

Science of the Total Environment

Evaluation of natural background levels of high mountain karst aquifers in complex geological settings. A Gaussian mixture model approach in the Port del Comte (SE, Pyrenees) case study

--Manuscript Draft--

Manuscript Number:	STOTEN-D-20-22426
Article Type:	Research Paper
Keywords:	High-mountain karst system; Alpine hydrogeology; Natural background levels; Compositional data; Model-based clustering; Gaussian mixing model
Corresponding Author:	Jorge Jódar, Ph.D. (a)Geological and Mining Institute of Spain Zaragoza, SPAIN
First Author:	Ignasi Herms, MSc
Order of Authors:	Ignasi Herms, MSc Jorge Jódar, Ph.D. Albert Soler, Ph.D. Luis Javier Lambán, Ph.D. Emilio Custodio, Ph.D. Joan Agustí Núñez, BSc Georgina Arnó, BSc Maribel Ortego, Ph.D. David Parcerisa, Ph.D. Joan Jorge, Ph.D.
Abstract:	The hydrogeological processes driving the hydrochemical composition of groundwater in the alpine pristine aquifer system of the Port del Comte Massif (PCM) are characterized through the multivariate statistical techniques Principal Component Analysis (PCA) and Gaussian Mixture Models (GMM) in the framework of Compositional Data (CoDa) analysis. Also, the groundwater Natural Background Levels (NBLs) for NO ₃ and SO ₄ and Cl are evaluated, which are specially important for indicating occurrence of groundwater contamination derived from the anthropic activities conducted in the PCM. The hydrogeochemical facies in the aquifer system of the PCM comprises low mineralized Ca-HCO ₃ water for the main Eocene karst aquifer, but also Ca-SO ₄ and highly mineralized Na-Cl water types for the minor aquifers discharging from the PCM. The NBL values of SO ₄ , Cl and NO ₃ obtained for the main karst aquifer are 14.33, 4.06 and 6.55 mg/L, respectively. These values are 35, 3 and 1.2 times lower than the respective official NBLs values that were determined by the water administration to be compared with in the case of conducting a polluting assessment characterization in the main karst aquifer. Official overestimation of NBLs can put important groundwater resources in the PCM at risk.
Suggested Reviewers:	Tibor Yvan Stigter, Ph.D. IHE Delft Institute for Water Education Department of Water Science & Engineering t.stigter@un-ihe.org Specialist in groundwater resources management Miguel Angel García-Vera, Ph.D. Head of water resources planning, Head of water resources planning in Confederación Hidrográfica del Ebro mgarciaive@chebro.es Specialist in water resources planning from the Water resources Administration point of view

	<p>Miguel Rodríguez-Rodríguez, Ph.D. Researcher, Universidad Pablo Olavide Departamento de Sistemas Físicos Químicos y Naturales: Universidad Pablo Olavide Departamento de Sistemas Físicos Químicos y Naturales mrodrod@upo.es Specialist in Karst hydrology</p>
	<p>Josep Mas-Pla, Ph.D. Researcher, Universitat de Girona josep.mas@udg.edu; Specialist in groundwater resources management</p>
	<p>Juan Antonio Barberá, Ph.D. Researcher, Universidad de Malaga jabarbera@uma.es Specialist in Karst hydrology</p>
	<p>Javier Heredia, Ph.D. Researcher, Instituto Geológico y Minero de España: Instituto Geológico y Minero de España j.heredia@igme.es Specialist in Karst hydrology</p>
	<p>Caoimhe Hickey, Ph.D. Researcher, Geological Survey of Ireland Caoimhe.Hickey@gsi.ie Specialist in Karst hydrology</p>
	<p>Vera Pawlowsky-Glahn, Ph.D. Researcher, Universitat de Girona vera.pawlowsky@udg.edu Specialist in Compositional data Analysis</p>
	<p>Juan José Egozcue, Ph.D. Researcher, Universitat Politècnica de Catalunya juan.jose.egozcue@upc.edu Specialist in Compositional data Analysis</p>
	<p>Mark Engle, Ph.D. Researcher, University of Texas maengle@utep.edu Specialist in Compositional data Analysis</p>
Opposed Reviewers:	

To whom concern:

Knowledge on natural background levels of concentrations in groundwater is important to identify contamination of groundwater by anthropogenic activities. This is ever more important in high mountain aquifers where water resources are mostly generated for the downgradient depending areas. This is the case of the Port del Comte karst aquifer system discharges into the Cardener River, the most important tributary of the populated Llobregat River basin that provides drinking water to the city of Barcelona (NE Spain), where the population is larger than 2 million people, and where more than 70% of the water supply derives from surface waters. These waters also feed one of the major groundwater reservoirs for the Barcelona area, the aquifers of the Llobregat Delta (Otero et al., 2008).

High mountain aquifers are very sensitive to polluting activities. Management and preservation of these systems requires a well defined hydrochemical picture to compare with and warning in case of aquifer polluting threats affecting the water resources depending systems.

Processes such as aquifer recharge and runoff generation which force and describe the hydrological systems response also belong to the main knowledge areas considered by Science of the Total Environment, by including hydrosphere, biosphere, lithosphere, and anthroposphere.

Barcelona, 12/09/2020

The authors: Ignasi Herms. Jorge Jódar, Albert Soler, Luis Javier Lambán, Emilio Custodio, Joan Agustí Núñez, Georgina Arnó, Maribel Ortego, David Parcerisa and Joan Jorge

1 **Evaluation of natural background levels of high mountain karst**
2
3
4
5
6
7
8
9
10
11
12 **aquifers in complex geological settings. A Gaussian mixture model**
13
14 **approach in the Port del Comte (SE, Pyrenees) case study**

15
16
17
18
19 Herms, I.^a, Jódar, J.^{b*}, Soler, A.^c, Lambán, L.J.^b, Custodio, E.^d, Núñez, J.A.^a, Arnó, G.^a, Ortego,
20
21 M.I.^e, Parcerisa, D.^f, Jorge, J.^f

22
23
24 (a) Àrea de Recursos Geològics. Institut Cartogràfic i Geològic de Catalunya (ICGC), Barcelona,
25
26 Spain

27
28 (b) Instituto Geológico Minero de España (IGME), Zaragoza, Spain

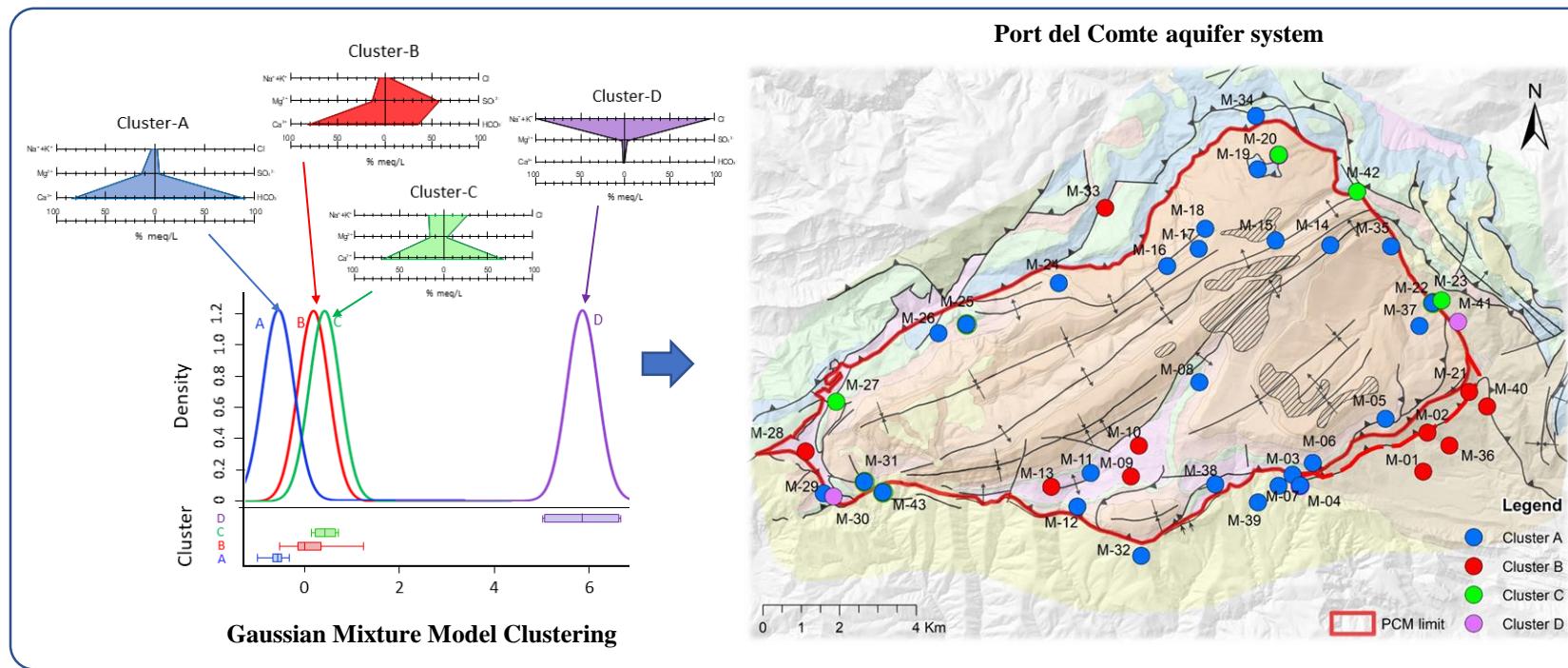
29
30 (c) Grup MAiMA, SGR Mineralogia Aplicada, Geoquímica i Geomicrobiologia, Departament de
31
32 Mineralogia, Petrologia i Geologia Aplicada, Facultat de Ciències de la Terra, Universitat de
33
34 Barcelona (UB), Barcelona, Spain

35
36 (d) Real Academia de Ciencias. Grupo de Hidrología Subterránea, Dept. de Ingeniería Civil y
37
38 Ambiental, Universitat Politècnica de Catalunya (UPC), Barcelona, Spain

39
40 (e) Compositional and Spatial Data Analysis (COSDA) Research Group. Department of Civil and
41
42 Environmental Engineering, Universitat Politècnica de Catalunya BarcelonaTech, Spain

43
44 (f) Departament d'Enginyeria Minera, Industrial i TIC. Universitat Politècnica de Catalunya
45
46 (UPC), Manresa, Spain

47
48 * Corresponding author: j.jodar@igme.es



The Gaussian Mixture Models are used to cluster springs by hydrochemical response

The clustering analysis groups the springs of the aquifer system into four families

The compositional analysis approach is required for hydrochemical data analysis

The hydrogeological characterization of Port del Compte Aquifer system is presented

Natural base levels of nitrate and sulphate in the Port del Compte Aquifers are given

1 **Evaluation of natural background levels of high mountain karst**
2 **aquifers in complex geological settings. A Gaussian mixture model**
3 **approach in the Port del Comte (SE, Pyrenees) case study**

4

5 Herms, I.^a, Jódar, J.^{b*}, Soler, A.^c, Lambán, L.J.^b, Custodio, E.^d, Núñez, J.A.^a, Arnó, G.^a,
6 Ortego, M.I.^e, Parcerisa, D.^f, Jorge, J.^f

7

8 (a) Àrea de Recursos Geològics. Institut Cartogràfic i Geològic de Catalunya (ICGC),
9 Barcelona, Spain

10 (b) Instituto Geológico Minero de España (IGME), Zaragoza, Spain

11 (c) Grup MAiMA, SGR Mineralogia Aplicada, Geoquímica i Geomicrobiologia,
12 Departament de Mineralogia, Petrologia i Geologia Aplicada, Facultat de Ciències de la
13 Terra, Universitat de Barcelona (UB), Barcelona, Spain

14 (d) Real Academia de Ciencias. Grupo de Hidrología Subterránea, Dept. de Ingeniería
15 Civil y Ambiental, Universitat Politècnica de Catalunya (UPC), Barcelona, Spain

16 (e) Compositional and Spatial Data Analysis (COSDA) Research Group. Department of
17 Civil and Environmental Engineering, Universitat Politècnica de Catalunya
18 BarcelonaTech, Spain

19 (f) Departament d'Enginyeria Minera, Industrial i TIC. Universitat Politècnica de
20 Catalunya (UPC), Manresa, Spain

21 * Corresponding author: j.jodar@igme.es (J.Jódar)

22

23 **Abstract**

24 The hydrogeological processes driving the hydrochemical composition of groundwater in
25 the alpine pristine aquifer system of the Port del Comte Massif (PCM) are characterized
26 through the multivariate statistical techniques Principal Component Analysis (PCA) and
27 Gaussian Mixture Models (GMM) in the framework of Compositional Data (CoDa)
28 analysis. Also, the groundwater Natural Background Levels (NBLs) for NO₃ and SO₄
29 and Cl are evaluated, which are specially important for indicating occurrence of
30 groundwater contamination derived from the anthropic activities conducted in the PCM.

31

32 The hydrogeochemical facies in the aquifer system of the PCM comprises low
33 mineralized Ca-HCO₃ water for the main Eocene karst aquifer, but also Ca-SO₄ and
34 highly mineralized Na-Cl water types for the minor aquifers discharging from the PCM.
35 The NBL values of SO₄, Cl and NO₃ obtained for the main karst aquifer are 14.33, 4.06
36 and 6.55 mg/L, respectively. These values are 35, 3 and 1.2 times lower than the
37 respective official NBLs values that were determined by the water administration to be
38 compared with in the case of conducting a polluting assessment characterization in the
39 main karst aquifer. Official overestimation of NBLs can put important groundwater
40 resources in the PCM at risk.

41

42 **Keywords:** High-mountain karst system; Natural background levels; Compositional
43 data; Model-based clustering; Gaussian mixing model.

44

45

46 1. Introduction

47 High mountain zones produce globally essential water resources that feed with fresh water
48 to the lowland depending ecosystems and a large portion of the world's population
49 ([Viviroli et al., 2020](#)). Mountain aquifers, specially those developed in karstifiable
50 carbonate rocks, storage the infiltrated precipitation, thus maintaining important
51 groundwater resources. These resources are typically released as groundwater through
52 large springs that regulates the hydro-ecological regime of the downstream rivers ([Kresic](#)
53 and [Stevanović, 2010](#)), and provide water resources during the dry season conspicuously
54 in semi-arid regions, where they are often the primary source of drinking water
55 ([Stevanović, 2019](#)).

56

57 Karst aquifers are much more vulnerable to pollution than other aquifers. Contaminants
58 may easily enter the subsurface into the karst system and rapidly spread in the conduit
59 system without any substantial attenuation ([Marín and Andreo, 2015](#)), threatening the
60 water resources of a region at large scale. These aquifers need special protection ([Drew](#)
61 and [Hötzl 1999](#), [Zwahlen, 2004](#)). In this line, the European Union enacted the Water
62 Framework Directive (2000/60/EC) ([WFD, 2000](#)) as an integrated approach to water
63 management to minimize the likelihood of pollution to enter the aquifers. The [WFD](#)
64 ([2000](#)) defines the rules that apply to the identification of the different groundwater bodies

65 (GWB), but also the criteria for chemical status assessment through defining pollutants
66 threshold values (TVs) and groundwater natural background values (NBLs). The former
67 are quality standards for pollutants in groundwater representative of those groundwater
68 bodies considered to be at risk. The latter provides the information regarding the
69 concentration of a given element, species or chemical substance present in solution which
70 is derived by natural processes from geological, chemical, biological and atmospheric
71 sources ([Müller et al., 2006](#)). In other words, NBLs are the corner stone to quantitatively
72 evaluate whether groundwater is significantly affected or modified by anthropogenic
73 influences ([Nieto et al., 2005; Custodio et al., 2007](#)).

74

75 It is not easy to define NBLs in high mountain karst aquifer systems (HMKS). For a given
76 aquifer and a certain element, the corresponding NBL value is obtained by averaging the
77 dissolved content of that element in groundwater discharge for the different springs
78 draining the aquifer. HMKS are usually embedded in complex geological structures as a
79 result of tectonic processes (e.g. faults, fold-and-thrust belts, wedge pinch out layers).
80 This often causes a strong compartmentalization ([Ballesteros et al, 2014](#)) that may involve
81 different lithologies (i.e. from carbonates to evaporites), thus generating a complex
82 aquifer system. The geological variability of such aquifer system influences the
83 hydrogeochemical signature of groundwater along the different flowlines, which typically
84 converge while mixing around springs. As a result, a different hydrochemical
85 composition than the expected may be obtained in the discharge of a spring given its
86 geological setting ([Lambán et al., 2015](#)), thus complicating a consistent NBLs
87 characterization for the different aquifers conforming the hydrogeological system.

88

89 To correctly define NBLs in HMKS it is necessary to have both a good hydrogeological
90 characterization of the aquifers at the local scale, and a good characterization of the
91 relevant hydrogeochemical finger prints describing the whole picture of the aquifer
92 system. In the framework of compositional data analysis (CoDA) ([Aitchison, 1986](#)), in
93 which the concentration of a given element dissolved in water is actually expressing a
94 part of a whole, regardless of the dimensions in which the component concentration is
95 expressed, either as weight per cent ratio [-] (e.g., %, mg/kg), or given as element mass
96 per unit of dissolution volume [ML^{-3}] (e.g., mg/L) ([Reimann et al., 2012; Buccianti and](#)
[Grunsky, 2014; Filzmoser et al., 2018; Egozcue and Pawlowsky-Glahn, 2005, 2006;](#)
[Pawlowsky-Glahn, et al. 2015](#)). For CoDa multivariate statistical analysis (MSA)

99 techniques/tools have shown a proven track record in characterizing complex
100 hydrogeological systems through the analysis of spatial variations in hydrochemical data.
101 The combination of MSA tools (e.g. principal component analysis and clustering
102 analysis) allow to investigate the factors controlling the processes taking place in aquifers
103 driving the hydrogeochemical content of groundwater (Otero et al., 2005; Puig et al.,
104 2011; Blake et al., 2016; Owen et al., 2016; Piña et al., 2018; Shelton et al., 2018).
105 Clustering analysis (CA) methods have been largely used to separate groundwater
106 samples, especially for large and/or complicated datasets, into homogeneous groups to
107 reflect different source contributions to groundwater in the sampled springs (see Suk and
108 Lee, 1999; Cloutier et al., 2008; Yidana, 2010; Kim et al., 2014; Yolcubal et al 2019,
109 among others). This faculty makes CA methods a promising tool to correctly define NBLs
110 in HMKS.

111

112 There are two mainstreams in CA, (1) the “hard clustering” methods like hierarchical
113 clustering and partitioning methods (k-means, k-medoids: PAM, CLARA) where each
114 data point (i.e. the sample) is assigned to one and only one cluster (hard assignment) and
115 (2) the “soft clustering” methods like model-based clustering (e.g. the Gaussian Mixture
116 Models – GMM) and fuzzy clustering where instead of assigning each data point into a
117 unique specific cluster, it is assigned to all the clusters with different probabilities or
118 weights (soft assignment) (Güler and Thyne 2004).

119

120 Soft clustering methods are getting more popular since they provide degrees of
121 membership at different hydrogeochemical clusters, rather than clear-cut distinctions. As
122 a result, they can better reflect the spatial continuity of a hydrological system while
123 providing a more rigorous framework to validate the clustering results (Kim et al., 2014,
124 2015; Wu, et al., 2017, Bondu et al., 2020). Moreover, in the framework of HMKA where
125 the limited number of observations often is a challenge, GMM clustering algorithms are
126 shown able to provide valuable insights into hydrochemical processes, delineating the
127 different groundwater sources imprinting the hydrochemical signature of the aquifer
128 system despite a sparse hydrochemical dataset (Wu et al., 2017). GMM are specially well
129 suited to provide a solid basement for NBLs determination in HMKS. Altough GMM
130 have been used for some authors to evaluate NBLs (Kim, et al. 2015), surprisingly, there
131 are no references in the scientific literature that use GMM in the framework of CoDA to
132 evaluate NBLs in HMKS.

133
134 This work aims at filling this gap. To that end, we characterize the hydrochemical
135 composition of the different aquifers associated to the alpine karst aquifer system of the
136 Port del Comte Massif (PCM) to evaluate in a consistent way the NBLs for the different
137 aquifers integrated in this HMKS. This is conducted through a MSA approach that
138 combines in a CoDA framework both PCA and GMM clustering analysis.

139

140

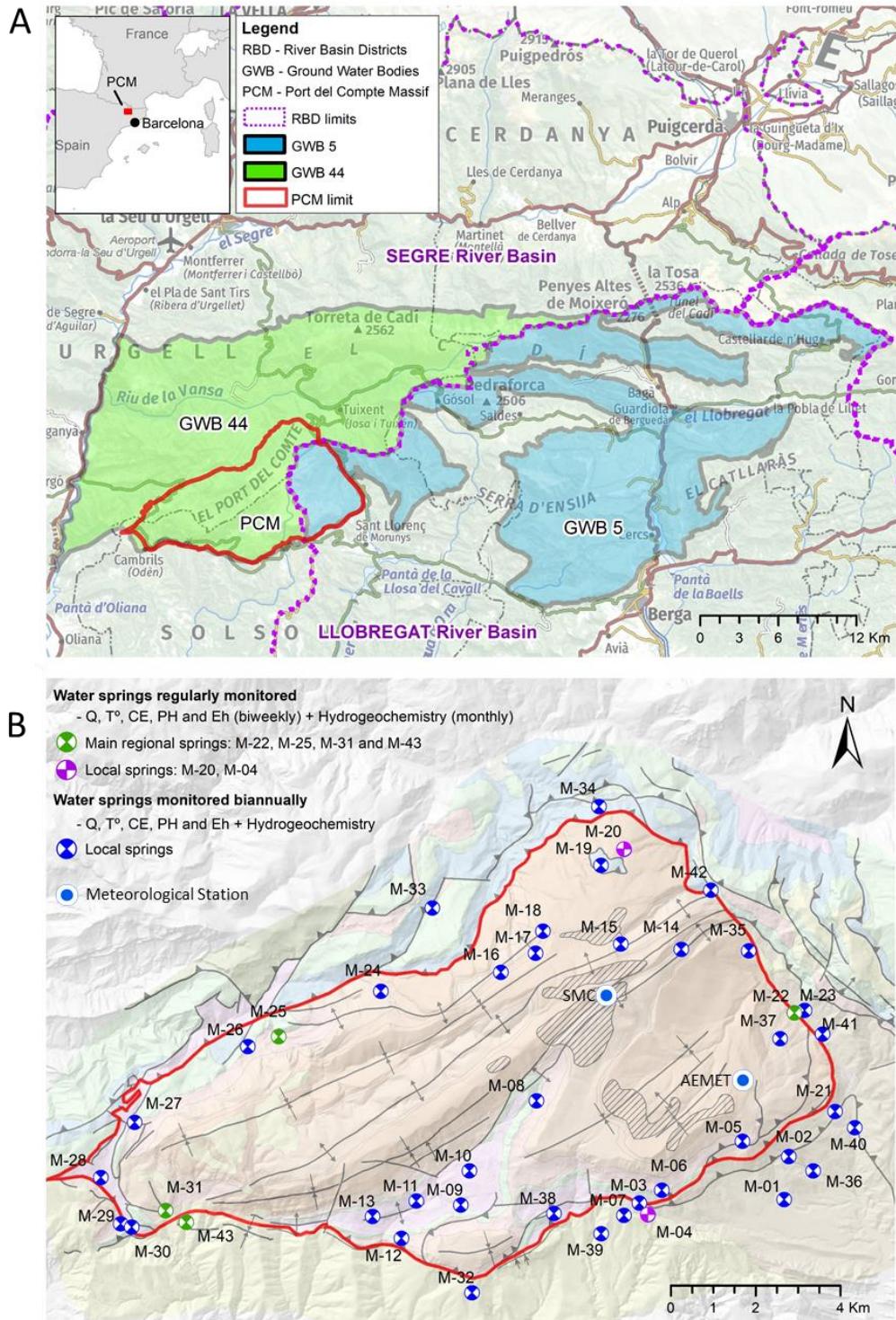
141 **2. Study area**

142 The PCM is located in the South-Central Catalan Pyrenees (north-east of Spain), which
143 constitute an orogenic system that runs along the boundary between the Iberian and
144 European plates. It was formed from the Late Cretaceous to Miocene times ([Muñoz, et al](#)
145 [2018](#)). The elevation of the mountainous massif ranges from 900 m a.s.l. to 2390 m a.s.l.
146 The massif constitutes an independent structural and hydrogeological system with a
147 surface area of 110 km². The highest peaks of the massif conform a water divide between
148 the upper Segre River basin to the NW and SW and the upper Cardener River basin (a
149 tributary of the LLobregat River) to the SE ([Fig. 1](#)).

150

151 According to the Köppen-Geiger classification ([Peel et al., 2007](#)), the study area is
152 characterized by a cold climate without a dry season and with a temperate summer. For
153 the period 2005-2019, the average annual precipitation (P), temperature (T) and potential
154 evapotranspiration (Hargreaves' method) at the SMC meteorological station located at
155 2315m a.s.l. ([Fig. 1](#)) are 1055 mm, 3.2° C and 525 mm, respectively. At elevations > 1800
156 m a.s.l. the snow covers the massif from December to March.

157



158

159

160 **Fig.1.** (A) Location map of the study area. (A). Delimitation of the groundwater bodies
161 affecting the PCM; GWB-44 belongs to the Segre river basin, and GWB-5 belongs to the
162 Llobregat river basin. (B) Location of the 43 monitored springs in the PCM

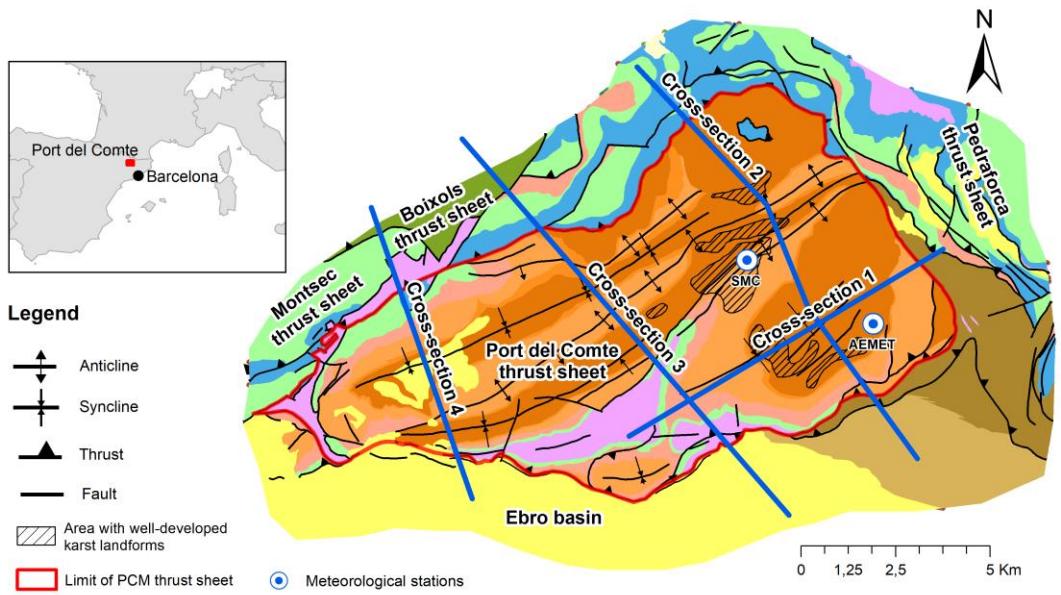
163

164 Geologically, the PCM constitutes an independent thrust sheet which presents complex
165 structural shapes in its boundaries ([Fig. 2](#)), with different thrust sheets individualizing the
166 whole domain in one independent structural system. The internal structure of the PCM is
167 formed by a set of folds and thrusts. These folds have a constant direction NE-SW parallel
168 to the NW limit ([Vergés, 1999](#)). The stratigraphic series contains limestones and
169 evaporites mainly from the Triassic, Cretaceous limestones, calcarenites, and shales and
170 Paleogene: Eocene-Oligocene limestones, sandstones and marls. The Jurassic marls,
171 limestones and dolomites only outcrops in the NW part of the sheet. The limestones
172 present a total thickness greater than 1300 m. From the geomorphological perspective,
173 the PCM presents a rounded-soft landscape in the highest domains with no vegetation
174 cover and almost no soil horizon development. The rest of the massif is covered by
175 mountain meadows and forest with a shallow soil depth up to medium development
176 ground cover. Many different karst forms appear progressively from 1950 m.a.s.l.
177 upwards, being well developed at 2050m a.s.l. (see [Fig 2](#), indicated as 'Area with well-
178 developed karst landforms') with sinkholes, dolines and karren fields. They underline the
179 heterogeneity of the karst system.

180

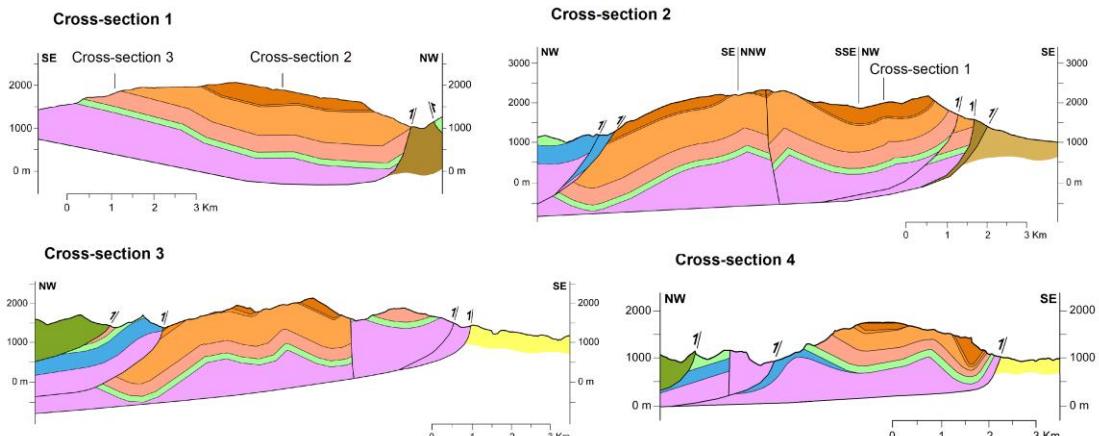
181 From the hydrogeological point of view, the PCM can be considered an independent unit
182 and a multi-aquifer system. The main aquifer is formed by Lower Eocene – fissured and
183 karstified limestones and dolomites and constitute one of the most important karst
184 aquifers of the Catalan Pyrenees. The other existing aquifers and aquitards in the system
185 are related to the levels of Cretaceous limestones, Triassic limestone and evaporites, other
186 Paleogene conglomerates and sandstones, and also to small quaternary aquifers (draining
187 small areas), which can be recharged locally at low or medium elevations. The lower
188 Upper Cretaceous/Paleocene (Garumnian facies) substrate materials composed by
189 siltstone and shales constitute an impervious layer for the overlaying Lower Eocene karst
190 aquifer. The complex geologic structure of the system strongly influences the location of
191 the existing karst springs, their groundwater geochemistry and their long-term hydrologic
192 behaviour.

193



Stratigraphy (Note: The sketch shows a simplified version of the official geological database 1: 50,000 - BG50M - of the ICGC (2007). The epigraphs in parentheses correspond to the geological units of the BG50M where the springs studied in this investigation are located).

Triassic - Shales, limestones, dolomites and evaporites (Tk, Tm)	Jurassic - Marl, bioclastic limestones and dolomites (TJb, TJcd)
Lower Cretaceous - Micritic limestone-marl alternations	Upper Cretaceous - limestone-marl alternations and calcarenites (Kat, KMca)
Garumnian (Upper Cretaceous-Lower Paleogene): shales, marls and limestone (Kgp), multicoloured clay deposits 'redbed' facies.	
Lower Eocene - Fissured/karstified alveoline limestones and dolomites (PPEc). Includes colluvial quaternaries that partially overlap (Qpe, Qvl)	Lower Eocene - Marl, sandstones and limestones (PEci) Lower Eocene - Fissured/karstified micritic and bioclastic limestones (PEcp1, PEcp2)
Middle Eocene - Sandstones, marls, conglomerates, limestones and evaporates (PEalb, PEm1, PEmb). Includes colluvial quaternary deposits (Qcoo) and alluvial (Qoo) partially overlap	Upper Eocene - Alluvial systems: conglomerates and sandstones
conglomerates breccias and sandstones (POcgs, POmlg)	Oligocene - Alluvial systems:



194

195 **Fig. 2.** Geological map and geological cross-sections of the PCM (modified from [ICGC](#),
196 [2007](#))

197 The hydrogeological conceptual model of the PCM aquifer system, as presented by
198 [Herms, et al. \(2019\)](#), considers that recharge is produced by infiltration of precipitation
199 as rainfall and snowmelt and it occurs both in a concentrated way through the local karst
200 conductive features, mostly situated at the top of the massif, and diffuse through the whole

201 domain. The infiltrated water percolates through the thick unsaturated zone (with more
202 than 1000 m at the top of the massif) towards the saturated zone, and discharges through
203 a large number of springs.

204

205 More than 100 springs were inventoried in the study zone. Nevertheless, only 43 of them
206 discharge throughout the year ([Fig. 1](#)). These springs were monitored during the period
207 September 2013 – October 2015. Most of them discharge small-scale local sub-surface
208 water flows with flow rates ranging between 0.1 L/s to 1 L/s. Nevertheless, there are four
209 ‘regional’ springs (M-22, M-25, M-31 and M-43) with flow rates comprised between 1
210 L/s and 900 L/s during the monitored period. These regional springs are recharged at
211 medium to high elevations, and drain the system discharging through the limestones
212 outcrops (M-31), quaternary deposits overlying the limestones (M-25, M-22), and also
213 through well-developed karst conduits in the conglomeratic materials of the Ebro Basin
214 (M-43). These conglomerates conform the southern foreland basin of the Pyrenees, which
215 is located just at the southern border of the PCM. There is also a diffuse groundwater
216 discharge through the ‘Riu Fred’ sub-basin to the North. With the exception of two
217 singular groundwater wells on the SW and E edges of the PCM, there are no other water
218 wells within the perimeter of the PCM that exploit the main karst aquifer. It is estimated
219 that the regional water table of the karst system is between 1000 and 1100 m a.s.l. ([Herms et al., 2019](#)).

221

222 Although the whole PCM massif belongs to the same geomorphological structure, the SE
223 sector has been assigned to GWB-5 (‘Conca Alta del Cardener i Llobregat’), whereas the
224 rest of the PCM was assigned to GWB-44 (‘Cadí Port del Comte’). [Table 1](#) summarizes
225 the natural background levels at the 90th percentile values (NBL90) determined through
226 the Pre-selection (PS) method described by the EU research project “BRIDGE” ([2007](#))
227 ([Müller et al., 2006](#)) using different control points for each GWB. It is worth noting the
228 high values determined for SO₄ contents in both GWBs. Considering an aquifer of pristine
229 waters related to the main Eocene karst aquifer of PCM composed only of limestone and
230 dolomites, these values are surprisingly high.

231

232 **Table 1.** NBL90 values for Cl, NO₃ and SO₄ in the GWB-5 and GWB-44.

NBL90		
Cl [mg/L]	SO ₄ [mg/L]	NO ₃ [mg/L]

GWB-5	12 ^a	485 ^a	-
GWB-44	36 ^b	609 ^b	8 ^b

(a) Data source: Agència Catalana de l'Aigua
(b) Data source: Confederación Hidrográfica del Ebro

233

234

235

3. Materials and methods

236

3.1. Sampling and analysis

237 In this work 43 springs were sampled twice per year (i.e. before snowfall and after
238 snowmelt seasons) between September 2013 and October 2015. Nevertheless, in six of
239 them (M-04, M-20, M-22, M-25, M-31 and M-43; Fig. 1) the groundwater sampling
240 frequency was higher, every three to four weeks, to study hydrogeochemical evolution of
241 the groundwater discharge. The springs M-22, M-25, M-31 and M-43 correspond to
242 regional discharge points of the karst system, whereas M-04 and M-20 are considered
243 representative of local small aquifers of the area (Herms et al., 2019).

244

245 A total of 288 groundwater samples were collected. Additionally, 10 snow samples (7
246 from natural snow and 3 from artificial snow produced in the existing sky resort in the
247 NE zone of the PCM) and 2 water samples from water ponds destined to artificial snow
248 production were collected. In all cases, the in situ physico-chemical parameters
249 Temperature (T), electrical conductivity (EC), pH, Eh and the total dissolved solid (TDS)
250 were measured. The geochemical analysis considered major cations and anions.

251

252 All samples were filtered using a 0.45µm membrane filter and stored in new 200-500 mL
253 polyethylene bottles washed with diluted nitric acid and rinsed with samples prior to
254 sampling. Samples for cation analysis were acidified with ultrapure HNO₃, to pH<2 to
255 prevent precipitation. Samples for anion analysis were not acidified. All water samples
256 were preserved at 4 °C before laboratory measurement. T, CE, pH, Eh and TDS were
257 measured by a portable Hanna meter (Multiparameter Water Quality Meter HI9829).
258 Total alkalinity was also determined in-situ by using the titration method in the first four
259 campaigns as well as with an Alkalinity Test Checker de (HI755) of Hanna Instruments.
260 The major cations (Ca, Mg, Na, K, NH₄) and anions (Cl, NO₃, HCO₃, CO₃, SO₄, and F)
261 were determined in the Laboratori Ambiental d'Aigües de Terrassa: the cations were

262 analysed by inductively coupled plasma atomic emission spectrometry – ICP-OES
263 (Agilent 5100 DV), and the anions by ion chromatography (Dionex, DX-120). Ionic
264 balance errors were calculated using the USGS software PHREEQC ([Parkhurst and](#)
265 [Appelo, 2013](#)) within the version PhreeqC Interactive (version 3.3.3 10424), and with the
266 phreeqc.dat database, except for the most salinized natural waters (M-30 and M-41)
267 related to deep flow through Keuper evaporates. The majority of samples had ionic
268 balance errors below the recommended standard of ±5% ([Appelo and Postma, 2005](#)).
269

270 **3.2 Data transformation using the CoDa approach**

271 Geochemical datasets are mostly composed by compositional variables . That is,
272 multivariate variables where the individual parts are parts of a whole ([Buccianti and](#)
273 [Grunsky, 2014](#)). Classical examples refer to constant sum variables, but recent definitions
274 of compositional data include all types of data representing parts of some whole. Ignoring
275 the compositional character of these geochemical variables may lead to misleading results
276 ([Pawlowsky-Glahn et al., 2015](#)). In this context, the CoDa methodology is used in this
277 work. In order to avoid the problems derived from the compositional data character, three
278 transformations, all based on log-ratios has been historically proposed, named as: additive
279 log-ratio (alr) transformation and centered log-ratio (clr) transformation ([Aitchison, 1986](#))
280 and isometric log-ratio (ilr) transformation ([Egozcue et al., 2003](#)).
281

282 In this study, the hydrochemical dataset was transformed using, firstly clr and secondly
283 ilr. The former transformation is described by
284

$$285 \text{clr}(\bar{x}_i) = \ln\left(\frac{x_i}{g(x_i)}\right); i = 1 \div D , \quad (1)$$

286 where $g(x) = \sqrt[D]{\prod_{i=1}^D x_i}$ is the geometric mean of all the considered parts (ions), and D
287 is the matrix dimension.
288

289 The ilr transformation allows to express hydrochemical compositions with respect to an
290 orthonormal basis. Their coordinates, called balances, may be easily obtained using a
291 Sequential binary partition (SBP) ([Egozcue et al., 2003; Egozcue and Pawlowsky-Glahn,](#)
292 [2005, 2006; Pawlowsky-Glahn et al. 2015](#)). The SBP has been widely used for many

293 authors on water chemistry studies (Engle and Rowan, 2013; Owen et al., 2016; Hee Kim
294 et al., 2019; Bondu et al., 2020). For a D column matrix, i.e. a D-part composition, D-1
295 balances are calculated from the SBP as

$$ilr(\bar{x}_i) = \sqrt{\frac{r_{i+} \cdot r_{i-}}{r_{i+} + r_{i-}}} \ln \frac{g(c_{i+})}{g(c_{i-})}; \quad i = 1 \div D - 1, \quad (2)$$

296
297 where c_{i+} and c_{i-} the groups of parts separated in the i^{th} step of the SBP; r_{i+} and r_{i-} are the
298 numbers of parts included in c_{i+} and c_{i-} , respectively.

299 There are different software tools that allow to perform these transformations. The called
300 CoDaPack v.2.0. program (Comas-Cufí and Thió-Henestrosa, 2011) is a software
301 developed by the Research Group in Statistics and Compositional Data Analysis at
302 University of Girona (UdG). This software can be freely downloaded from
303 <http://ima.udg.edu/codapack>. It allows performing the log-ratio transformations and to
304 prepare different kind of plots to show the results. In this research, all statistical analyses
305 were done using the statistics program R version 3.6.1 (2019-07-05) (R Development
306 Core Team 2004), which is available for free under the GNU-public License and for all
307 platforms from <http://www.cran.r-project.org>, through the software RStudio, a graphical
308 user interface for R . For multivariate statistical analysis (MSA) using the CoDa approach,
309 the following packages for R software were used: {stats} version 3.6.1. (R-core R-
310 core@R-project.org); {compositions} version 1.40-5 (Van den Boogaart and Tolosana-
311 Delgado, 2008) {zCompositions} version 1.3.4 (Palarea-Albaladejo and Martín-
312 Fernández, 2015).

313 Water samples with solute dissolved concentrations lower than the detection limit (the
314 so-called ‘left-censored values’) put an extra challenge when addressing MSA
315 techniques. The censored data can be either removed, or replaced or imputed (e.g. values
316 below detection limit are rounded as zeros) (Carranza, 2011). Following the criteria used
317 for several authors (Reimann and Filzmoser, 2000; Farnham et al., 2002), in this work,
318 left-censored values were excluded from the MSA when they represented > 25% of the
319 total number of samples (i.e. when the variable had a ‘medium–high’ level of nondetects
320 according to Palarea-Albaladejo and Martín-Fernández, 2014). Different algorithms can
321 be applied within the {zCompositions} package for R for imputing these values (like
322 multRepl, multLN, lrEM and lrDA methods).

323

324 **3.3 Principal Component Analysis (PCA) and Model-based clustering**

325 The first step to apply any MSA, is to check the presence of left-censored data and the
326 imputation of values. The function ‘zPatterns’ {zCompositions} is used to find and display
327 patterns of zeros/missing values in the whole dataset (see pattern diagrams at [Fig.SM.2.1](#)
328 [of Supp. Mat.](#)). In this work, the left-censored detected values were imputed using the
329 ‘lrDA’ (log ratio Data Argumentation) function. It is based on the log ratio Markov Chain
330 Monte Carlo Data Argumentation algorithm ([Palarea-Albaladejo and Martín-Fernández,](#)
331 [2015](#)), and it has been already used by different authors to delineate water types (e.g.
332 [Owen et al 2016; Hee Kim et al., 2019](#)). Following the commented procedure two data
333 matrices were prepared:

334

- 335 • Dataset Matrix (300x8), corresponding to 300 water samples (288 groundwater
336 samples and 12 snow and water ponds samples) and 8 variables (HCO₃, Cl, SO₄,
337 NO₃, Ca, Mg, Na, K).

338

- 339 • Dataset Matrix (43x8), corresponding to 43 springs and 8 variables (HCO₃, Cl,
340 SO₄, NO₃, Ca, Mg, Na, K), which represent the median hydrochemical
341 composition of groundwater evaluated for each spring ([Table SM.1](#). Supp. Mat.)
342 The consideration of “median composition” of time series follows the
343 requirements to estimate NBL’s using the PS method (see [section 3.4](#)).

344

345 [Table SM.3.1 \(Supp. Mat.\)](#) shows the list of parameters ‘included’ and 'excluded' for the
346 MSA and their justification.

347

348 PCA is a very common method that is based on dimensionality reduction of datasets. It
349 helps deciphering hydrogeochemical patterns and to infer the controlling variables of the
350 water chemistry ([Merchán et al., 2015; Moya et al., 2015](#)). In order to perform the PCA
351 it is necessary to calculate the ‘*variation matrix*’ of the dataset ([Aitchison, 1986](#)) as a first
352 step to obtain a measure of the dependence of the different variables, that is, the parts of
353 the composition. Each component of the variation matrix, τ_{ij} , describes the log-
354 relationship between two of the composition x_i and x_j (in this case chemical species), and
355 it's defined as

356

$$\tau_{ij} = \text{var}\left(\ln\frac{x_i}{x_j}\right) = \frac{1}{N-1} \sum_{n=1}^N \ln^2\left(\frac{x_{ni}}{x_{nj}}\right) - \ln^2\left(\frac{g_i}{g_j}\right), \quad (3)$$

357

358 where N is the number of observations and g_i , g_j are the geometric mean values for the
 359 two variables in question. A small value of τ_{ij} (which is equivalent to τ_{ji}) implies a good
 360 proportionality between the two variables. The variation matrix, τ_{ij} , is obtained using the
 361 R function ‘summary.acomp’ of the package {compositions}.

362

363 Once the variation matrix is obtained, then the correlation between the variables x_i and x_j
 364 is estimated through the ‘*index of proportionality*’ function, ρ_{ij} (Eq. 4) (Aitchison, 1986).
 365 The stronger the correlation between x_i and x_j the closer to 1 is the value of ρ_{ij} .

366

$$\rho_{ij} = \exp\left(\frac{-\tau_{ij}^2}{2}\right) \quad (4)$$

367

368 Data transformation following the CoDa approach is applied before using any MSA tool.
 369 In this case, PCA is applied using clr-transformed data (Eq. 1), with the function ‘clr’ of
 370 the {compositions} R package. The method provides a new matrix of standardized
 371 coordinates for each sample called ‘the scores’, and also a new matrix of variable
 372 ‘loadings’ with columns representing the principal components of the (clr-transformed)
 373 data.

374

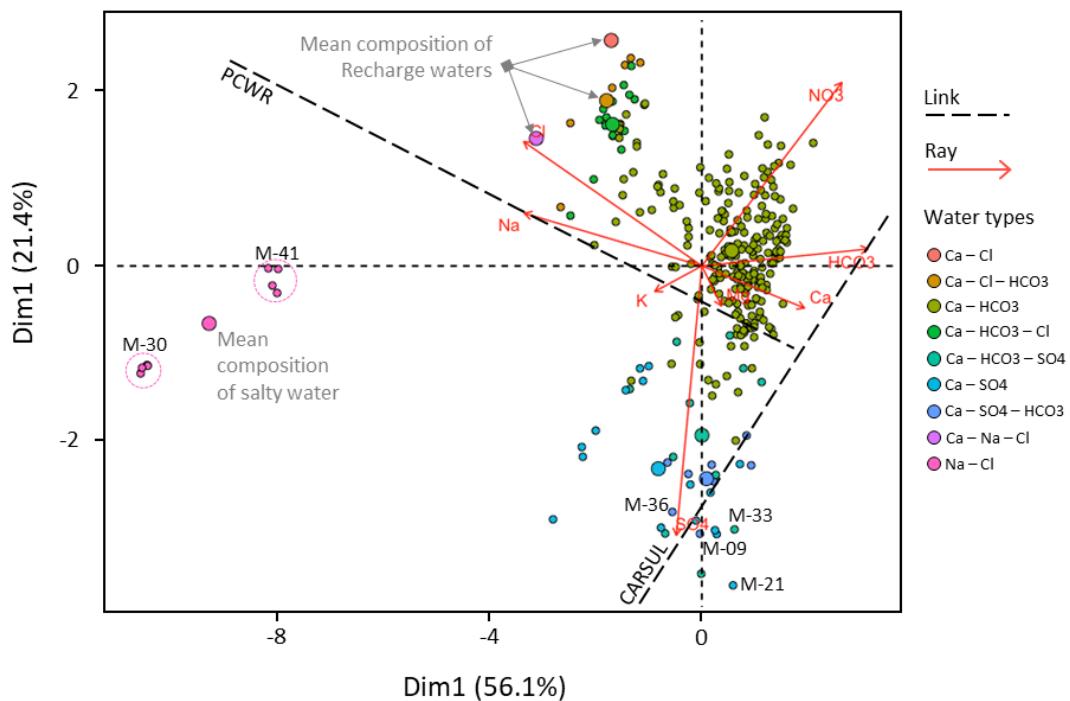
375 The graphical representations of the PCA results of clr-transformed data were done using
 376 the well-known biplot graphic (Gabriel, 1971) (Fig. 3), where the individuals are
 377 expressed as dots and the variables as rays. However, the interpretation of the clr-biplot
 378 differs from the interpretation of the classical biplot. Therefore, the clr-biplot
 379 interpretation is conducted by following the criteria proposed by Aitchison and Greenacre
 380 (2002), which is well suited for analyzing compositional data (Otero et al., 2005; Engle
 381 and Rowan, 2013; Blake et al., 2016; Piña et al., 2018). The criteria can be summarized
 382 as:

383

- 384 • The length of a link (i.e. black shaded line) between the rays (red arrows) defining
 385 $\text{clr}(x_i)$ and $\text{clr}(x_j)$ is proportional to the variance of $\ln(x_i/x_j)$.

- 386 • If two rays lay near each other, their quotient might be almost constant, and they
 387 might be proportional.
 388 • If two links between four different clr-variables are orthogonal, then the
 389 corresponding pairwise quotients may be independent.
 390 • If three or more vectors lie on the same link, the corresponding sub-composition
 391 might have one single degree of freedom.
 392 • If two links between four separate clr-variables are orthogonal then the
 393 corresponding pairs of variables may vary independently of each other.

394



395

Fig. 3. clr-Biplot of the Principal Components PC1 and PC2 for the dataset Matrix (300x8). The label of the axes indicates the percentage of the variance explained by PC1 and PC2, respectively. The PCWR dashed line indicates the link between Pristine waters and groundwater with Water-Rock interaction. The CARSUL dashed line indicates the link between CARbonate and SULphate waters. The smaller circles correspond to the different water samples and their color indicates their corresponding water type, whereas the larger circles represent the average composition of the different water types. To illustrate this the groundwater samples from springs M-30 and a M-41 are indicated as well as the corresponding mean composition.

405

406 The cluster analysis is applied to group observations into several homogeneous clusters.
 407 It is based upon similarities between the observations and provides insights regarding the
 408 multivariate geochemistry characteristics (Bondu et al., 2020; Templ et al., 2008). In this
 409 work it is used the ‘soft’ model-based clustering method. One of the main advantages is
 410 that it uses a probability-based approach. Therefore the obtained partition can be
 411 interpreted from a statistical point of view unlike the classical ‘hard’ - or heuristic-based
 412 - algorithms (k-means, hierarchical clustering, etc.) (Bouveyron and Brunet-Saumard,
 413 2014). The model-based clustering approach used was the finite mixtures of multivariate-
 414 normal or Gaussian distributions known as Gaussian Mixture Model (GMM), which is
 415 included as {Mclust} (Fraley and Raftery, 2002; Fraley et al. 2012; Scrucca et al., 2016),
 416 in R version: 5.4.6 (Raferty et al., 2020). It assumes that observed data come from a
 417 mixture of underlying probability distributions representative of two or more clusters.
 418

419 The GMM assumes that following probability distribution function (PDF)

$$f(x) = \sum_{k=1}^K \omega_k f_k(x|\mu_k, D_k) , \quad (5)$$

420
 421 where ω_k represents the weight or mixing proportion ($0 \leq \omega_k \leq 1$; $\sum_{k=1}^K \omega_k = 1$) or
 422 probability that an observation comes from the k^{th} mixture component, K is total number
 423 of components (i.e., groups or clusters), and f_k is the PDF of the observations for the k^{th}
 424 variable. Each component is usually modeled by a normal distribution (Eq. 6) with mean
 425 μ_k and covariance matrix D_k .

$$f_k(x|\mu_k, D_k) = \frac{1}{(2\pi)^{\frac{p}{2}} \cdot |D_k|^{\frac{1}{2}}} \exp \left[-\frac{(x-\mu_k)^T (x-\mu_k)}{2 \cdot D_k} \right] \quad (6)$$

426
 427 Taking into account Eq. 5 the conditional probability of assigning one observation to a
 428 given cluster is given by
 429

$$P(\text{cluster } k|x) = \frac{\omega_k f_k(x|\mu_k, D_k)}{f(x)} \quad (7)$$

430
 431 The greater the value of P the closer the association of sample x with the PDF
 432 corresponding to the cluster k. By definition, those samples for which $P > 0.5$ for PDF k
 433 constitute a “cluster”.

434

435 For the different components K , the model parameters ω_k , μ_k , and D_k are estimated using
436 the expectation–maximization (EM) algorithm (Dempster et al., 1977). The covariance
437 matrix D_k describes the geometry of the clusters with its volume, shape and orientation
438 The different combinantions of these parameters allows to define 14 multivariate mixture
439 models grouped in three main families: spherical, diagonal, and, ellipsoidal which are
440 included in the version used of {Mclust} package. In the other hand, this package uses
441 the Bayesian Information Criterion (BIC) to find the optimum number of clusters and
442 identify from those 14 multivariate mixing models, the one that best characterizes the data
443 while maximizing BIC. More details of the GMM, BIC and EM mathematical approach,
444 can be found on Biernacki and Govaert (1999); Fraley and Raftery (2002, 2012) and
445 Raferty et al. (2020). In this study, model-based clustering has been applied to the dataset
446 Matrix (43x8) of major ion data (HCO_3 , Cl, SO_4 , NO_3 , Ca, Mg, Na, K), represented in
447 this case using ilr-coordinates (Eq. 2).

448

449 **3.4. Determination of Natural Background Levels (NBLs) and Threshold**

450 **Values (TV)**

451 After identifying the number of underlying clusters in the data set in hand based on MSA
452 tools, the NBL and TV values for Cl, SO_4 and NO_3 are determined, which are the most
453 common solutes causing specific groundwater pollution issues in HMKS, are determined.
454 In this work the PS-method developed in the framework of the EU “BRIDGE” (2007)
455 project (Müller et al., 2006) is applied since it has been successfully proven in many
456 studies (Coetsiers et al., 2009; Ducci and Sellerino, 2012; Hinsby et al., 2008; Marandi
457 and Karro, 2008; Parrone et al., 2019; Preziosi et al., 2010; Wendland et al., 2008; Zabala
458 et al., 2016). The PS-method considers the following criteria for the data preparation
459 before estimating the NBL’s:

460

- 461 • Time series should be replaced by median averaging (i.e. all sampling sites
462 contribute equally to the NBL estimation).
- 463 • Samples with incorrect ion balance (exceeding 10%) and samples with median
464 NO_3 contents >10 mg/L must be rejected.
- 465 • Salt waters (i.e. NaCl) exceeding 1 g/L must not be considered.
- 466 • If samples are anaerobic ($\text{O}_2 < 1 \text{ mg/L}$) or denitrification occurs, the dataset needs
467 to be evaluated for the aerobic and anaerobic samples separately.

468

469 To obtain the NBL, the 90th percentile is advisable for small datasets ($N \leq 60$ sampling
470 points) or when human impact cannot be excluded from the data, which is the case of the
471 case study in this research. For $n > 60$ the 97.7th percentile is preferred. Once the NBLs
472 are defined then the TVs are obtained following the final methodology suggested by the
473 EU “BRIDGE” project:

474

$$TV = \begin{cases} \frac{1}{2} \cdot (NBL + Ref); & NBL \leq Ref \\ NBL; & NBL > Ref \end{cases}, \quad (9)$$

475

476 where *Ref* is the reference value. In case of the Spanish Royal Decree 140/2003 of 7
477 February, laying down the health criteria for the quality of water intended for human
478 consumption, the values of *Ref* for SO₄, NO₃ and Cl are 250 mg/L, 50 mg/L and 200
479 mg/L, respectively.

480

481

482 **4. Results and discussion**

483 **4.1. Exploratory analysis of data and general water chemistry**

484 In order to explore the internal structure of the dataset Matrix (43x8) of major ions (HCO₃,
485 Cl, SO₄, NO₃, Ca, Mg, Na, K), different EDA plots (a combination of histogram, density
486 trace, and one-dimensional scatterplot in one view) (Reimann et al., 2008), and boxplot
487 in just one plot, were prepared (Fig. 4). Having this in mind, the ilr coordinates are
488 adapted to the univariate case with the package {StatDa} (Filzmoser et al, 2009, 2009b).
489 The variable of interest x (i.e. in this case Cl, NO₃ and SO₄) is single ilr-transformed (Eq.
490 11):

491

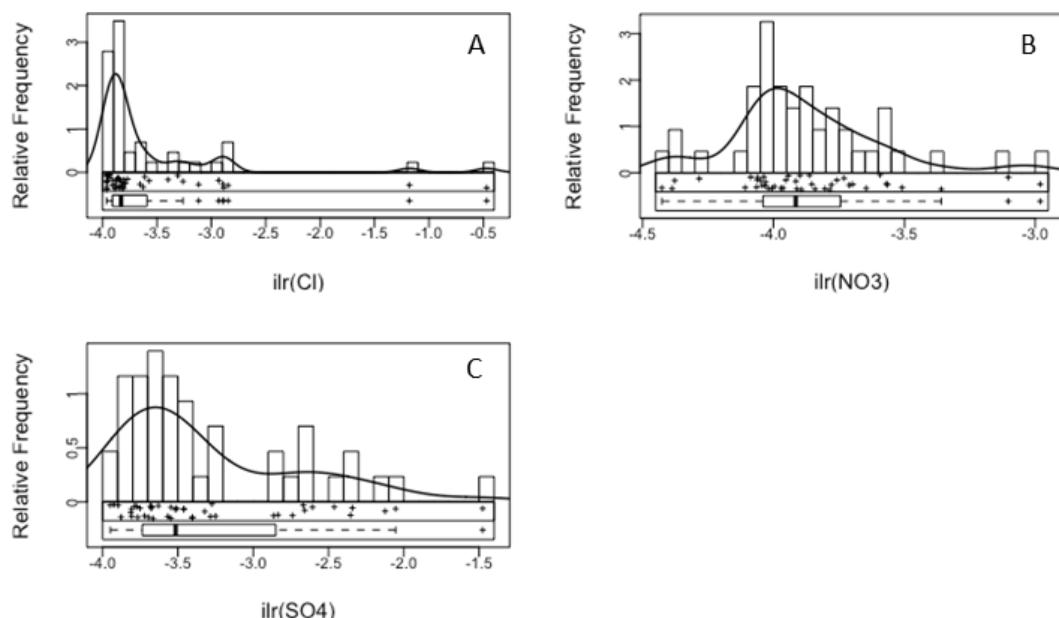
$$z = \frac{1}{\sqrt{2}} \cdot \ln \left(\frac{x}{1-x} \right) \quad (11)$$

492

493 The resulted histograms show multi-model shapes in all the cases (i.e. major ions)
494 suggesting that different populations are superimposed. In order to explain the dataset,
495 and considering the geological setting of the area, a hypothetical mixture model with
496 multiple components of different natural geogenic origin (possibly affected with local

anthropogenic sources) must be considered. Thus, a simply bi-modal distribution composed of natural vs anthropogenic contamination cannot be considered to establish the NBLs without taking into account the multivariate character of the data. Thus, the first step is to address separating the chemical groups or clusters.

501

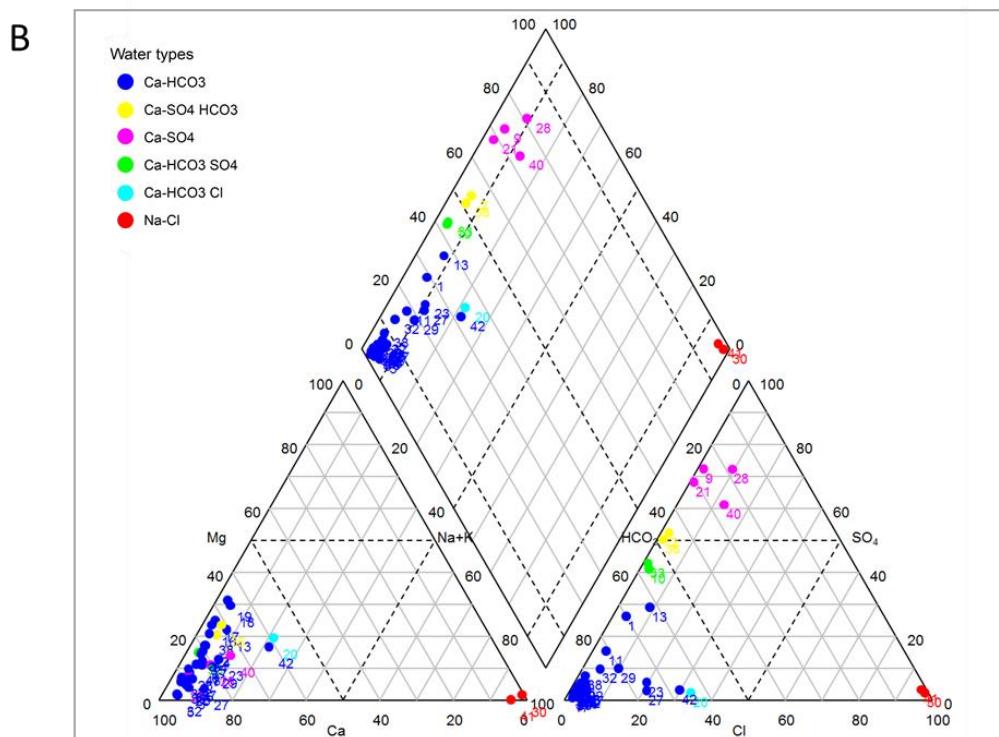
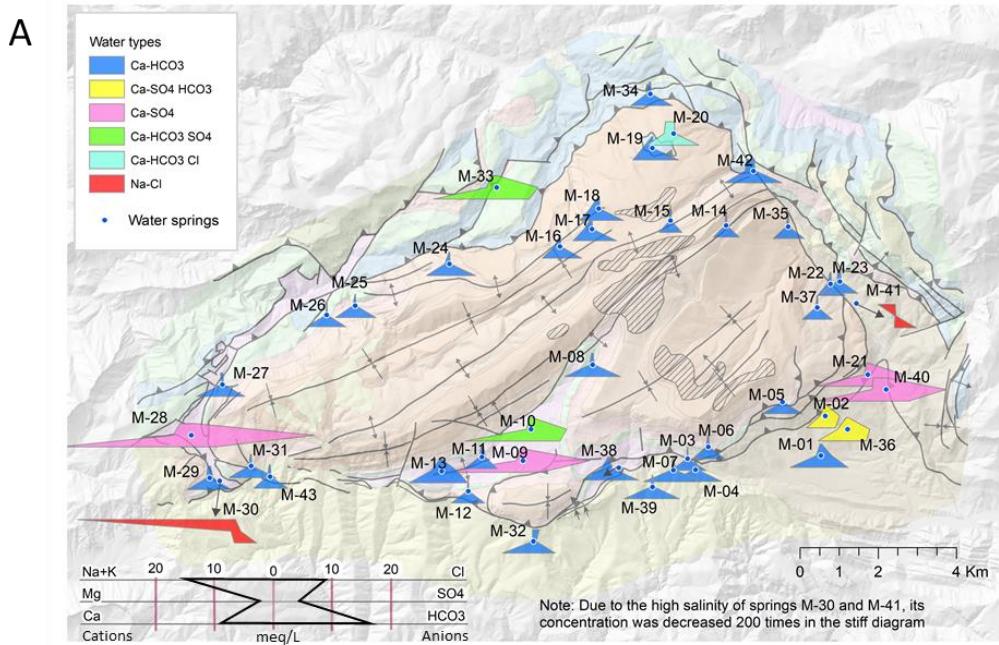


502

503 **Fig. 4** EDA plots for Cl (A), NO₃ (B), and SO₄ (C).

504 Classical graphical methods for the classification of water chemistry data, such as Piper
505 and modified Stiff diagrams were used as a first step to analyse the whole dataset (except
506 those water samples from pluviometers (i.e. in total 288 samples). The [Fig. 5](#) shows a
507 map with the Stiff diagrams distribution over the PCM and also the corresponding
508 modified Piper diagram. Based on that information it is possible to initially aggregate the
509 groundwater discharge from the 43 springs into 6 types of hydrogeochemical facies
510 ([Table 2](#)):

511



512

513 **Fig. 5.** Hydrochemical diagrams. (A) Stiff diagram map and (B) Piper diagram associated
 514 to the selected springs in the PCM. In both cases, for every spring the ion dissolved
 515 content values correspond to the median value associated to all samples taken from that
 516 spring. The springs are classified by their hydrochemical facies.

517

518 **Table 2.** Identified water types

Water type	Num. Springs	Geological units ^a
Ca-HCO ₃	32	Cretaceous (KMca, Kgp, Kat) Paleogene-Eocene (PEab, PEci, PEcp1, PEm1) Paleogene-Oligocene (POcgs, POmlg, PPEc) Quaternary (Qpe, Qt0, Qv1) Triassic-Jurassic (TJb, TJcd) Triassic Muschelkalk (Tm)
Ca-HCO ₃ -Cl	1	Paleogene-Eocene (PEcp2)
Ca-SO ₄ _	4	Quaternary (Qcoo) Triassic-Keuper (Tk)
Ca-HCO ₃ -SO ₄	2	Triassic-Keuper (Tk)
Na-Cl	2	Triassic-Keuper (Tk)
Ca-SO ₄ -HCO ₃	2	Paleogene-Eocene (Pemb)

(a) For a given Water type, the geological units based on [ICGC, \(2007\)](#) ordered by number of springs

519

520 At the first glance, the results show that diverse springs outcropping from different
 521 geological units ([ICGC, 2007](#)) show similar groundwater facies, or also the same facies
 522 can be shown from different points located at different geological units. In this context,
 523 these graphical techniques should not be considered determinant alone to discriminate
 524 between hydrochemical groups and therefore, their results should be considered only
 525 preliminary. [Table SM.1 \(Supp. Mat.\)](#) shows the summary of the major ions content of
 526 the 43 monitored springs (expressed as median values of time series for the period
 527 September 2013 – October 2015) and also the water facies associated to them.

528 **4.2. PCA and dataset matrix size**

529 The variation matrix for the dataset Matrix (300x8) ([Table 3](#)) shows strong correlations
 530 between different pairs of variables such as Ca and HCO₃, Na and Cl, and Mg and HCO₃.
 531 Besides, NO₃ shows a high correlation with Ca and HCO₃, whereas almost no correlation
 532 with SO₄. This result indicates that the most groundwater samples affected by nitrate
 533 pollution are those from the Eocene karst aquifer with a Ca-HCO₃ hydrochemical
 534 composition.

535

536 **Table 3.** The upper triangle over the main diagonal shows the ‘*index of proportionality*’
 537 ([Eq. 4](#)) of the dataset Matrix (300x8). The lower triangle over the main diagonal shows

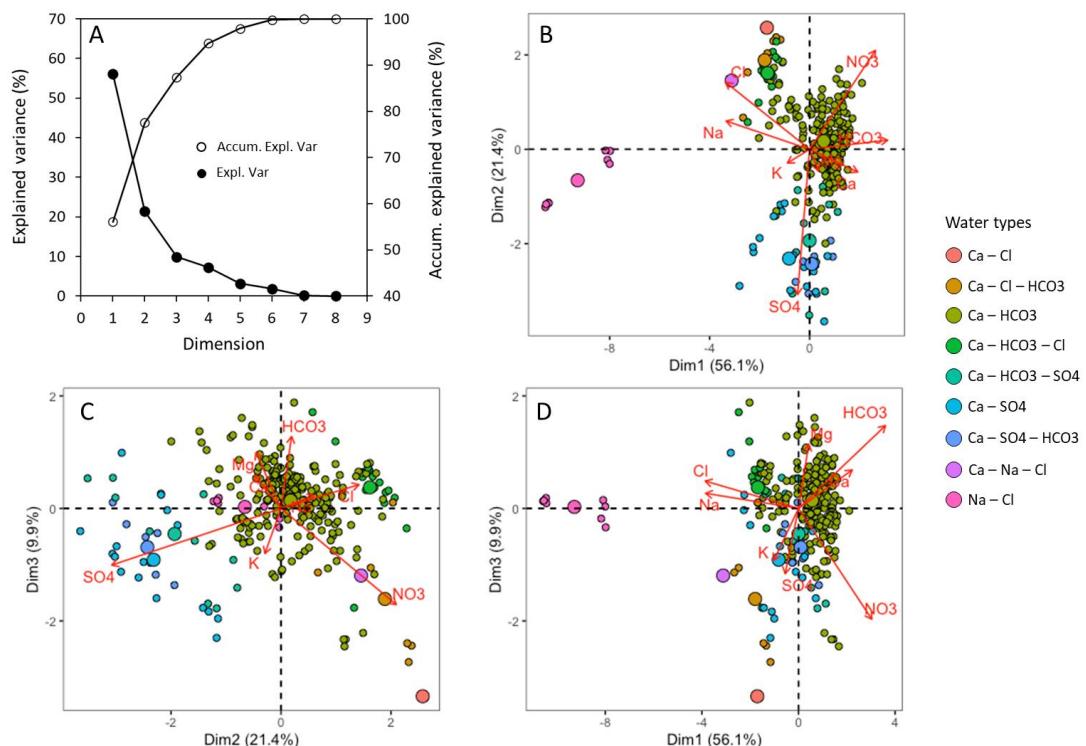
538 in italic the ‘*index of proportionality*’ of the dataset Matrix (43x8). In both cases, the
 539 correlation values larger than 0.5 are shaded in blue.

	Ca	Mg	Na	K	HCO ₃	Cl	NO ₃	SO ₄
Ca	--	0.88	0.06	0.51	0.98	0.04	0.57	0.44
Mg	0.76	--	0.34	0.69	0.67	0.26	0.27	0.52
Na	0.01	0.10	--	0.54	0	0.96	0	0.15
K	0.49	0.75	0.58	--	0.15	0.36	0.13	0.43
HCO ₃	0.94	0.41	0	0.07	--	0	0.56	0.06
Cl	0.01	0.07	0.99	0.40	0	--	0	0.05
NO ₃	0.55	0.07	0	0.02	0.59	0	--	0.01
SO ₄	0.17	0.24	0.07	0.25	0	0.02	0	--

540

541

542 Initially the PCA is conducted with the whole dataset (N=300), including the
 543 hydrochemical composition of natural and artificial snow, water from ponds and
 544 groundwater samples. The PCA with clr transformed data shows that only with three
 545 principal components, the 87.4 % of total variance can be explained (Fig. 6). The PCA is
 546 affected by the presence of natural outliers, in our case from the Na-Cl hydro-facies, that
 547 completely distorts the shape of the biplots (Fig. 6B, 6C and 6D). The scores are classified
 548 according to the distinguished nine water types when considering the complete dataset.
 549



550

551

552 **Fig. 6** (A) Scree-plot of dataset Matrix (300x8) showing the explained (solid circles)
553 variance associated to every PC of the PCA, and the accumulated explained variance
554 (empty circles) as the different PCs are accounted in the PCA. (B) Compositional biplot
555 PC1 vs PC2 (C) Compositional biplot PC2 vs PC3 and (D) Compositional biplot PC1 vs
556 PC3 showing scores (circles) and loadings (arrows) for clr transformed data. In the
557 biplots, the bigger points represent the mean clr-value for each water type.

558 From the distribution of the water samples in the clr-biplots several subgroups of waters
559 with clear similarities can be read. The biplot between PC2 and PC3 clearly separates
560 sulfate waters. Moreover, looking at the biplot between PC1 and PC2 closely (Fig. 3),
561 different hydrochemical spatial trends, likely associated with changes in terms of bedrock
562 lithology, can be observed. In fact, it can be inferred that: (1) The highest clr-variances
563 are shown for SO₄, Cl and NO₃, followed by Na and HCO₃. The lowest clr-variances are
564 shown for Ca, K, and Mg; (2) The PCA has situated the saltiest waters (M-30 and M-41)
565 separately in the western quadrant of the biplot. As can be shown, using clr-transformed
566 data allows to correctly separate characteristic points of the domain, which correspond to
567 the deepest drainage from the Keuper materials; (3) The groundwater samples from the
568 remaining springs are located in the eastern north and south quadrants: the freshest waters
569 that are more related to the upper Eocene karst aquifer are situated at the eastern-north
570 quadrant and present some correlation with NO₃. The samples related to Cretaceous and
571 Triassic materials appear to be more disperse, being most of them at the eastern-south
572 part of the biplot with extreme values in springs M-21, M-9, M-33, M-36, among others.

573 Taking into account the specific rules for interpreting clr-biplots, the following aspects
574 can be highlighted:

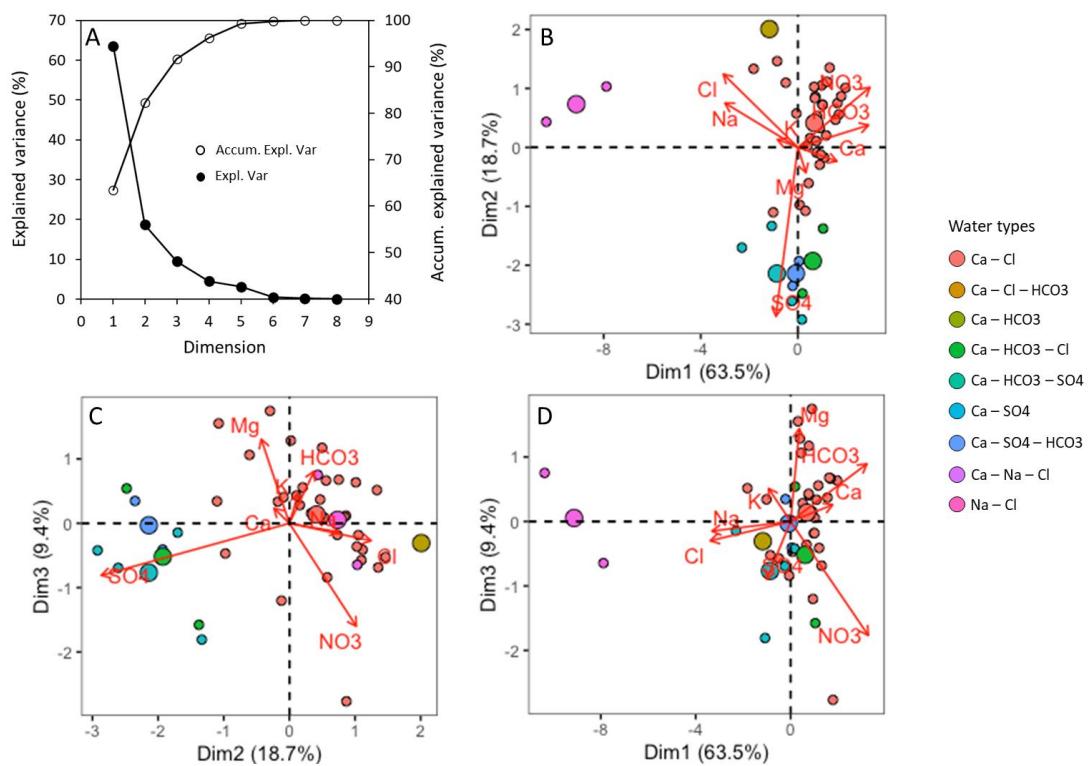
- 575 ▪ It is possible to draw a link between the vertices of Na, K and Mg indicating that
576 these variables may form a sub-composition with a single degree of freedom.
577 ▪ The vertices of SO₄, Ca and HCO₃ lie almost on a common link. This link is also
578 almost orthogonal to the link drawn between Na, K and Mg suggesting that these
579 two sub-compositions may vary independently of each other.
580 ▪ The two indicated links can be interpreted as two independent set of processes on
581 the hydrochemistry of the springs: (1) The “Pristine Character/Water-Rock
582 interaction” link PCWR [Na, K, Mg] which represents as one end-member, the
583 groundwaters influenced NaCl contributions derived from Keuper materials but

584 also to recharged waters ($\text{Ca}-\text{Cl}-\text{HCO}_3$; $\text{Ca}-\text{Cl}$, $\text{Ca}-\text{NaCl}$) at the upper part of the
 585 PCM, and representing the other end member the waters that have interacted
 586 longer with the Tertiary karst system materials. (2) The “CARbonate/SULfate
 587 dissolution” link ‘CARSUL’ [SO_4 , Ca , HCO_3] representing the dissolution of
 588 different types of carbonate and sulfate rocks (HCO_3 as one end member of the
 589 link and SO_4 as the other one).

- 590 ▪ Samples in the eastern-south quadrant of the biplot are more disperse having a
 591 stronger association with the SO_4 vertices.

592 In the case of the dataset Matrix (43x8) the variation matrix (Table 4) is consistent with
 593 that of the dataset Matrix (300x8), showing strong correlations between the same pairs of
 594 variables, and even with similar correlation values. The PCA with clr transformed data
 595 shows that when considering two or three PCs, it can be explained 87.4% and 91.7% of
 596 total variance, respectively (Fig. 7). Besides, the resulting clr-biplots are similar in shape
 597 to those of Matrix (300x8). As it can be shown, the reduction of the dataset matrices from
 598 (300x8) to (43x8) in the PCA does not introduce any relevant change in the final inference
 599 regarding the geochemical characteristics of groundwater. This is convenient from the
 600 perspective of dimensionality issues.

601



602

603 **Fig. 7** (A) Scree-plot of dataset Matrix (43x8) showing the explained variance (solid
 604 circles) associated to every PC of the PCA, and the accumulated explained variance
 605 (empty circles) as the different PCs are accounted for in the PCA. (B) Compositional
 606 biplot PC1 vs PC2 (C) Compositional biplot PC2 vs PC3 and (D) Compositional biplot
 607 PC1 vs PC3 showing scores (circles) and loadings (arrows) for clr transformed data. In
 608 the biplots, the bigger points represent the mean value for each water type.

609

610 4.3 Clustering analysis

611 The principal aim of cluster analysis is to split a number of observations into groups that
 612 are similar in their characteristics or behaviour ([Reimann et al. 2008](#)). To this end, the
 613 dataset Matrix (43x8) is used. Before conducting the ilr transformation, an intuitive
 614 sequential binary partition (SBP) is used to characterize the hydrochemical variability
 615 within the domain wherein non-overlapping groups of parts, known as balances, are
 616 defined. According to [Egozcue and Pawlowsky-Glahn \(2005; 2006\)](#), two methods for
 617 performing SBP can be applied: (1) directly a PCA, and (2) by experienced judgment. In
 618 this work SBP is based on knowledge of the groundwater chemistry in the study area and
 619 on the resulting compositional biplot ([Fig 7](#)). As a result, seven groundwater partitions
 620 are considered ([Table 4](#)): the ilr_1 balance separates the Ca-HCO₃ waters (mostly affected
 621 by NO₃) from the rest; the ilr_2 separates those waters affected/non-affected by NO₃
 622 pollution; the ilr_3 separates the contribution of calcite and dolomite to groundwater; the
 623 ilr_4 separates Ca from HCO₃; the ilr_5 separates SO₄ waters from most salty waters; the
 624 ilr_6 separates K from Na/Cl; and finally the ilr_7 separates Na and Cl.

625 **Table 4.** SBP of a 7-part composition (ilr_1, ilr_2, ..., ilr_7) for describing isometric log
 626 ratio (ilr) coordinates based on the separation of anions and cations related to the
 627 hydrochemical composition of natural groundwaters for the clustering analysis.

ilr	Ca ²⁺	Mg ²⁺	Na ⁺	K ⁺	HCO ₃ ⁻	Cl ⁻	NO ₃ ⁻	SO ₄ ²⁻
ilr_1	+1	+1	-1	-1	+1	-1	+1	-1
ilr_2	+1	+1	0	0	+1	0	-1	0
ilr_3	+1	-1	0	0	+1	0	0	0
ilr_4	+1	0	0	0	-1	0	0	0
ilr_5	0	0	+1	+1	0	+1	0	-1
ilr_6	0	0	+1	-1	0	+1	0	0
ilr_7	0	0	+1	0	0	-1	0	0

628

629 In order to inspect the sensibility of considering ‘hard’ clustering methods to determine
630 the optimal number of clusters (k), a first and preliminary analysis was performed using
631 the {clValid} (Brock et al. 2008) and {NbClust} (Charrad et al. 2014) R packages and
632 using clr and ilr coordinates. These packages contain a set of functions for both validating
633 the results of different heuristic-based clustering algorithms (e.g., ‘*hierarchical*’, ‘*k-*
634 *means*’, ‘*diana*’, ‘*pam*’ and ‘*clara (k-medoids)*’), and determining the relevant number of
635 clusters in a dataset. The clustering models may account for different linkage methods
636 (i.e., ‘*complete*’, ‘*average*’, ‘*single*’ and ‘*ward*’) and dissimilarity metrics (‘*Euclidean*’
637 and ‘*Manhattan*’, among others). In the case of the dataset Matrix (43x8), the best results
638 have been obtained for the *hierarchical* and *k-medoids* models, and for k clusters between
639 2 and 4. Nevertheless, given the large number of methods and their corresponding options
640 for setting the models it is actually difficult to select the most suitable one.

641 To avoid the clustering method selection issues, the cluster analysis was performed by
642 solving a GMM using the {Mclust} R package (Fraley and Raftery, 2002; Fraley et al.,
643 2012; Scrucca et al., 2016) and considering the ilr-transformed data. The obtained results
644 indicates that in the case of the dataset Matrix (43x8), the best multivariate clustering
645 option is obtained applying the ‘EEI’ model while considering a total of 4 clusters (see
646 Fig. SM.4.1 in Suppl. Mat.).

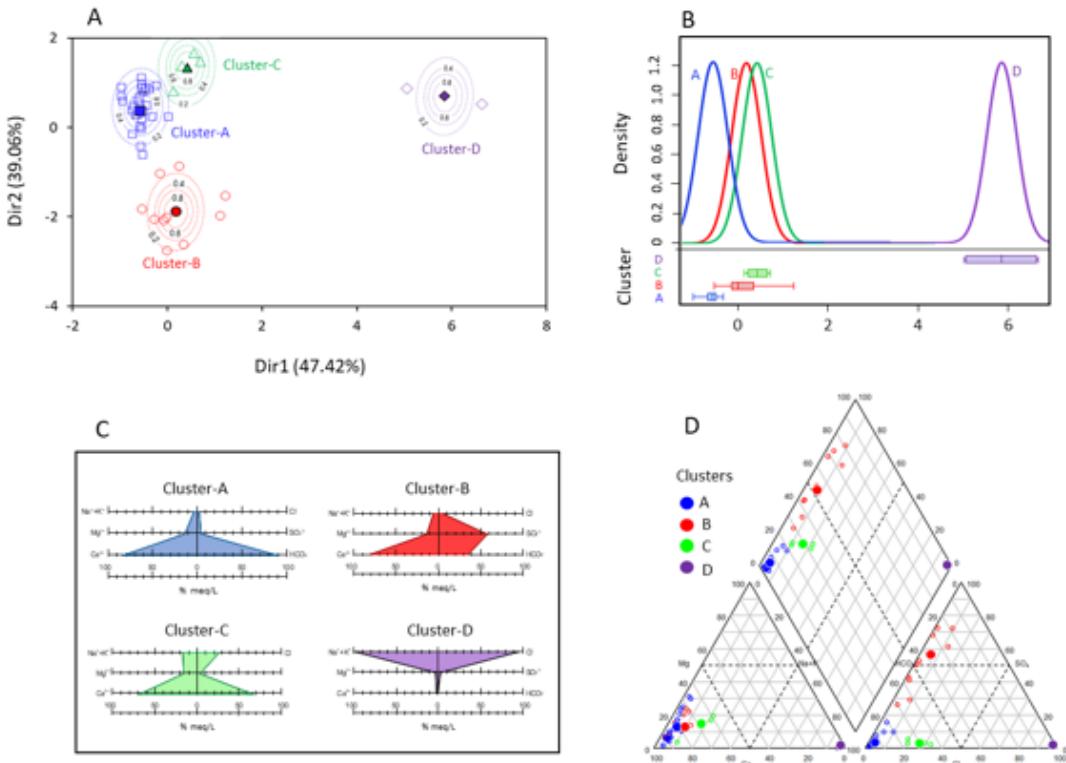
647 The scatterplot matrix obtained with the model-based clustering process using the seven
648 ilr coordinates. Being D the previous dimension of the original dataset matrix (43x8), $D - 1$
649 coordinates can be shown in Fig. SM.4.2 (Suppl. Mat.). In order to visualize the clusters
650 in a most suitable way, the dimension reduction function ‘MclustDR’ (Scrucca, L. 2010)
651 for visualizing the classification structure obtained from the finite mixture of Gaussian
652 densities of the {Mclust} package is used to reduce the dimensionality of the ilr matrix
653 and estimate the principal components. Table. SM.4.1 and Fig. SM.4.3 (Suppl. Mat.)
654 provide the scores of the reduced ilr-matrix and their representation in a scatterplot,
655 respectively. As can be shown, the two main principal components explain 86,42% of the
656 total variance. As a result, and only a glance at the scatterplot of PC1 and PC2 (Fig. 8A)
657 the cluster division for the different springs reveals in a clear way, and each cluster can
658 be described by the corresponding PDFs (Fig. 8B). Worth to point out the similarity
659 between the distributions of samples in the 2D space (albeit in a symmetric plane). The
660 Fig. 8C presents the mean hydrochemical composition of each cluster (Table SM.4.2 in

661 Suppl. Mat.) in terms of modified Stiff diagrams, and Fig. 8D shows in a Piper diagram
662 how the mean hydrochemical composition of the clusters is representative of the
663 composition of the corresponding springs.

664 The probabilistic GMM framework estimates the optimal number of clusters and provides
665 for every spring the probability of belonging to these clusters (soft assignment). This
666 approach is more interesting than the classical clustering approaches, in which the number
667 of clusters is assumed fixed, and every spring is assigned to one and only one of the
668 previously assumed clusters (hard assignment) (Kim et al., 2014). From an hydrochemical
669 point of view, the soft assignment often provides the more interesting interpretation
670 because the method reveals if one observation is influenced by several factors (Templ et
671 al., 2008). Moreover, Wu et al., (2017) show how the probabilistic GMM clustering
672 provides insights into hydrochemical processes affecting groundwater, even though with
673 a limited number of observations, which is a common situation in high mountain karst
674 aquifers such as PCM.

675

676 The conditional probabilities (P) of assigning one observation to a given cluster (Eq. 7)
677 are given in Table SM.4.3 (Suppl. Mat.). In all cases, springs are assigned to one cluster
678 with a probability > 0.95 , and specifically, more than 83% of the springs reach the
679 probability of ‘1’. The smaller probabilities occur in M-01 ($P = 0.911$ cluster A) and M-
680 13 ($P = 0.969$ cluster B). Spring M-01 discharges from the Eocene karstic limestones.
681 Nevertheless, this discharge might be affected by weak contributions of Tertiary sulfates
682 (which are related to the locally known as ‘Beuda gypsum Formation’). The discharge in
683 M-13 shows a Ca-HCO₃ hydrogeochemical composition despite discharging from the
684 Triassic (Muschelkalk) limestone aquifer. In this case the groundwater discharge is
685 weakly affected by the underlying Keuper materials.



686 **Fig. 8.** (A) clr-Biplot PC1 vs PC2 with the distribution of each cluster after dimension
687 reduction for dataset Matrix (43x8) ilr-transformed data. The dashed lines correspond to
688 the probability zones of belonging a certain cluster in the subspace PC1-PC2. Solid
689 symbols correspond to the mean hydrochemical composition of the clusters (B) PDF's of
690 the resulting 4 clusters. (C) Modified Stiff diagram associated to the mean hydrochemical
691 composition of the clusters. (D) Piper diagram associated to the selected springs in the
692 PCM classified by their corresponding cluster to which they belong. Solid symbols
693 correspond to the mean hydrochemical composition of the clusters.

694 The hydrogeochemical description of each groundwater cluster can be summarized as:

- 695 • **Cluster A** is characterized by low mineralization and dominated by Ca-HCO₃
696 water type and slightly alkaline pH. In total 27 springs are grouped in this cluster
697 which correspond to 203 groundwater samples collected in the study from the
698 total of 288. All the springs drain directly or indirectly (i.e. covered by local
699 quaternary deposits) the Tertiary Eocene upper karst aquifer of the PCM ([Fig. 9](#))
700 and from the higher parts of the mountain (944 - 2144 m a.s.l.). They are almost
701 mainly found inside the structural limits of the PCM sheet and at its boundaries
702 except some of them localized in Quaternary deposits or discharging karstic

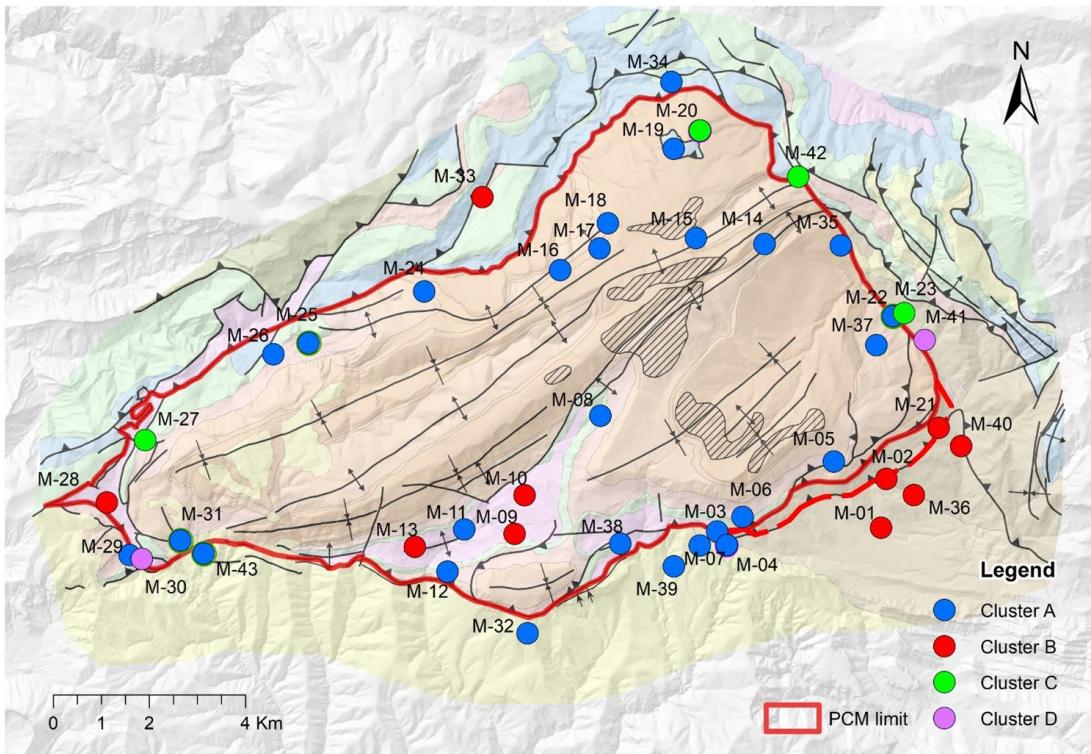
703 conduits through the Oligocene carbonate karstic conglomerates situated just in the
704 front of the thrust sheet (e.g. M-03, M-04, M-07, M-39, M-32 and M-43, which is
705 one of the most important karst springs of the system). Another special case is the
706 spring M-06 which lies over garumian shales, marls and limestones (Kgp)
707 outcropping materials. In this zone, a fault affecting the stratigraphy might allow
708 the hydrological connection between the lower Eocene limestones (PPEc) and
709 Kgp formations. This connection would explain the CA–HCO₃ water type
710 associated to spring M-06, and also its classification in the cluster A, thus pointing
711 the groundwater discharge origin as the Eocene Tertiary aquifer. Finally, the
712 spring M-29 actually drains a Eocene limestone level situated at the west of the
713 PCM boundary.

714 Cluster A presents the lower values of EC, which ranges between 186 and 486
715 $\mu\text{S}/\text{cm}$ and the minimum values of groundwater temperatures. The concentrations
716 of Cl and SO₄ are very low ranging between 2.5 and 15 mg/L and between 2.6
717 and 25.3 mg/L respectively. In 13 samples, the concentration of NO₃ is above 10
718 mg/L, and in one specific spring (M-32) exceeds in all samples the legal limit for
719 potable water (50 mg/L). The average Saturation Indices (SI) estimated with the
720 Phreeqc program ([Parkhurst and Appelo, 2013](#)) for calcite, gypsum and halite are
721 0.23, -2.67 and -9.68, respectively. The groundwaters are representative of the
722 recharge of the karst system in the highest altitudes of the massif, where the
723 dissolution of carbonates is the dominant geochemical process controlling
724 groundwater chemistry.

- 725 • **Cluster B** encompasses water types from Ca–HCO₃ to Ca–HCO₃–SO₄, Ca–SO₄–
726 HCO₃ and Ca–SO₄, which are characterized by moderate mineralization and
727 slightly alkaline pH. This group includes 10 springs. A total of 40 groundwater
728 samples collected in the study would correspond to this cluster. The springs
729 related to Cluster B are situated either inside or outside the internal structural
730 limits of the PCM thrust sheet. The springs situated inside (M-9, M-10 and M-13)
731 occur mostly in (1) Cretaceous and Triassic (Keuper) materials outcropping in the
732 area. These materials underly the principal aquifer of the massif (the Eocene
733 carbonate karstic system), and (2) local shallow granular aquifers. The springs M–
734 01, M-02, M-21 and M-36 are related to sediments with high content of Tertiary

735 gypsum from the Beuda Formation, which outcrops in small pinched belts located
736 in front of the southeastern part of the PCM thrust sheet. Springs are located at the
737 lowest parts of massif (altitudes ranging between 867 - 1456 m a.s.l.). The EC
738 varies between 493 and 2102 μ S/cm. The SO_4 concentration is quite high and it
739 ranges between 88 and 989 mg/L exceeding in most cases the legal limit for
740 potable water (250 mg/L). The concentration of Cl ranges between 3.8 and 94.5
741 mg/L. The average SI for calcite, gypsum and halite are 0.32, -0.99 and - 8.61
742 respectively.

- 743 • **Cluster C** includes water types from Ca–HCO₃ and Ca–HCO₃-Cl water types.
744 This group includes 4 springs and a total of 37 groundwater samples from which
745 26 of them correspond to the spring M-20 (located at 1858m a.s.l.). Except the
746 spring M-20, the rest (M-23, M-27, M-42) are located at the boundaries of the
747 PCM sheet. The EC varies between 332 and 747 μ S/cm. Although they have SO_4
748 concentration similar to cluster A with 9.7- 15.3 mg/L, the content of Cl is much
749 higher ranging between 24 and 82 mg/L. These higher values compared to cluster
750 A are interpreted as related with groundwater flow through areas with the presence
751 of relict halite or salty water in closed pores in the Keuper materials. In the case
752 of M-20 (which is located inside the PCM sheet) the salt is related to a klippe of
753 Jurassic delineated into the geological map. Besides, in the catchment area of this
754 spring, there are small outcrops of Keuper materials detected during the fieldwork.
755 The average SI for calcite, gypsum and halite are 0.24, -2.32 and -7.42
756 respectively.
- 757 • **Group D** is the most evident and special. corresponding to Na–Cl type waters
758 ([Fig. 8](#)). This group is composed of 2 salt springs (M-41 and M-30) located at the
759 993 and 1023 m a.s.l. at the East and West boundaries of the PCM sheet
760 respectively. They are characterized by very high mineralization discharging from
761 Keuper confined bedrocks and interpreted as the contribution of deep groundwater
762 flow with elevated transit times that allows a significant solute diffusion. The
763 waters are slightly acidic to near-neutral pH. The M-41 and M-30 samples
764 presents EC values of 57.2 and 247.1 mS/cm, Cl concentrations of 21 and 178.2
765 g/L, and SO_4 concentrations of 1.2 and 8.1 g/L respectively. The M-30 spring can
766 currently be considered the saltiest spring of natural origin in Catalonia.



767

768

769 **Fig. 9.** Spatial distributions of the 43 clustered springs over the geological map of the
 770 PCM based on the GMM. The description of the different geological materials is the same
 771 presented in Fig. 2.

772

773 In the framework of multivariate statistics data analysis (e.g. PCA and data clustering),
 774 specially dealing with compositional data (i.e. data that carry only information about the
 775 relative abundance of each component on the whole, such as the hydrogeochemical data
 776 sets), it is important to suitably transform the dataset using the CoDa approach (e.g., Eq.1
 777 or Eq.2) before conducting any analysis, otherwise it is very likely to obtain wrong results
 778 ([Otero et al., 2005](#)). Moreover, uninterpretable results are also obtained when applying
 779 the classical standardization methodology known as “z-score” on compositional data,
 780 which considers logarithms and then subtracts the mean and divides it by the standard
 781 deviation to scale ([Blake et al., 2016](#)). To illustrate the importance of using correct CoDA
 782 data transformations, the dataset Matrix (43x8) is used to apply the same MSA analysis
 783 (PCA and the model-based clustering GMM) but using the classical standardization
 784 approach (or Z-score normalization) instead of using the CoDa transformation approach.

785 As can be shown (see at the Annex.2) neither the PCA variable loadings nor the soft
786 clustering results make any clear hydrogeological sense.

787

788 **4.4 NBLs and TVs values.**

789 Once the groundwater clusters are defined for the PCM, the NBL and TV's for NO₃, SO₄
790 and Cl have been obtained applying the PS-method ([Müller et al., 2006](#)). Taking into
791 account the criteria required for data to be accounted when estimating the NBL's with
792 this method ([section 3.4](#)), the following observations apply:

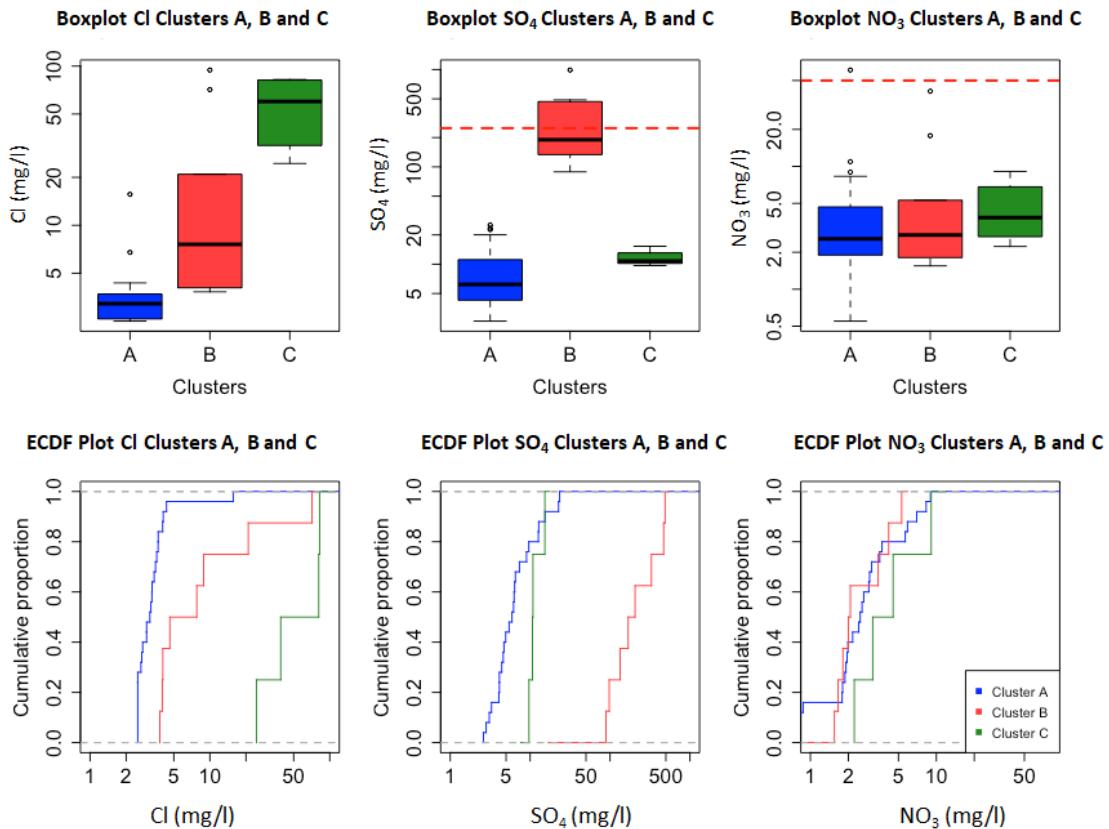
793

- 794 ▪ The groundwater samples from M-30 and M-41 (the whole cluster D) present Cl
795 concentrations of geogenic origin above the drinking water limit (>200 mg/L).
796 Therefore, these samples are not considered in the NBL determination.
- 797 ▪ NO₃ concentrations above the drinking water limit (>50 mg/L) are mostly
798 observed in M-32 spring (cluster A). Besides, the springs M-32, M-10, M-11 and
799 M-28 present NO₃ concentrations > 10 mg/L. Following the PS-method criteria,
800 these springs have been excluded of the NBL determination.

801

802 The NBLs for the remaining groundwater samples belonging to the clusters A, B and C
803 are obtained taking into account the 90th percentiles (P90) of the corresponding cluster
804 ECDF plots ([Fig. 10B](#)). The obtained NBL90 and TVs are presented in [Table 5](#). The
805 results indicate that Tertiary Eocene karst aquifer (cluster A), which is the principal
806 aquifer inside the PCM presents the lowest NBL90 values for Cl, SO₄ and NO₃. The
807 cluster B, which is related to the aquifers on the Cretaceous and specially the Triassic
808 Keuper materials presents the highest NBL90 value for SO₄, and the Cluster C, which is
809 generally related to local small aquifers located at the boundaries of the PCM presents the
810 highest values of NBL90 for both Cl and NO₃.

811



812

813 **Fig. 10.** (A) boxplots of the clusters A, B and C for SO₄, Cl and NO₃. The dashed red
 814 lines indicate the reference limits established in the Spanish Royal Decree 140/2003 (B)
 815 ECDF plots.

816

817 Comparing the obtained NBL90 values with those officially assigned to GWB-5 and
 818 GWB-44 (Table 1), it looks that the NBL official values of SO₄ assigned to both GWBs
 819 (485 and 609 mg/L, respectively) are likely conditioned by the interaction between fresh
 820 groundwater and most probably evaporites of the Upper Triassic (Keuper facies). These
 821 evaporites appear very often at the boundaries of many thrust sheets throughout the
 822 Southern Pyrenean zone. Additionally, the official NBL90 value of Cl assigned to GWB-
 823 44 is similar to that obtained for cluster B, which is related to the Keuper deposits.
 824 Likewise, the obtained NBL90 value of NO₃ for the Cluster C is similar to the official one
 825 for GWB-44. As can be shown, none of the official NBL90 values defined for GWB-5
 826 and GWB-44 correspond to those values obtained for the Lower Eocene limestones and
 827 dolomites, which constitute by large the main aquifer of the PCM.

828

829 **Table 5.** Summary results of the NBL and TV's values derived from the PS-method
 830 (BRIDGE, 2007) for clusters A, B and C for the solutes Cl, SO₄ and NO₃.

Cl [mg/L]			SO ₄ [mg/L]		NO ₃ [mg/L]	
Clusters	NBL _{90%}	TVs	NBL _{90%}	TVs	NBL _{90%}	TVs
A	4.06 ± 2	8.12 ± 2	14.33 ± 2	29.66 ± 2	6.55 ± 2	13.1 ± 2
B	35.98 ± 2	71.96 ± 2	471.71 ± 2	471.71 ± 2	4.51 ± 2	9.02 ± 2
C	81.92 ± 2	140.96 ± 2	13.96 ± 2	27.92 ± 2	7.73 ± 2	15.46 ± 2

831

832

833 It is well known that high mountain karst aquifers generate highly valuable water
 834 resources for the downstream water depending ecosystems. Their protection and rational
 835 management is of utmost importance to sustain such ecosystems and satisfying their water
 836 demands ([Kazakis et al., 2018](#)). In this framework, NBLs provide an objective scale to
 837 compare with when the quality status of the aquifer is assessed. Nevertheless, these
 838 aquifers are often immersed in complex geological settings as happens in the axial zone
 839 of the Central Pyrenees ([Lambán et al., 2015](#)), in the Picos de Europa massif ([Ballesteros](#)
 840 [et al., 2015](#)), in the Jura Mountains ([Luetscher and Perrin, 2005](#)) and the Hochifen–
 841 Gottesacker Alps ([Goldscheider, 2005](#)), among others. The NBLs are obtained as a
 842 function of the hydrochemical content measured in the different springs discharging the
 843 system. Nevertheless, in geological complex zones it is difficult to assert if one certain
 844 spring is discharging groundwater from the aquifer of interest or not, because the
 845 geographical location of the spring may suggest an origin for the sampled groundwater
 846 while hiding mixing relations between groundwater flow lines from other local aquifers
 847 with different hydrogeochemical fingerprint ([Lambán et al., 2015](#); [Barbieri et al., 2017](#);
 848 [Sánchez et al., 2017](#)).

849

850 The European Union Water Framework Directive ([WFD, 2000](#)) defines a general
 851 framework for integrated river basin management in Europe to ensure their “good water
 852 status”. Nevertheless, the river basin is often a entity hard to manage because the larger
 853 the size of the basin the larger is (1) the number of water bodies enclosed and (2) the
 854 likelihood of political-administrative boundaries issues to appear. To avoid such
 855 problems, instead of looking at river basins the WFD refocussed on the smaller scale
 856 “river basin districts”, for which administrative structures were defined to correctly
 857 manage the corresponding bodies, thus ensuring -hopefully- the right management of
 858 whole river basin ([Boeuf and Fritsch, 2016](#)). In this line, the WFD includes the guidelines
 859 that apply to define the groundwater bodies (GWB). Even in this case, some scale issues

860 may arise when considering the definition of the GWB in mountain zones. By definition,
861 the GWB are assumed to belong to a certain river basin. Despite of that, it is well known
862 that groundwater basins, specially in mountain zones, may extend throughout several
863 river basins ([Struckmeier et al., 2006; Serianz et al., 2020](#)). As a result, GWBs may
864 include several aquifers or even only parts of them as it happens in the PCM, whose
865 discharge contributes to both the Ebro and the Llobregat rivers through GWB-44 and
866 GWB-5, respectively. This is the reason why there are two different sets of NBL applying
867 for the same aquifer ([Table 1](#)).

868

869 The WFD recognises the importance of having well defined NBLs. Their characterization
870 must be based on a consistent and rigorous hydrochemical criteria, given that they will be
871 used to quantitatively assess whether or not anthropogenic pollution is taking place in the
872 corresponding aquifer ([Nieto et al., 2005](#)). The hydrogeological fingerprint of each
873 aquifer belonging to the same GWB may be different. Therefore, the criterion of defining
874 a single set of NBLs for the whole GWB may have no sense. Moreover, such criterium
875 may be counterproductive from a safety perspective, given that one may assume for the
876 GWB some quantities of species or chemical substances present in solution as normal,
877 when actually those concentrations may be already indicating the existence of a polluting
878 issue in some aquifers of the GWB. This is even worst when only one of these aquifers
879 play a relevant role from a water resources perspective, as happens in the PCM. Here, the
880 Lower Eocene karst aquifer generates an overall mean groundwater discharge that
881 represents 15% of the mean annual water consumption in the city of Barcelona ([Herms et](#)
882 [al., 2019](#)). Therefore, from a water resources management perspective, it might worth
883 defining NBLs at the local scale for each aquifer. In this line, the methodology presented
884 in this work to “complement” the sample pre-selection method is a useful tool to
885 objectively reel off the NBL of the different high mountain aquifers belonging to a given
886 GWB. Besides, the proposed methodology provides the GWBs managing authorities a
887 full-sense hydrochemical criteria to better protect the high mountain pristine and strategic
888 aquifers, while ensuring the good status of the associated high mountain river basins.

889

890

891 **6. Conclusions**

892

893 The PCM is a complex hydrogeological system composed by a main Eocene karst aquifer
894 that drives the hydrodynamical discharge response of the massif. The PCM also includes
895 small aquifers whose discharge present a different hydrochemical composition. The
896 discrepancies between the official NBLs of the GWBs associated to the PCM reveal the
897 disparities in the hydrochemical composition of groundwater from the different sampled
898 springs belonging the GWBs. To estimate correctly the NBLs associated to one aquifer it
899 is necessary to consider only samples from springs discharging groundwater from the
900 aquifer of interest. In high complex hydrogeological settings this selection is not easy and
901 must be guided by a consistent and objective clustering method.

902

903 In the case of the PCM, four compositional groups have been identified by means of
904 GMM clustering analysis. Most of the analysed springs are dominated by Ca–HCO₃ water
905 type coming from the main aquifer of the area. There are some springs dominated by Ca–
906 HCO₃, Ca–HCO₃–SO₄, Ca–SO₄–HCO₃, Ca–SO₄, Ca–HCO₃–Cl, Na–Cl water types
907 derived from other small/local aquifers. Determination of NBLs values in the area must
908 take into account the four groups defined in this study.

909

910 In complex aquifer systems the proposed soft clustering approach, which is based on
911 probabilistic Gaussian mixture models, provides the optimal number of clusters for the
912 sampled springs only based upon the observed compositional data while estimating the
913 probability of belonging to everyone of these clusters for each spring. The presented
914 clustering approach relies on multivariate statistics methods. In this framework it is
915 essential to transform the dataset using the CoDa rules, specially when dealing with
916 hydrochemical compositions. Otherwise, uninterpretable results will be likely obtained.

917

918 In the case of different existing aquifers with discrepant hydrochemical fingerprints in the
919 same GWBs, it would be reasonable to evaluate the NBLs in all of them rather than having
920 a single set of NBLs for the whole GWB. Otherwise, errors may appear when estimating
921 the quality status of some of these aquifers, even if the overall assessed quality status of
922 the GWB appears to be correct.

923

924

925 **Acknowledgements:**

926 This research has been supported by Agencia Estatal de Investigación (AEI) from the
927 Spanish Government and the European Regional Development Fund (FEDER) from EU
928 through PACE-ISOTEC (CGL2017-87216-C4-1-R) projects, the EFA210/16/
929 PIRAGUA project which is co-founded by the European Regional Development Fund
930 (ERDF) through the Interreg V Spain-France-Andorre Programme (POCTEFA 2014-
931 2020) of the European Union, the Catalan Government projects to support consolidated
932 research groups MAG (Mineralogia Aplicada, Geoquimica i Geomicrobiologia,
933 2017SGR-1733) from Universitat de Barcelona (UB) and GREM (Grup de Recerca de
934 Minería Sostenible) from the Universitat Politècnica de Catalunya (UPC), the Ministerio
935 de Ciencia, Innovación y Universidades through the METHods for COnpositional
936 analysis of DAta (CODAMET) project (Ref: RTI2018-095518-B-C22, 2019-2021). We
937 thank the Hydrogeology and Geothermics Unit Team of the Institut Cartogràfic i
938 Geològic de Catalunya (ICGC) by their helpful collaboration. We acknowledge the
939 Confederacion Hidrográfica del Ebro (CHE) and Agència Catalana de l'Aigua (ACA) for
940 providing the official NBLs for WGB44 and WGB5, respectively. Meteorological data
941 have been kindly provided by the Spanish Meteorological Agency (AEMET) and the
942 Meteorological Service of Catalonia (SMC).

943

944

945 **References**

- 946 Aitchison, J., 1986. The Statistical Analysis of Compositional Data. Monographs on
947 Statistics and Applied Probability Chapman and Hall, London, New York (416 pp.).
- 948 Aitchison, J., and Greenacre, M. 2002. Biplots of compositional data. Journal of the Royal
949 Statistical Society: Series C (Applied Statistics), 51(4), 375–392.
950 doi:10.1111/1467-9876.00275
- 951 Appelo, C., Postma, D. 2005. Geochemistry, Groundwater and Pollution. 2nd edition.
952 London: CRC Press, 683 pp. <https://doi.org/10.1201/9781439833544>
- 953 Ballesteros, D., Malard, A., Jeannin, P. Y., Jiménez-Sánchez, M., García-Sansegundo, J.,
954 Meléndez-Asensio, M., & Sendra, G., 2015. KARSYS hydrogeological 3D
955 modeling of alpine karst aquifers developed in geologically complex areas: Picos
956 de Europa National Park (Spain). Environmental Earth Sciences, 74(12), 7699-
957 7714. <https://doi.org/10.1007/s12665-015-4712-0>

- 958 Barbieri, M., Nigro, A., Petitta, M., 2017. Groundwater mixing in the discharge area of
959 San Vittorino Plain (Central Italy): geochemical characterization and implication
960 for drinking uses. *Environmental Earth Sciences*, 76(11), 393.
961 <https://doi.org/10.1007/s12665-017-6719-1>
- 962 Biernacki, C., Govaert, Gérard., 1999. Choosing models in model-based clustering and
963 discriminant analysis. *J. Stat. Comput. Simul.* 64, 49–71.
964 <https://doi.org/10.1080/00949659908811966>
- 965 Blake, S., Henry, T., Murray, J., Flood, R., Muller, M.R., Jones, A.G., Rath, V., 2016.
966 Compositional multivariate statistical analysis of thermal groundwater provenance:
967 A hydrogeochemical case study from Ireland. *Appl. Geochem.* 75, 171–188.
968 <https://doi.org/10.1016/j.apgeochem.2016.05.008>
- 969 Boeuf, B., Fritsch, O., 2016. Studying the implementation of the Water Framework
970 Directive in Europe: a meta-analysis of 89 journal articles. *Ecology and Society*,
971 21(2):19. <http://dx.doi.org/10.5751/ES-08411-210219>
- 972 Bondu, R., Cloutier, V., Rosa, E., Roy, M., 2020. An exploratory data analysis approach
973 for assessing the sources and distribution of naturally occurring contaminants (F,
974 Ba, Mn, As) in groundwater from southern Quebec (Canada). *Appl. Geochem.* 114,
975 104500. <https://doi.org/10.1016/j.apgeochem.2019.104500>
- 976 Bouveyron, C., Brunet-Saumard, C., 2014. Model-based clustering of high-dimensional
977 data: A review. *Comput. Stat. Data Anal.* 71, 52–78.
978 <https://doi.org/10.1016/j.csda.2012.12.008>
- 979 BRIDGE, 2007. Background cRiteria for the IDentification of Groundwater Thresholds.
980 <https://cordis.europa.eu/project/id/6538>.
- 981 Brock, G., Pihur, V., Datta, S., Datta, S. 2008. cIValid: An R Package for Cluster
982 Validation. *Journal of Statistical Software* 25(4). doi: 10.18637/jss.v025.i04
- 983 Buccianti and Grunsky, 2014. Compositional data analysis in geochemistry: Are we sure
984 to see what really occurs during natural processes? *J. Geochem. Explor.* 141, 1–5.
985 <https://doi.org/10.1016/j.gexplo.2014.03.022>
- 986 Carranza, E.J.M., 2011. Analysis and mapping of geochemical anomalies using logratio-
987 transformed stream sediment data with censored values. *J. Geochem. Explor.* 110,
988 167–185. <https://doi.org/10.1016/j.gexplo.2011.05.007>
- 989 Charrad, M., Ghazzali, N., Boiteau, V., Niknafs, A. 2014. NbClust: an R package for
990 determining the relevant number of clusters in data set. *J Stat Soft.* 61:1–36.

- 991 Cloutier, V., Lefebvre, R., Therrien, R., Savard, M.M., 2008. Multivariate statistical
992 analysis of geochemical data as indicative of the hydrogeochemical evolution of
993 groundwater in a sedimentary rock aquifer system. *J. Hydrol.* 353, 294–313.
994 <https://doi.org/10.1016/j.jhydrol.2008.02.015>
- 995 Coetsiers, M., Blaser, P., Martens, K., Walraevens, K., 2009. Natural background levels
996 and threshold values for groundwater in fluvial Pleistocene and Tertiary marine
997 aquifers in Flanders, Belgium. *Environ. Geol.* 57, 1155–1168.
998 <https://doi.org/10.1007/s00254-008-1412-z>
- 999 Comas-Cuffí M, Thió-Henestrosa S. CoDaPack 2.0: a stand-alone, multi-platform
1000 compositional software. In: Egozcue JJ, Tolosana-Delgado R, Ortego MI, eds.
1001 CoDaWork'11: 4th International Workshop on Compositional Data Analysis. Sant
1002 Feliu de Guíxols; 2011.
- 1003 Custodio, E.; Nieto, P.; Manzano, M. 2007. Natural groundwater quality: policy
1004 considerations and European opinion. *The Natural Baseline Quality of Groundwater*
1005 (eds. W.M. Edmunds & P. Shand). Blackwell Publ., Oxford. Chap. 8: 178–194.
1006 ISBN: 978-14051-5675-2. Dempster, A. P. , Laird, N. M.; Rubin, D. B. 1977.
1007 Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the*
1008 *Royal Statistical Society. Series B (Methodological)*, Vol. 39, No. 1. pp. 1-38. DOI:
1009 10.2307/2984875
- 1010 Drew, D., Hötzl, H. (eds.), 1999. Karst Hydrogeology and Human Activities. Impacts,
1011 Consequences and Implications. – International Contributions to hydrogeology
1012 (IAH) 20, 322 p.
- 1013 Ducci, D., Sellerino, M., 2012. Natural background levels for some ions in groundwater
1014 of the Campania region (southern Italy). *Environ. Earth Sci.* 67, 683–693.
1015 <https://doi.org/10.1007/s12665-011-1516-8>
- 1016 Egozcue J, Pawlowsky-Glahn V, Mateu-Figueras F, Barceló-Vidal C. 2003. Isometric
1017 logratio transformations for compositional data analysis. *Math Geol*;35:279–300.
- 1018 Egozcue, J.J. and Pawlowsky-Glahn, V. 2005. Groups of parts and their balances in
1019 compositional data analysis. *Mathematical Geology*, 37(7), 795-828.
- 1020 Egozcue J, Pawlowsky-Glahn V. 2006. Simplicial geometry for compositional data. In:
1021 Buccianti A, Mateu-Figueras G, Pawlowsky-Glahn V, editors. *Compositional data*
1022 *analysis in the geosciences: from theory to practice*. Bath, UK: Geological Society
1023 Publishing House; p. 67–77.

- 1024 Engle, M.A., Rowan, E.L., 2013. Interpretation of Na–Cl–Br Systematics in Sedimentary
1025 Basin Brines: Comparison of Concentration, Element Ratio, and Isometric Log-
1026 ratio Approaches, *Math Geosci* (2013) 45:87–101 DOI 10.1007/s11004-012-9436-
1027 z
- 1028 Farnham, I.M., Sinh, A.k., Stetzenbach, K.J., Johannesson, K.H., 2002. Treatment of
1029 nondetects in multivariate analysis of groundwater geochemistry data.
1030 *Chemometrics and Intelligent Laboratory Systems*. Volume 60, Issues 1–2, 28 pp
1031 265–281. [https://doi.org/10.1016/S0169-7439\(01\)00201-5](https://doi.org/10.1016/S0169-7439(01)00201-5)
- 1032 Filzmoser P, Steiger B. 2009. StatDA: statistical analysis for environmental data. R
1033 package version 1.1. <http://cran.at.r-project.org/web/packages/StatDA/index.html>.
- 1034 Filzmoser P, Hron K, Reimann C. 2009b Univariate statistical analysis of environmental
1035 (compositional) data: problems and possibilities. *Sci Total Environ* 407:6100–8.
1036 <https://doi.org/10.1016/j.scitotenv.2009.08.008>
- 1037 Filzmoser, P., Hron, K., & Templ, M. 2018. Applied Compositional Data Analysis. With
1038 Worked Examples in R. Springer Series in Statistics. doi:10.1007/978-3-319-
1039 96422-5
- 1040 Fraley C. and Raftery A. E. 2002. Model-based clustering, discriminant analysis and
1041 density estimation, *Journal of the American Statistical Association*, 97/458, pp.
1042 611–631
- 1043 Fraley C., Raftery A. E., Murphy T. B. and Scrucca L. 2012. mclust Version 4 for R:
1044 Normal Mixture Modeling for Model-Based Clustering, Classification, and Density
1045 Estimation. Technical Report No. 597, Department of Statistics, University of
1046 Washington
- 1047 Gabriel, K.R., 1971. The biplot-graphic display of matrices with application to principal
1048 component analysis. *Biometrika* 58, 453e467.
- 1049 Goldscheider, N. 2005. Fold structure and underground drainage pattern in the alpine
1050 karst system Hochifen-Gottesacker. *Eclogae geol. Helv.* 98, 1–17.
1051 <https://doi.org/10.1007/s00015-005-1143-z>
- 1052 Güller, C., Thyne, G.D. 2004. Delineation of hydrochemical facies distribution in a
1053 regional groundwater system by means of fuzzy c- means clustering. *Water Resour
1054 Res* 40:W12503. <https://doi.org/10.1029/2004WR003299>
- 1055 He Kim, S.H., Choi, B., Lee, G., Yun, S., Kim S. 2019. Compositional data analysis and
1056 geochemical modeling of CO₂–water–rock interactions in three provinces of

- 1057 Korea. Environ Geochem Health 41, 357–380. <https://doi.org/10.1007/s10653-017-0057-9>
- 1059 Herms, I., Jódar, J., Soler, A., Vadillo, I., Lambán, L. J., Martos-Rosillo, S., Núñez, J.A.,
1060 Arnó, G., Jorge, J., 2019. Contribution of isotopic research techniques to
1061 characterize high-mountain-Mediterranean karst aquifers: The Port del Comte
1062 (Eastern Pyrenees) aquifer. Science of The Total Environment, 656, 209-230.
1063 <https://doi.org/10.1016/j.scitotenv.2018.11.188>
- 1064 Hinsby, K., Condesso de Melo, M.T., Dahl, M., 2008. European case studies supporting
1065 the derivation of natural background levels and groundwater threshold values for
1066 the protection of dependent ecosystems and human health. Sci. Total Environ. 401,
1067 1–20. <https://doi.org/10.1016/j.scitotenv.2008.03.018>
- 1068 ICGC, 2007. Mapa Geològic Comarcal de Catalunya 1:50,000. Full Alt Urgell
1069 (BDGC50M).<http://www.icgc.cat/ca/Administracio-i-empresa/Descarregues/Cartografia-geologica-i-geotematica/Cartografia-geologica/Mapa-geologic-comarcal-de-Catalunya-1-50.000/Mapa-geologic-comarcal-de-Catalunya-1-50.000>.
- 1073 Kazakis, N., Chalikakis, K., Mazzilli, N., Ollivier, C., Manakos, A., Voudouris, K. 2018.
1074 Management and research strategies of karst aquifers in Greece: Literature
1075 overview and exemplification based on hydrodynamic modelling and vulnerability
1076 assessment of a strategic karst aquifer, Sci. Total Environ. 643, 592-609,
1077 <https://doi.org/10.1016/j.scitotenv.2018.06.184>.
- 1078 Kim, K.-H., Yun, S.-T., Park, S.-S., Joo, Y., Kim, T.-S., 2014. Model-based clustering of
1079 hydrochemical data to demarcate natural versus human impacts on bedrock
1080 groundwater quality in rural areas, South Korea. J. Hydrol. 519, 626–636.
1081 <https://doi.org/10.1016/j.jhydrol.2014.07.055>
- 1082 Kim, K.-H., Yun, S.-T., Kim, H.-K., Kim, J.-W., 2015. Determination of natural
1083 backgrounds and thresholds of nitrate in South Korean groundwater using model-
1084 based statistical approaches. J. Geochem. Explor. 148, 196–205.
1085 <https://doi.org/10.1016/j.gexplo.2014.10.001>
- 1086 Kresic, N., Stevanović, Z., 2010. Groundwater hydrology of springs: engineering, theory,
1087 management, and sustainability. Butterworth-Heinemann, Oxford
- 1088 Lambán, L. J., Jódar, J., Custodio, E., Soler, A., Saprizá, G., & Soto, R. (2015). Isotopic
1089 and hydrogeochemical characterization of high-altitude karst aquifers in complex
1090 geological settings. The Ordesa and Monte Perdido National Park (Northern Spain)

- 1091 case study. *Science of the Total Environment*, 506, 466-479.
1092 <http://dx.doi.org/10.1016/j.scitotenv.2014.11.030>
- 1093 Luetscher, M., Perrin, J. 2005. The Aubonne karst aquifer (Swiss Jura). *Eclogae
1094 Geologicae Helvetiae*, 98(2), 237-248. <https://doi.org/10.1007/s00015-005-1156-7>
- 1095 Marandi, A., Karro, E., 2008. Natural background levels and threshold values of
1096 monitored parameters in the Cambrian-Vendian groundwater body, Estonia.
1097 *Environ. Geol.* 54, 1217–1225. <https://doi.org/10.1007/s00254-007-0904-6>
- 1098 Marín, A.I., Andreo, B., 2015. Vulnerability to Contamination of Karst Aquifers. In:
1099 Stevanović Z. (eds) *Karst Aquifers—Characterization and Engineering.*
1100 Professional Practice in Earth Sciences. Springer, Cham.
1101 https://doi.org/10.1007/978-3-319-12850-4_8
- 1102 Martín-Fernández, J.A., Barceló-Vidal, C., Pawlowsky-Glahn, V., 2003. Dealing with
1103 zeros and missing values in compositional data sets using nonparametric
1104 imputation. *Mathematical Geology* 35 (3), 253–278.
- 1105 Merchán, D., Auqué, L.F., Acero, P., Gimeno, M.J., Causapé, J., 2015. Geochemical
1106 processes controlling water salinization in an irrigated basin in Spain: Identification
1107 of natural and anthropogenic influence. *Sci. Total Environ.* 502, 330–343.
1108 <https://doi.org/10.1016/j.scitotenv.2014.09.041>
- 1109 Moya, C.E., Raiber, M., Taulis, M., Cox, M.E., 2015. Hydrochemical evolution and
1110 groundwater flow processes in the Galilee and Eromanga basins, Great Artesian
1111 Basin, Australia: A multivariate statistical approach. *Sci. Total Environ.* 508, 411–
1112 426. <https://doi.org/10.1016/j.scitotenv.2014.11.099>
- 1113 Müller, D., Blum, A., Hart, A., Hookey, J., Kunkel, R., Scheidleder, A., Tomlin, C.,
1114 Wendland, F., 2006. Final proposal for a methodology to set up groundwater
1115 threshold values in Europe. In: Report to the EU project “BRIDGE”, Deliverable
1116 D18.
- 1117 Muñoz, J.A., Mencos, J., Roca, E., Carrera, N., Gratacós, A., Ferrer, O. Fernández, O.
1118 2018. The structure of the South-Central-Pyrenean fold and thrust belt as
1119 constrained by subsurface data. *Geologica Acta*, Vol.16, No 4, 439-460 DOI:
1120 10.1344/GeologicaActa2018.16.4.7
- 1121 Nieto, P., Custodio, E., Manzano, M., 2005. Baseline groundwater quality: a European
1122 approach. *Environmental Science & Policy*, 8(4), 399-409.
1123 <https://doi.org/10.1016/j.envsci.2005.04.004>

- 1124 Otero, N., Tolosana-Delgado, R., Soler, A., Pawlowsky-Glahn, V., Canals, A., 2005.
1125 Relative vs. absolute statistical analysis of compositions: A comparative study of
1126 surface waters of a Mediterranean river. *Water Res.* 39, 1404–1414.
1127 <https://doi.org/10.1016/j.watres.2005.01.012>
- 1128 Owen, D.Des.R., Pawlowsky-Glahn, V., Egozcue, J.J., Buccianti, A., Bradd, J.M., 2016.
1129 Compositional data analysis as a robust tool to delineate hydrochemical facies
1130 within and between gas-bearing aquifers: COMPOSITIONAL DATA ANALYSIS
1131 TO DELINEATE WATER TYPES. *Water Resour. Res.* 52, 5771–5793.
1132 <https://doi.org/10.1002/2015WR018386>
- 1133 Palarea-Albaladejo J, Martin-Fernandez JA, Olea, RA. A bootstrap estimation scheme for
1134 chemical compositional data with nondetects. *Journal of Chemometrics* 2014; 28:
1135 585-599. <https://doi.org/10.1002/cem.2621>
- 1136 Palarea-Albaladejo and Martín-Fernández, 2015. zCompositions — R package for
1137 multivariate imputation of left-censored data under a compositional approach.
1138 *Chemometrics and Intelligent Laboratory Systems*, Volume 143, pp 85-96, ISSN
1139 0169-7439, <https://doi.org/10.1016/j.chemolab.2015.02.019>.
- 1140 Parkhurst, D.L., Appelo, C.A.J. 2013. Description of Input and Examples for PHREEQC
1141 Version 3—A Computer Program for Speciation, Batch-Reaction, One-
1142 Dimensional Transport, and Inverse Geochemical Calculations: U.S. Geological
1143 Survey Techniques and Methods, Book 6, Chap. A43: 1–497. (Available only at
1144 <https://pubs.usgs.gov/tm/06/a43>. Last access 28 August 2020).
- 1145 Parrone, D., Ghergo, S., Preziosi, E., 2019. A multi-method approach for the assessment
1146 of natural background levels in groundwater. *Sci. Total Environ.* 659, 884–894.
1147 <https://doi.org/10.1016/j.scitotenv.2018.12.350>
- 1148 Pawlowsky-Glahn, V., Egozcue, J.J., Tolosana-Delgado, R. 2015. Modeling and Analysis
1149 of Compositional Data. ed. John Wiley & Sons Ltd, The Atrium, Southern Gate,
1150 Chichester, West Sussex, PO19 8SQ, United Kingdom. 272 pages. ISBN:
1151 9781118443064
- 1152 Peel, M.C., Finlayson, B.L., McMahon, T.A., 2007. Updated world map of the Köppen–
1153 Geiger climate classification. *Hydrol. Earth Syst. Sci.* 11, 1633–1644.
1154 <https://doi.org/10.5194/hess-11-1633-2007>.
- 1155 Piña, A., Donado, L.D., Blake, S., Cramer, T., 2018. Compositional multivariate
1156 statistical analysis of the hydrogeochemical processes in a fractured massif: La

- 1157 Línea tunnel project, Colombia. Appl. Geochem. 95, 1–18.
1158 <https://doi.org/10.1016/j.apgeochem.2018.05.012>
- 1159 Preziosi, E., Giuliano, G., Vivona, R., 2010. Natural background levels and threshold
1160 values derivation for naturally As, V and F rich groundwater bodies: a
1161 methodological case study in Central Italy. Environ. Earth Sci. 61, 885–897.
1162 <https://doi.org/10.1007/s12665-009-0404-y>
- 1163 Puig, R., Tolosana-Delgado, R., Otero, N., Folch, A., 2011. Combining isotopic and
1164 compositional data: a discrimination of regions prone to nitrate pollution. In V.
1165 Pawlowsky-Glahn and A. Buccianti (Eds.), Compositional Data Analysis: Theory
1166 and Applications 390.
- 1167 Raftery, AE., Scrucca, L., Brendan, T., Fop, M. 2020. Gaussian Mixture Modelling for
1168 Model-Based Clustering, Classification, and Density Estimation. Package ‘mclust’.
1169 Version 5.4.6. <https://mclust-org.github.io/mclust/>
- 1170 Reimann, C., Filzmoser, P. 2000. Normal and lognormal data distribution in
1171 geochemistry: death of a myth. Consequences for the statistical treatment of
1172 geochemical and environmental data. Environmental Geology 39, 1001–1014.
1173 <https://doi.org/10.1007/s002549900081>
- 1174 Reimann C, Filzmoser P, Garrett RG, Dutter R. 2008. Statistical data analysis explained.
1175 Applied environmental statistics with R. Wiley, Chichester, UK, 362 pp. ISBN:
1176 978-0-470-98581-6
- 1177 Reimann, C., Filzmoser, P., Fabian, K., Hron, K., Birke, M., Demetriadis, A., Dinelli, E.,
1178 Ladenberger, A., 2012. The concept of compositional data analysis in practice —
1179 Total major element concentrations in agricultural and grazing land soils of Europe.
1180 Sci. Total Environ. 426, 196–210. <https://doi.org/10.1016/j.scitotenv.2012.02.032>
- 1181 Sánchez, D., Antonio Barberá, J., Mudarra, M., Andreo, B., & Martín, J. F. 2017.
1182 Hydrochemical and isotopic characterization of carbonate aquifers under natural
1183 flow conditions, Sierra Grazalema Natural Park, southern Spain. Geological
1184 Society, London, Special Publications, 466(1), 275–293. doi:10.1144/sp466.16
- 1185 Scrucca, L. 2010. Dimension reduction for model-based clustering. Statistics and
1186 Computing, 20(4), pp. 471-484.
- 1187 Scrucca, L., Fop, M., Murphy, T.B., Raftery, A.E., 2016. Mclust 5: clustering,
1188 classification and density estimation using Gaussian finite mixture models, The R
1189 Journal, 8/1, pp. 289-317.

- 1190 Serianz, L., Cerar, S. & Šraj, M. 2020. Hydrogeochemical characterization and
1191 determination of natural background levels (NBL) in groundwater within the main
1192 lithological units in Slovenia. Environ Earth Sci 79, 373.
1193 <https://doi.org/10.1007/s12665-020-09112-1>
- 1194 Shelton, J.L., Engle, M.A., Buccianti, A., Blondes, M.S., 2018. The isometric log-ratio
1195 (ilr)-ion plot: A proposed alternative to the Piper diagram. J. Geochem. Explor. 190,
1196 130–141. <https://doi.org/10.1016/j.gexplo.2018.03.003>
- 1197 Stevanović, Z., 2019. Karst waters in potable water supply: a global scale overview.
1198 Environ Earth Sci 78, 662. <https://doi.org/10.1007/s12665-019-8670-9>
- 1199 Struckmeier WF, Gilbrich WH, Gun Jvd, Maurer S, Puri S, Richts A, Winter P, Zaepke
1200 M. 2006. WHYMAP and the groundwater resources map of the world at the scale
1201 of 1:50 000 000. Special edition for the 4th world water forum, Mexico City, March
1202 2006. BGR Hannover/UNESCO, Paris.
- 1203 Suk, H., and K.K. Lee. 1999. Characterization of a groundwater hydrochemical system
1204 through multivariate analysis: clustering into groundwater zones. Groundwater v.
1205 37, no. 3pp. 358-366.
- 1206 Templ, M., Filzmoser, P., Reimann, C., 2008. Cluster analysis applied to regional
1207 geochemical data: Problems and possibilities. Appl. Geochem. 23, 2198–2213.
1208 <https://doi.org/10.1016/j.apgeochem.2008.03.004>
- 1209 Tolosana-Delgado, R., Otero, N., Soler, A. 2005. A compositional approach to stable
1210 isotope data analysis. Conference: CODAWORK'05
- 1211 Van den Boogaart, K.G. and Tolosana-Delgado, R. (2008) "compositions": a unified R
1212 package to analyze Compositional Data, Computers & Geosciences. 34 (4), 320-
1213 338. <https://doi.org/10.1016/j.cageo.2006.11.017>
- 1214 Vergés, J. 1999. Estudi geològic del vessant sud del Pirineu oriental i central. Evolució
1215 cinemàtica en 3D. PhD Thesis. University of Barcelona (UB), Faculty of Geology,
1216 180 pp.
- 1217 Vivenzio, D., Kummu, M., Meybeck, M., Kallio, M., Wada, Y. 2020. Increasing
1218 dependence of lowland populations on mountain water resources. Nature
1219 Sustainability, 1-12. <https://doi.org/10.1038/s41893-020-0559-9>
- 1220 Wendland, F., Blum, A., Coetsiers, M., Gorova, R., Griffioen, J., Grima, J., Hinsby, K.,
1221 Kunkel, R., Marandi, A., Melo, T., Panagopoulos, A., Pauwels, H., Ruisi, M.,
1222 Traversa, P., Vermooten, J.S.A., Walraevens, K., 2008. European aquifer typology:
1223 a practical framework for an overview of major groundwater composition at

- 1224 European scale. Environ. Geol. 55, 77–85. <https://doi.org/10.1007/s00254-007-0966-5>
- 1225
- 1226 WFD 2000, Water Framework Directive, 2000. Directive 2000/60/CE of the European
1227 Parliament (ECOJ 22 December 2000).
http://www.bygg.ntnu.no/borsanyi/eamn_web/documents/wfd-es.pdf.
- 1228
- 1229 Wu, X., Zheng, Y., Zhang, J., Wu, B., Wang, S., Tian, Y., Li, J., Meng, X., 2017.
1230 Investigating Hydrochemical Groundwater Processes in an Inland Agricultural
1231 Area with Limited Data: A Clustering Approach. Water 9, 723.
1232 <https://doi.org/10.3390/w9090723>
- 1233 Yidana, S.M. 2010. Groundwater classification using multivariate statistical methods:
1234 Southern Ghana. Journal of African Earth Sciences 57(5):455-469 doi:
1235 10.1016/j.jafrearsci.2009.12.002
- 1236 Yolcubal, İ., Gündüz, Ö.C.A., Kurtuluş, N., 2019. Origin of salinization and pollution
1237 sources and geochemical processes in urban coastal aquifer (Kocaeli, NW Turkey).
1238 Environmental Earth Sciences, 78(6), 181. <https://doi.org/10.1007/s12665-019-8181-8>
- 1239
- 1240 Zabala, M.E., Martínez, S., Manzano, M., Vives, L., 2016. Groundwater chemical
1241 baseline values to assess the Recovery Plan in the Matanza-Riachuelo River basin,
1242 Argentina. Sci. Total Environ. 541, 1516–1530.
1243 <https://doi.org/10.1016/j.scitotenv.2015.10.006>
- 1244 Zwahlen, F. (ed.), 2004. Vulnerability and risk mapping for the protection of carbonate
1245 (karst) aquifers, final report COST Action 620. EUR 20912, European
1246 Commission, Brussels, 297 p.
- 1247
- 1248

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

The authors: Ignasi Herms. Jorge Jódar, Albert Soler, Luis Javier Lambán, Emilio Custodio, Joan Agustí Núñez, Georgina Arnó, Maribel Ortego, David Parcerisa and Joan Jorge