Author contributions

All co-authors contributed to the conceptualization of the work, which was led by R.Z.M. and E.I.V.. F.N., M.W.G., and O.K. prepared the first draft of the tutorials, and wrote the initial draft of this manuscript. F.N. suggested changes to the co- and tri-clustering algorithms. R.Z.M., E.I.V., and S.G. led the design of experiments to test and improve the CGC package. R.Z.M. and E.I.V. helped to improve the tutorials. S.G. provided the required computational resources to run the experiments and tutorials and made suggestions for code optimizations. All co-authors reviewed and edited the final document.

References

- Ailem, M., Role, F., & Nadif, M. (2015). Co-clustering document-term matrices by direct maximization of graph modularity. *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, 1807–1810. <https://doi.org/10.1145/2806416.2806639>
- Banerjee, A., Dhillon, I., Ghosh, J., Merugu, S., & Modha, D. S. (2007). A generalized maximum entropy approach to bregman co-clustering and matrix approximation. *Journal of Machine Learning Research*, 8(67), 1919–1986. <http://jmlr.org/papers/v8/banerjee07a.html>
- Barkow, S., Bleuler, S., Prelić, A., Zimmermann, P., & Zitzler, E. (2006). BicAT: a biclustering analysis toolbox. *Bioinformatics*, 22(10), 1282–1283. <https://doi.org/10.1093/bioinformatics/btl099>
- Cheng, Y., & Church, G. M. (2000). Biclustering of expression data. *Proceedings of the Eighth International Conference on Intelligent Systems for Molecular Biology*, 93–103. ISBN: 1577351150

- 207 Dask Development Team. (2016). *Dask: Library for dynamic task scheduling*. <https://dask.org>
- 208
- 209 Dhillon, I. S. (2001). Co-clustering documents and words using bipartite spectral graph parti-
210 tioning. *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge*
211 *Discovery and Data Mining*, 269–274. <https://doi.org/10.1145/502512.502550>
- 212 Dhillon, I. S., Mallela, S., & Modha, D. S. (2003). Information-theoretic co-clustering. *Pro-*
213 *ceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery*
214 *and Data Mining - KDD '03*, 89. <https://doi.org/10.1145/956750.956764>
- 215 Eren, K., Deveci, M., Küçüktunç, O., & Çatalyürek, Ü. V. (2012). A comparative analysis
216 of biclustering algorithms for gene expression data. *Briefings in Bioinformatics*, 14(3),
217 279–292. <https://doi.org/10.1093/bib/bbs032>
- 218 Govaert, G., & Nadif, M. (2018). Mutual information, phi-squared and model-based co-
219 clustering for contingency tables. *Advances in Data Analysis and Classification*, 12(3),
220 455–488. <https://doi.org/10.1007/s11634-016-0274-6>
- 221 Harris, C. R., Millman, K. J., Walt, S. J. van der, Gommers, R., Virtanen, P., Cournapeau,
222 D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., Kerkwijk,
223 M. H. van, Brett, M., Haldane, A., Río, J. F. del, Wiebe, M., Peterson, P., ... Oliphant,
224 T. E. (2020). Array programming with NumPy. *Nature*, 585(7825), 357–362. <https://doi.org/10.1038/s41586-020-2649-2>
- 225
- 226 Hartigan, J. A. (1972). Direct Clustering of a Data Matrix. *Journal of the American Statistical*
227 *Association*, 67(337), 123–129. <https://doi.org/10.1080/01621459.1972.10481214>
- 228 Hoyer, S., & Hamman, J. J. (2017). Xarray: N-D labeled Arrays and Datasets in Python.
229 *Journal of Open Research Software*, 5, 10. <https://doi.org/10.5334/jors.148>
- 230 Jupyter Project, Bussonnier, M., Forde, J., Freeman, J., Granger, B., Head, T., Holdgraf, C.,
231 Kelley, K., Nalvarte, G., Osheroff, A., Pacer, M., Panda, Y., Perez, F., Ragan-Kelley, B.,
232 & Willing, C. (2018). *Binder 2.0 - Reproducible, interactive, sharable environments for*
233 *science at scale*. 113–120. <https://doi.org/10.25080/Majora-4af1f417-011>
- 234 Kluger, Y., Basri, R., Chang, J. T., & Gerstein, M. (2003). Spectral Biclustering of Microarray
235 Data: Coclustering Genes and Conditions. *Genome Research*, 13(4), 703–716. <https://doi.org/10.1101/gr.648603>
- 236
- 237 Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B., Bussonnier, M., Frederic, J., Kelley, K.,
238 Hamrick, J., Grout, J., Corlay, S., Ivanov, P., Avila, D., Abdalla, S., Willing, C., & Jupyter
239 development team. (2016). Jupyter notebooks - a publishing format for reproducible
240 computational workflows. In F. Loizides & B. Schmidt (Eds.), *Positioning and power in*
241 *academic publishing: Players, agents and agendas* (pp. 87–90). IOS Press. [https://](https://eprints.soton.ac.uk/403913/)
242 eprints.soton.ac.uk/403913/
- 243 Kullback, S., & Leibler, R. A. (1951). On Information and Sufficiency. *The Annals of Math-*
244 *ematical Statistics*, 22(1), 79–86. <https://doi.org/10.1214/aoms/1177729694>
- 245 Lam, S. K., Pitrou, A., & Seibert, S. (2015). Numba: A LLVM-based python JIT compiler.
246 *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*. <https://doi.org/10.1145/2833157.2833162>
- 247
- 248 Merugu, S., & Banerjee, A. (2004). *Bregman Co-clustering code in Matlab*. [http://www.](http://www.ideal.ece.utexas.edu/software.html)
249 [ideal.ece.utexas.edu/software.html](http://www.ideal.ece.utexas.edu/software.html)
- 250 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel,
251 M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau,
252 D., Brucher, M., Perrot, M., & Duchesnay, É. (2011). Scikit-learn: Machine learning in
253 python. *Journal of Machine Learning Research*, 12(85), 2825–2830. [https://www.jmlr.](https://www.jmlr.org/papers/v12/pedregosa11a.html)
254 [org/papers/v12/pedregosa11a.html](https://www.jmlr.org/papers/v12/pedregosa11a.html)

- 255 Role, F., Morbieu, S., & Nadif, M. (2019). **CoClust** : A *python* Package for Co-Clustering.
256 *Journal of Statistical Software*, 88(7). <https://doi.org/10.18637/jss.v088.i07>
- 257 Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of
258 cluster analysis. *Journal of Computational and Applied Mathematics*, 20, 53–65. [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7)
- 260 Shi, G., Ren, L., Miao, Z., Gao, J., Che, Y., & Lu, J. (2018). Discovering the trading pattern
261 of financial market participants: Comparison of two co-clustering methods. *IEEE Access*,
262 6, 14431–14438. <https://doi.org/10.1109/ACCESS.2018.2801263>
- 263 Wu, X., Zurita-Milla, R., Izquierdo-Verdiguier, E., & Kraak, M.-J. (2018). Triclustering
264 Georeferenced Time Series for Analyzing Patterns of Intra-Annual Variability in Tem-
265 perature. *Annals of the American Association of Geographers*, 108(1), 71–87. <https://doi.org/10.1080/24694452.2017.1325725>
- 267 Wu, X., Zurita-Milla, R., & Kraak, M.-J. (2015). Co-clustering geo-referenced time series:
268 Exploring spatio-temporal patterns in Dutch temperature data. *International Journal of*
269 *Geographical Information Science*, 29(4), 624–642. <https://doi.org/10.1080/13658816.2014.994520>
- 271 Wu, X., Zurita-Milla, R., & Kraak, M.-J. (2016). A novel analysis of spring phenological
272 patterns over Europe based on co-clustering: Co-Clustering European Spring Phenology.
273 *Journal of Geophysical Research: Biogeosciences*, 121(6), 1434–1448. <https://doi.org/10.1002/2015JG003308>
- 274

DRAFT