

A spectral approach for the clustering of source rocks

V. I. Makri^{a,b}, D. Pasadakis^c

^aInstitute of GeoEnergy (IG)-FORTH, Chania, Greece.

^bTechnical University of Crete, School of Mineral Resources Engineering, Chania, Greece.

^cUniversità della Svizzera italiana, Institute of Computing, Lugano, Switzerland.



- S. White and P. Smyth, *A spectral clustering approach to finding communities in graphs*, in Proceedings of the 2005 SIAM International Conference on Data Mining (SDM), pages 274-285, 2005.
- D. Pasadakis, C. L. Alappat, O. Schenk, and G. Wellein, *Multiway p-spectral graph cuts on Grassmann manifolds*, Machine Learning 111, pages 791-829, 2022.
- K. Peters, C. Walters, and J. Moldowan, *Geochemical correlation and chemometrics*, In The Biomarker Guide, pages 475-482, Cambridge University Press, 2004.

Spectral multiway clustering

For a graph $\mathcal{G}(V, E, \mathbf{A})$ with n vertices:

- Adjacency: $\mathbf{A} \in \mathbb{R}^{n \times n}$, degree $\mathbf{D} \in \mathbb{R}^{n \times n}$.
- Graph Laplacian matrix:

$$\mathbf{L} = \mathbf{D} - \mathbf{A} = \begin{bmatrix} \sum_j \mathbf{A}_{1j} & -\mathbf{A}_{12} & 0 & -\mathbf{A}_{14} \\ -\mathbf{A}_{12} & \sum_j \mathbf{A}_{2j} & -\mathbf{A}_{23} & -\mathbf{A}_{24} \\ 0 & \mathbf{A}_{23} & \sum_j \mathbf{A}_{3j} & -\mathbf{A}_{34} \\ -\mathbf{A}_{14} & -\mathbf{A}_{24} & -\mathbf{A}_{34} & \sum_j \mathbf{A}_{4j} \end{bmatrix}$$

Normalized variant: $\mathbf{L}_{rw} = \mathbf{D}^{-1}\mathbf{L}$, with

$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$, $\lambda_1 = 0$ for a connected graph.

Identify k clusters using the smallest k eigenvectors of \mathbf{L}_{rw} .

Rayleigh quotient minimization

$$\min_{U \in \mathbb{R}^{n \times k}} F(\mathbf{U}) = \text{Tr}(\mathbf{U}^\top \mathbf{L} \mathbf{U}),$$

subject to $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$.

Graph construction

- Input data $\mathbf{Y} \in \mathbb{R}^{n \times m}$, n : source rock samples, m : n-alkane peak areas.
 - Connectivity \mathbf{G} via a k-nearest neighbors (kNN) routine.
 - Similarity \mathbf{S} via a Gaussian kernel.
- $\mathbf{A} = \mathbf{G} \odot \mathbf{S}$ for a sparse graphical representation of the source rocks.

Spectral graph coordinates

$$\begin{array}{c|cccc} \mathbf{U} & u_1 & \dots & u_k \\ \hline \mathbf{U}_1 & u_{11} & \dots & u_{1k} \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{U}_n & u_{n1} & \dots & u_{nk} \end{array}$$

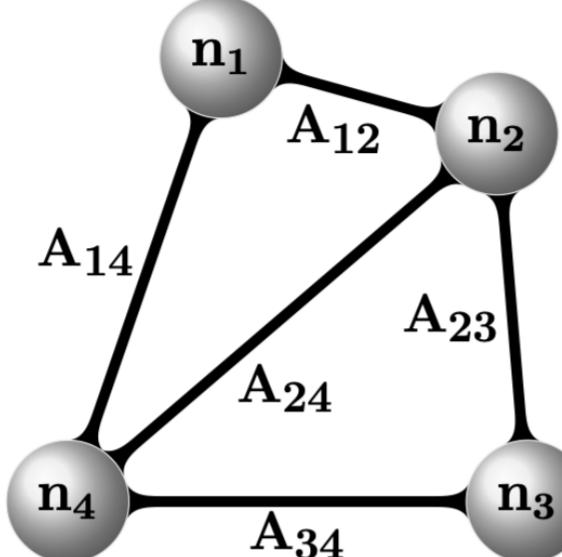
✓ Direct multiway clustering.

✓ Global information.

Maximize modularity function

$$Q = \sum_i \left(\mathbf{D}_{ii} - \left(\sum_j \mathbf{D}_{ij} \right)^2 \right)$$

→ High values: clear community structure.

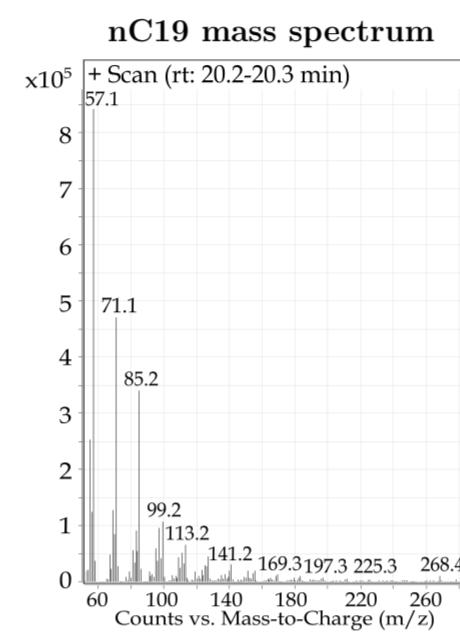


Data background

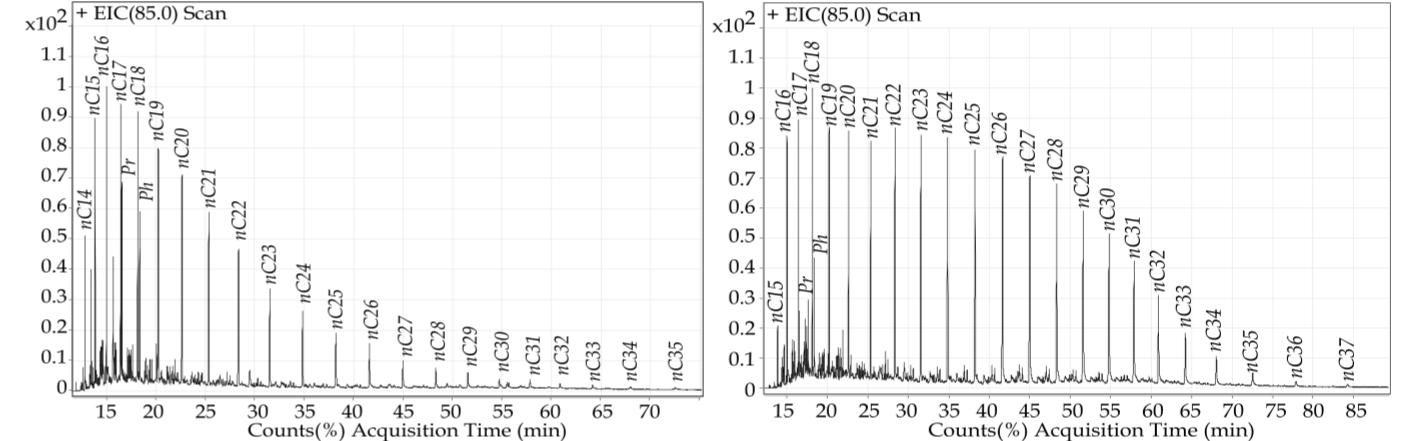
For this study, a set of 83 rock samples from Western Greece potential source rocks was used. These samples are of multiple geological ages, located in the so-called Ionian geotectonic zone. They were treated in the lab with the following experimental procedure:

Experimental procedure

- Soxhlet extraction using a DCM:MeOH mixture (90:10 v/v).
- De-asphalting by n-pentane (7ml) and filtration through Teflon syringe filters (0.45μm).
- SARA fractionation technique on silica-alumina chromatographic column for the separation of maltenes into saturated, aromatic and NSO fractions.
- GC-MS analysis on Agilent 7890A Gas Chromatograph (HP-5MS UI (60m x 250μm x 0.25μm column)) coupled to an Agilent 5975E Mass Spectrometer.
- Identification of peak areas of n-alkanes and isoprenoids by their relative retention times and mass spectra.
- Peak areas of nC15-nC35 were used for the multiway spectral clustering.



GC-MS chromatograms (m/z 85) of saturated fractions

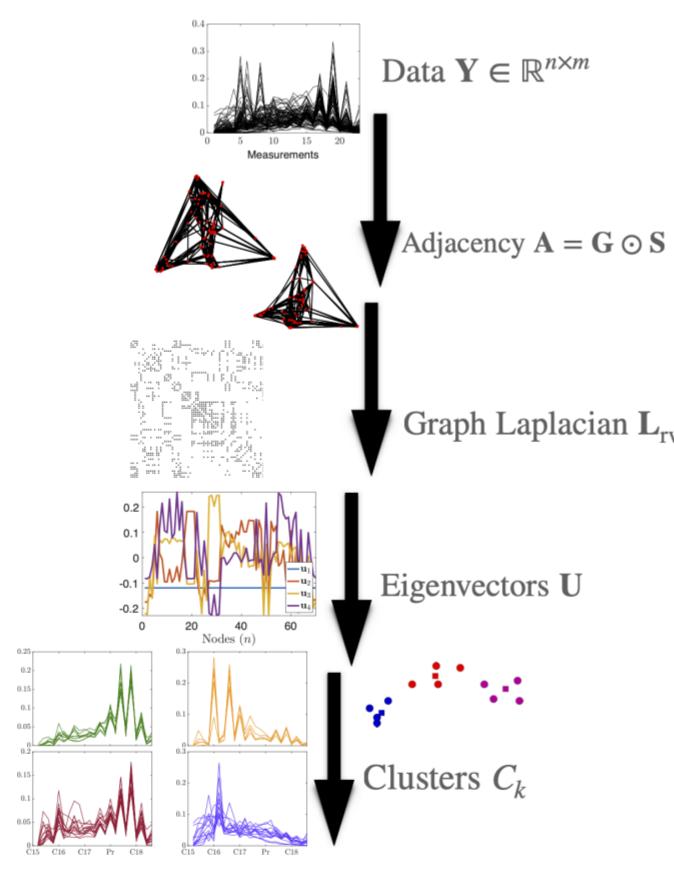


Key algorithmic components

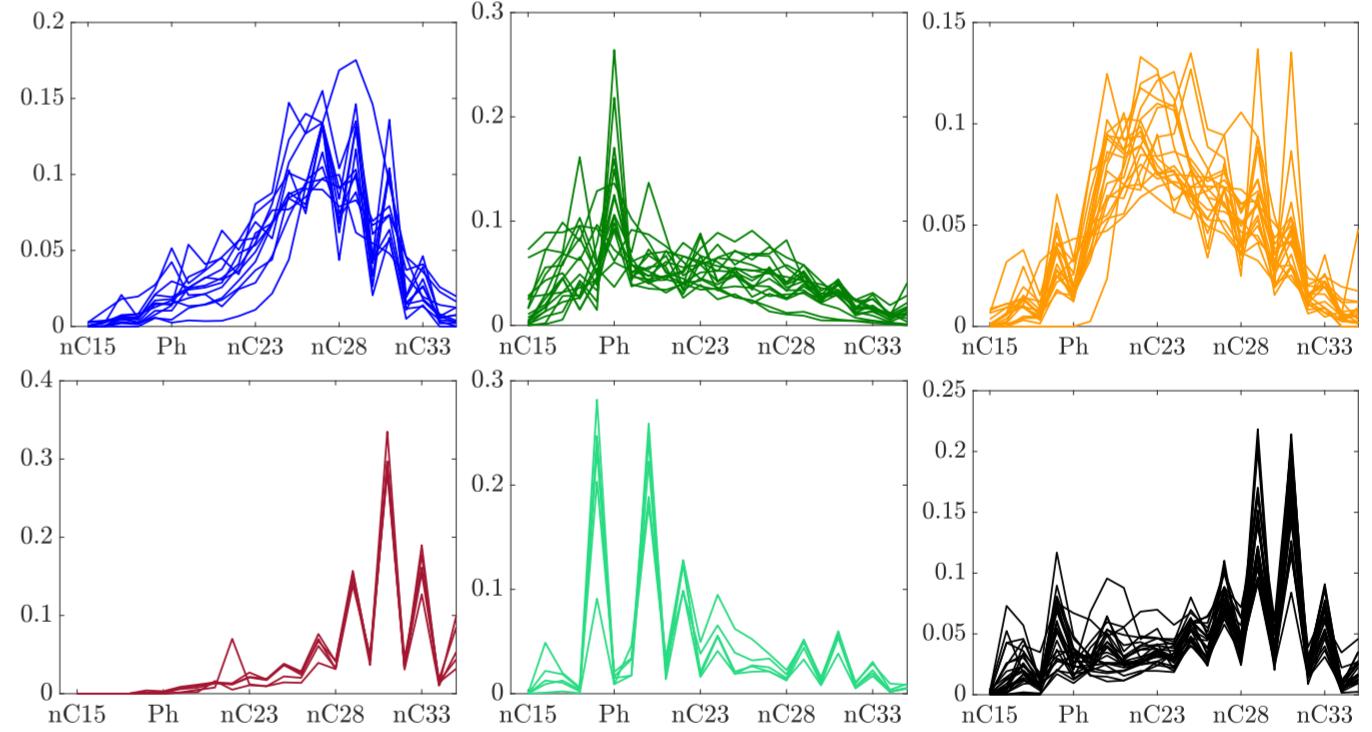
Spectral clustering of source rocks

Input: Data $\mathbf{Y} \in \mathbb{R}^{n \times m}$ ▶ Source rock measurements
Output: C_k ▶ Source rock clusters

- Data normalization ▶ Sum of peak areas
- Build adjacency matrix \mathbf{A} ▶ $\mathbf{A} = \mathbf{G} \odot \mathbf{S}$
- Construct the graph Laplacian \mathbf{L}_{rw}
- Rayleigh quotient min for the first $n/5$ eigenvectors \mathbf{U}
- while** $k \in [2, n/5]$ **do**
- Cluster \mathbf{U} in k groups ▶ Using kmeans
- Compute the modularity Q_k of the k clusters
- end while**
- Select the clusters C_k for which $Q = \max(Q_k)$

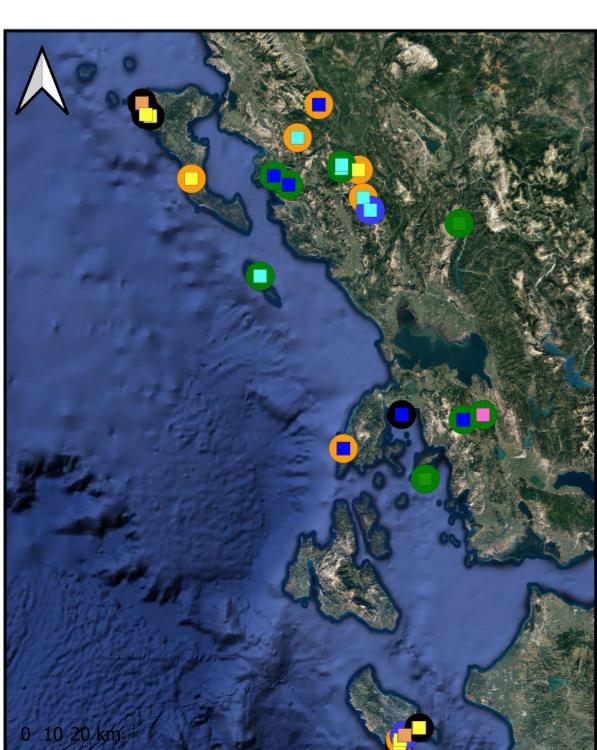


Numerical results



The direct multiway spectral clustering routine identifies 6 clusters in the dataset. The maximum value of modularity reported is $Q = 0.71$. Inspecting the chromatograms of each cluster we observe clear distinctions between the groups.

Map distribution



Legend

Sample Age
■ Early Cretaceous
■ Lower Jurassic
■ Mid-Late Jurassic
■ Miocene
■ Oligocene
■ Pliocene
■ Triassic

Cluster number
● 1
● 2
● 3
● 4
● 5
● 6

Interpretation and conclusions

Conclusions

- The 6 “naturally existing” groups showcase distinct concentration profiles of the n-alkanes of the samples.
- A geographical view of the clusters is shown on the map.
- Robust algorithm that can handle datasets of varying size and complexity.
- Higher accuracy compared to traditionally employed methods in Geochemistry (like hierarchical clustering and PCA)
- The reliability of this technique to support the classification of source rocks is proved.
- This technique may be employed in studies involving chromatographic analytical data such as oil fingerprinting, oil-oil and oil-source correlations, as well as bio-informatics.

Acknowledgements

This study is supported by Hellenic Petroleum S.A. (HELPE).

Estimating the number of clusters

Maximizing the modularity of the graph results in high quality clusters and leads to high intra-cluster and low inter-cluster connectivity.

