
K nearest neighbors classification of water masses in the western Alboran Sea using the sigma-pi diagram

Belattmania Ayoub ¹, El Arrim Abdelkrim ¹, Ayouche Adam ²,*, Charria Guillaume ⁵, Hilmi Karim ³, El Mounni Bouchta ⁴

¹ Faculté des Sciences et Techniques, Tanger, Morocco

² Laboratory for Ocean Physics and Satellite Remote Sensing (LOPS), UMR6523, Ifremer, Univ. Brest, CNRS, IRD, Brest, France

³ Institut National de Recherche Halieutique, Casablanca, Morocco

⁴ Université Abdelmalek Essaâdi, Tétouan, Morocco

⁵ Laboratory for Ocean Physics and Satellite Remote Sensing (LOPS), UMR6523, Ifremer, Univ. Brest, CNRS, IRD, Brest, France

* Corresponding author : Adam Ayouche, email address : adam.ayouche@ensta-bretagne.org

Abstract :

Different classification techniques of water masses have been developed using the potential temperature-salinity (θ -S) diagram and its volumetric analysis. In this study, we propose a new method to automatically classify water masses via a supervised machine learning algorithm based on the K nearest neighbors (Knn), in the potential density and potential spicity (σ - π) coordinates. This method is applied to temperature and salinity data collected in the western side of the Alboran Sea during a glider mission, dedicated to sample the Western Alboran Gyre (WAG) in late winter 2021. The water masses in the studied region were classified into five different categories following a supervised learning process, based on ocean profile databases available on the region of interest. The results corroborate previous studies of the spatial distribution of water masses in the Alboran Sea, inferred from traditional method based on the expert analysis of the (θ -S) diagram, and suggest that this methodology is efficient and reliable for water masses classification. Compared to a classical clustering computation (herein k-means), this method is more appropriate in a region where the characteristics of the water masses change considerably in both space and time.

Highlights

- ▶ High spatial resolution glider profiles of θ -S in the western Alboran sea.
- ▶ Water masses derived on a (σ - π) diagram using Knn algorithm.
- ▶ Classification results confirm earlier derived circulation schemes.
- ▶ The proposed method outperforms classical clustering analysis in delineating water mass boundaries.

Keywords : Alboran Sea, Western Alboran Gyre, water masses, (σ - π) diagram, K nearest neighbor classification.

1. Introduction

2 A water mass is a volume of oceanic water with horizontal and vertical ex-
3 tensions, and having specific physical characteristics. In general, most of the
4 water masses are formed by atmosphere-ocean exchanges, however some oth-
5 ers acquire their characteristics (e.g minimum salinity) through biochemical or
6 physical processes (e.g convection). The signature of such characteristics are
7 represented by tracers such as the potential temperature and salinity. These
8 tracers are important to understand the oceanic circulation at different global
9 and regional scales, as the thermohaline circulation (Broecker, 1991). The ther-
10 mohaline circulation plays a key role in the climate regulation by the transport
11 of heat, carbon and oxygen across the different basins around the world (Clark
12 et al., 2002).

13 In this context, Pantiulin (2002) sketches a brief history about the genesis of
14 the concept of water masses, depending on the evolution of the in-situ observa-
15 tions of temperature and salinity. Indeed, the definition, classification and first
16 principles of water masses appeared for the first time in the monograph called
17 the Norwegian Sea in 1909 (Hansen and Nansen). The latter was followed by
18 the introduction of the potential temperature-salinity (θ -S) diagram as a tool
19 to analyze water masses properties, in a Norwegian study after the first world
20 war (Hansen, 1916). He showed on a wide area of the eastern Atlantic ocean
21 that the variations in the (θ -S) diagram can be attributed to the intrusion of
22 offshore water masses.

23 Since then, the (θ -S) diagram has been used widely in physical oceanog-
24 raphy and by numerous authors across different fields. Major progress in the
25 water mass analysis was the introduction of the volumetric (θ -S) diagram which
26 was used in different studies that includes the Pacific ocean, the Indian ocean,
27 the Atlantic ocean and the Global ocean (Cochrane, 1958; Pollak, 1958; Mont-
28 gomery, 1958). In these studies, the quantity of volumetric units for standard
29 levels of depths were estimated statistically. This estimation was based on a

30 division of the oceans in bi-variate classes defined by their temperature and
31 salinity.

32 Other studies followed the previous ones based on the volumetric (θ -S) sta-
33 tistically analysis methodology. They improved and reworked this methodology
34 for the sake of understanding the water masses distribution in a volumetric (θ -S)
35 diagram (Miller and Stanley, 1961; Wright and Worthington, 1970; Worthing-
36 ton, 1981).

37 Besides the volumetric (θ -S) diagram analysis, other techniques of water
38 mass classification have been used such as the cluster analysis where the data
39 are grouped on the basis of a set of measured parameters. The objective of
40 this method is to find an optimal data distribution which minimizes a certain
41 metric that define the similarity within the clusters. For example, Kim et al.
42 (1991) applied a clustering analysis based on the average linkage between groups
43 for the temperature and salinity to identify the water masses in the Yellow sea
44 and the East China sea. The metric used for their clustering analysis is the
45 squared Euclidean distance defined as the normalized temperature and salinity
46 differences between points.

47 Hur et al. (1999) studied the yellow and east china seas for over 40 years
48 (1950-1992) using historical data of temperature and salinity. They included in
49 their study the geographical distance and the depth separation in computing the
50 distance for the clustering method. Naranjo et al. (2015) examined the distri-
51 bution and spatio-temporal evolution of water masses in the strait of Gibraltar
52 using clustering analysis. These authors used historical values of potential tem-
53 perature, salinity and potential density for each water mass as initial centroids
54 for the classification. Roseli et al. (2015) applied the k-means algorithm on
55 temperature and salinity data from CTD casts in two different seasons (fall and
56 summer), to classify water masses at the Shallow Sunda Shelf of Southern South
57 China Sea. Recently, Gao et al. (2020) proposed a novel and robust method to
58 identify the frontiers between water masses in the Northern South China sea.
59 Their identification of the water masses center is based on ranges and standard
60 deviations of the potential spicity π in different potential density layers, and

61 water volumetric distributions in the bi-dimensional plan (σ - π).

62 The (θ - S) diagram and its different techniques for analyzing water masses
 63 have been developed for several oceans; but it is also interesting to conduct such
 64 studies for regions where several water masses from different oceans can interact,
 65 such as the Alboran Sea : the westernmost Mediterranean sub-basin where
 66 Atlantic and Mediterranean waters interact through the strait of Gibraltar. In
 67 our knowledge, only traditional water masses analysis based on (θ - S) diagram,
 68 have been previously used in this region (Bryden et al., 1982; Pistek et al., 1985;
 69 Gascard and Richez, 1985; Parrilla et al., 1986; Parrilla and Kinder, 1987; Millot
 70 et al., 2006; Millot, 2009; Renault et al., 2012; Millot, 2014).

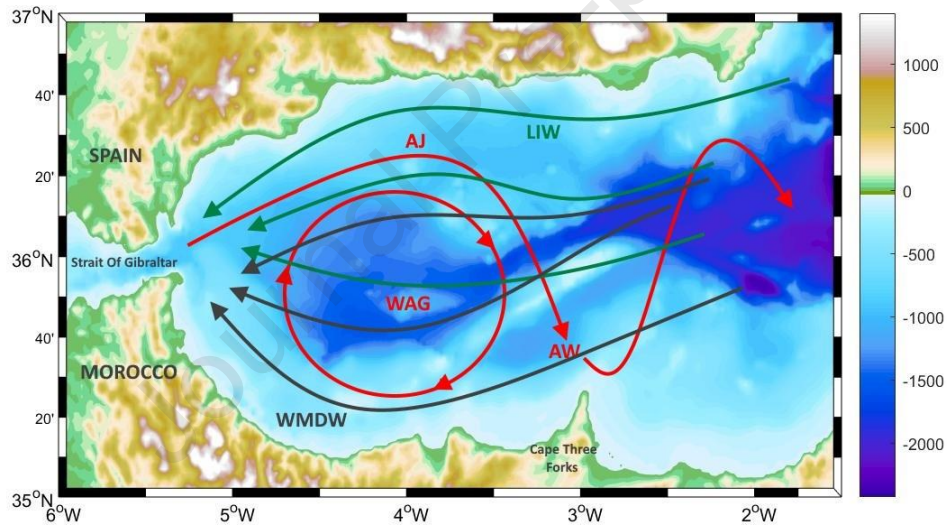


Figure 1: Map of the Western Alboran Sea sketching the bathymetric depths and topographic elevations in meters (m) relative to the mean sea level. Red arrows show the general surface circulation of Atlantic Water (AW), showing the Western Alboran Gyre (WAG) as well as the Atlantic Jet (AJ). The green and black arrows represent respectively the intermediate and deep circulation of Mediterranean waters (LIW and WMDW).

71 The related circulation schemes have been sketched for each water mass of
 72 the Atlantic ocean and the Mediterranean sea (figure 1) where their properties
 73 and distributions are summarized in table 1 and can be described as follows:

74 The Atlantic Water (AW), located in the western side of the strait of Gibralt-
75 tar, is injected in the Alboran Sea (top 200 m depth). It is subject to different
76 variations through its cyclonic path (Coriolis effect) at the surface due to its in-
77 teraction with the atmosphere and the surface mixing with older Atlantic water.
78 This water becomes saltier (~ 38 psu) and progressively cooler in winter (~ 13
79 $^{\circ}\text{C}$) and therefore this results in increased density. Then, this water is called the
80 Modified Atlantic Water (MAW). At the surface layer, quasi-homogeneous light
81 waters are observed with a salinity $S = 36.6$ psu: the Surface Atlantic Water
82 (SAW). A second layer, the North Atlantic Central Waters (NACW), is char-
83 acterized by a minimum of salinity separating the SAW from the MAW. This
84 separation is progressively dissipated through mixing in the strait of Gibraltar
85 and the Alboran Sea ($T=11-17$ $^{\circ}\text{C}$, $S=35.6-36.5$ psu).

86 Previous studies about the Mediterranean Waters (MWs) in the Alboran
87 Sea suggest the presence of Winter Intermediate Water (WIW), Levantine In-
88 termediate Water (LIW), Western Mediterranean Deep Waters (WMDW) and
89 the Tyrrhenian Deep Water (TDW). The LIW and WMDW were considered
90 as the main contributors for the outflow. The WIW results from AW cooling
91 along the continental shelf of the Liguro-Provencal sub-basin and is generated
92 periodically in the Alboran Sea near the Spanish continental shelf. The WIW
93 can be identified by its minimum potential temperature ($12.9-13$ $^{\circ}\text{C}$) between
94 100 and 350 m depth and between 28 and 29 $\text{kg}\cdot\text{m}^{-3}$ isopycnals. The LIW from
95 the Western Mediterranean sea generated by winter convection is the most salty
96 and warmest water mass encountered at mid depth (200-600 m) in the Albo-
97 ran Sea. The LIW is mostly concentrated in the north and center sides of the
98 Alboran Sea and absent along the African coast. The LIW is characterized
99 by temperature and salinity maximum ($13.1-13.3$ $^{\circ}\text{C}$, $38.47-38.52$ psu). The
100 WMDW is generated in the gulf of Lion by deep convection and is cold (< 12.9
101 $^{\circ}\text{C}$) and relatively salty (> 38.4 psu) water. The WMDW is considered as the
102 most dense water in the Mediterranean sea (at 800 m depth in the central part
103 of the Alboran Sea). The TDW is the result of mixing between ancient WDMW
104 in the Tyrrhenian sea and the LIW coming from the Western Mediterranean

105 sea through the strait of Sicily. The TDW is slightly denser than the LIW
 106 and lighter than the WMDW and lies between these two water masses. In the
 107 Alboran Sea, the temperature and salinity values of the TDW are respectively
 108 within the range 13-13.1 °C and 38.41-38.51 psu.

Water mass	Description	Reference
SAW	Quasi homogenous salinity layer ($S \approx 36.6$) and a constant temperature gradient.	(Gascard and Richez, 1985; Parrilla et al., 1986; Vélez-Belchi et al., 2005)
NACW	The Separation layer between the SAW and MAW. It's characterized by a salinity minimum (35.5-36.6) that attenuated quite rapidly after entering the Mediterranean Sea.	(Gascard and Richez, 1985; Parrilla et al., 1986; Vélez-Belchi et al., 2005)
MAW	A mixture layer of Atlantic (16°C-36.5) and Mediterranean waters (12.9°C-38.45)	(Gascard and Richez, 1985; Parrilla et al., 1986; Vélez-Belchi et al., 2005)
LIW	The warmest and saltiest Mediterranean waters, easily recognised anywhere in the sea. Concerning the western Alboran, it is characterised by ($T=13.1-13.2^{\circ}\text{C}$ and $S=38.5$).	(Gascard and Richez, 1985; Parrilla and Kinder, 1987; Millot et al., 2006; Millot, 2009, 2014)
WIW	Results from the AW wintertime cooling in the northern part of the western basin and characterised by a Temperature minimum (12.9-13°C).	(Millot, 2009, 2014)
TDW	Results from mixing between ancient WDMW in the Tyrrhenian sea and the LIW coming from the Western Mediterranean sea. Its core characteristics are in ranges ($T = 13.0-13.1^{\circ}\text{C}$ and $S = 38.48-38.51$).	(Millot et al., 2006; Millot, 2009, 2014)
WMDW	Formed in the Liguro-Provencal mainly from an AW-LIW mixture by wintertime convection processes. It is Cold ($< 12.9^{\circ}\text{C}$) and relatively salty (> 38.4).	(Gascard and Richez, 1985; Parrilla and Kinder, 1987; Millot et al., 2006; Millot, 2009, 2014)

Table 1: Summary of water masses definitions with their respective references.

110 The application of clustering analysis methods for the purpose of automatically
111 classify water masses, has yielded encouraging results in many regions. Nev-
112 ertheless, these techniques have revealed many shortcomings in region with a
113 high spatio-temporal variability and could not exactly identify the water mass
114 boundary (Gao et al. (2020)). Clustering analysis is particularly relevant to dis-
115 tinguish water masses with similar salinity and temperature variance (Naranjo
116 et al. (2015)). This is not the case in the Alboran Sea, where SAW is widely
117 variable in temperature and the MWs range much more in temperature than in
118 salinity.

119 Within this context of challenges to be solved in automatic water masses
120 classification notably in region where intense mixing occurs, in this paper, we
121 propose a novel methodology that classify automatically water masses in the
122 Alboran Sea, based on machine learning supervised algorithm, applied on curvi-
123 linear potential density and potential spicity (σ - π) diagram. Two datasets,
124 described in section 2, have been used for the study. The first concerns the
125 global database of temperature and salinity vertical profiles used as a training
126 dataset of the algorithm. The second is relative to the glider in-situ observations
127 collected in the Western Alboran Sea, on which the classification algorithm is
128 applied. Section 2 also describes the classification methodology, starting with
129 the labeling process and ending with the sensitivity test of the employed method.
130 Water mass classification results in the glider transects are provided in section
131 3 and are discussed in section 4. Finally, conclusions are drawn in section 5.

132 2. Materials and Methods

133 2.1. Database

134 To build the training water mass classes, we assemble the available in-situ
135 observations from oceanographic databases such as World Ocean Database 2018
136 'WOD18' (Boyer et al., 2019) and the Global Data Assembly Centers 'GDACs'
137 (Argo, 2021) in a given geographic domain (Figure 2).

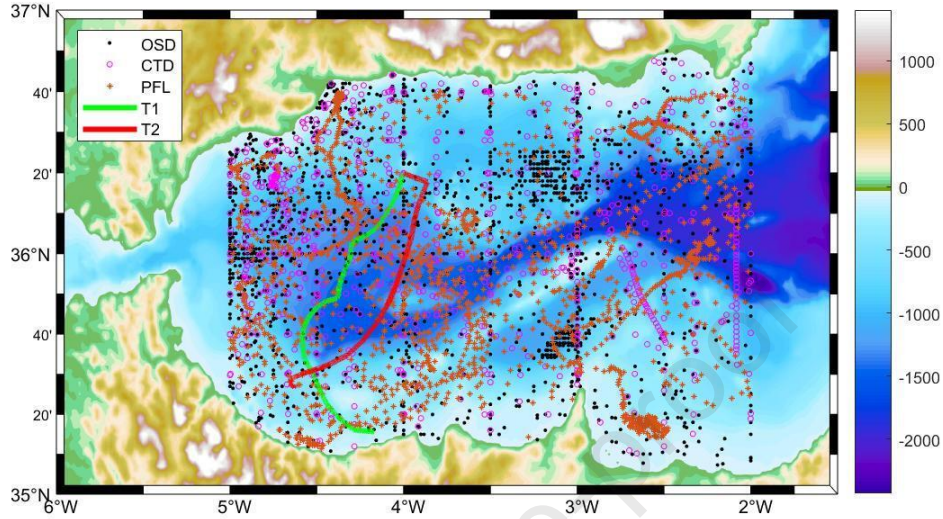


Figure 2: Map of the Western Alboran Sea sketching the bathymetric depths and topographic elevations in meters (m) relative to the mean sea level. The black dots and magenta circles indicate the localization of the vertical profiles of WOD18 related to Ocean Station Data (OSD) dataset and Conductivity Temperature Depth (CTD) dataset respectively (Table 2). Brown asterisks represent Argo Profiling Floats (PFL) trajectories from GDACs (Table 2). The first glider transect (T1) is sketched in green and the second transect (T2) in red.

138 The resultant product is based on 5068 vertical profiles of temperature and
 139 salinity including 1759 sampling cycles of Argo floats. The in-situ data are
 140 gathered over a broad range of temporal scales between 1951 and 2020. Table
 141 2 summarizes the key informations about the mentioned databases.

142

143 2.2. Glider data

144 In this study, our analysis is focused on the second mission performed by
 145 the Moroccan association AGIR (Leader of the Marine Observatory of Al Ho-
 146 ceima) in the Western Alboran Sea, as part of the European project ODYSSEA
 147 (<https://odysseaplatform.eu/fr/home-fr/>). This mission was conducted
 148 after a first one in late fall 2020 (from 10 November to 11 December) in the

Dataset	Description	Temporal range	# of casts	# of TS observations
OSD	measurements made from a stationary vessels using reversing thermometers mounted on special bottles including LVR CTD rosette system, LVR STD and LVR XCTD (XCTD is collected from moving vessels).	1951-2011	2344 stations	26568
CTD	data from a stationary vessels using HVR CTD rosette system, STD (The salinity S is computed from the conductivity) data measured at high frequency with respect to depth as well as HVR XCTD (XCTD is collected from moving vessels).	1975-2018	965 stations	344881
PFL	contains temperature and salinity data collected from drifting profiling floats of the Argo project.	2006-2020	1759 cycles	249788

Table 2: Database information used for the supervised learning. The acronyms XCTD, STD, LVR and HVR stand respectively for: eXpandable Conductivity Temperature Depth, Salinity Temperature Depth, Low Vertical Resolution and High Vertical Resolution. All casts with a depth increment less than two meters are considered High Resolution otherwise, the casts are considered as Low Resolution.

149 same region (Nibani et al., 2021). The second mission occurred in late winter
150 – early spring (from 11 February to 23 March). During this mission, a Sea-
151 Explorer glider (manufactured and commercialized by ALSEAMAR in France),
152 equipped by a Seabird CTD, performed a total of 873 cycles from the surface
153 to approximately 500 m depth with a sampling rate of 4 seconds. Only a
154 part of these cycles (during the 11 first days of the mission) was dedicated to
155 sample the WAG and have been studied herein (Figure 2). In this paper, only
156 the classification of water masses in the WAG and the ambient environment will
157 be discussed.

158 Isotherms and isohalines sketched hereafter in all vertical sections as continuous
 159 lines, represent the interpolated temperature and salinity on a grid of (horizontal
 160 and vertical) resolution $dx=1.1\text{km}$ and $dz=1\text{m}$. The interpolation is performed
 161 using the optimal spatial kriging. In order to remove high frequencies, the inter-
 162 polated data were smoothed using a gaussian filter with a width corresponding
 163 to the radius of deformation in the studied region (Bosse et al., 2015). The
 164 parameter $p \in \{T, S\}$ of each transect is transformed in a smoothed parameter
 165 $\tilde{p} \in \{\tilde{T}, \tilde{S}\}$ by a convolution product:

$$(1) \quad \tilde{p}(x, z) = \int_{x_{min}}^{x_{max}} p(x, z) \times \exp \frac{-x^2}{2L^2} dx$$

166 Where x is the distance along the section, x_{min} and x_{max} the section limits,
 167 z the depth and L the standard deviation. Taking $L = 15 \text{ km}$ is sufficient to
 168 conserve the signal linked to the WAG.

169 In addition to the glider data described previously, and in order to further
 170 test the performance of our method, more examples of data and their related
 171 classification results are represented in Appendix A.

172 2.3. Data single-labeling

173 To build a training dataset with a unique labeling, each sample of the
 174 database has been attributed to a water mass from those described in the intro-
 175 duction $\{\text{SAW}, \text{NACW}, \text{MAW}, \text{WIW}, \text{LIW}, \text{TDW}, \text{WMDW}\}$. To keep a clear
 176 physical sense, the approximate boundaries between the water masses, have
 177 been defined manually by specifying polygons in the θ - S plane (Figure 12a).
 178 the separation interface has been characterized in such a way to present the
 179 water masses as objectively as possible on the basis of the values of θ , defined
 180 in the various studies cited in the introduction. The large seasonal variability
 181 of SAW, the intermittency of NACW as well as the occasional direct mixing
 182 of dense MWs with AW were taken into account during this process. Then
 183 each sample labeled on the θ - S plane is projected into the coordinate system,
 184 potential density and potential spicity (σ - π) (Figure 12b). The reason why the
 185 labeling was not directly done on the (σ - π) diagram is explained by the fact that

186 the characteristics of the water masses in the study area are defined in previous
 187 studies via the θ -S diagram and that the equivalent potential spicity properties
 188 will only be deduced after the projection of the labeled samples into the $(\sigma$ - π)
 189 plane.

190 The constructed training dataset is therefore $A = \{(\sigma_s, \pi_s, \lambda_s)\}_{s=1}^N$ where σ_s , π_s
 191 and λ_s are respectively the potential density anomaly, the potential spicity and
 192 the water mass label of a sample s at a given longitude, latitude and depth.
 193 The choice of the $(\sigma$ - π) diagram for this classification study is justified in the
 194 next part of this section (2.4.3). It's worth mentioning that the terminologies
 195 of spicity and spiciness are used by several authors with different definitions
 196 to describe a 'spice' type variable in physical oceanography. Some authors
 197 have chosen to derive such a variable, called potential spicity, so that its con-
 198 tours are orthogonal to those of potential density (Veronis, 1972; Huang et al.,
 199 2018). Other studies are based on the non-orthogonal functions, called spici-
 200 ness (Jackett and McDougall, 1985; Flament, 2002; McDougall and Krzysik,
 201 2015). In our study, the potential spicity (π) is calculated on the basis of its
 202 definition as a function whose contours are orthogonal to those of the potential
 203 density (Huang et al., 2018) via the MATLAB subroutine `gsw_pspi(SA, CT,`
 204 `pr)`, also provided by (Huang et al., 2018), where (SA, CT, pr) is the absolute
 205 salinity ($\text{g}\cdot\text{kg}^{-1}$), conservative temperature ($^{\circ}\text{C}$) and reference pressure (db)
 206 (<https://github.com/lanlankai/Spicity-JGR>). The pressure value of $pr =$
 207 0 (the sea surface pressure) was taken as a reference level. Another remark
 208 concerns the WIW and WMDW : No traces of these two water masses were
 209 detected in the glider transects. In these cases, they will be excluded from the
 210 training dataset to avoid the distortion of the classification results. The choice
 211 to eliminate WMDW from the study was based on the 12.85°C potential tem-
 212 perature isoline used by (Millot, 2014) as an unambiguous definition of WMDW.
 213 Thus $\lambda_s \in \{\text{SAW, NACW, MAW, LIW, TDW}\}$. In analogy with Millot (2009)
 214 and Millot (2014) we differentiate hereafter, for convenience, a lower-TDW from
 215 an upper-TDW that will behave more like WMDW and LIW, respectively.

216 2.4. K nearest neighbors classification

217 2.4.1. Problem statement

218 The classification using the nearest neighbor search (Cover and Hart, 1967;
 219 Fix and Hodges, 1989) is a well known decision procedure, non parametric for
 220 automatic learning. It is used in this study to evaluate the presence and preva-
 221 lence of each water mass sampled by the glider in the different transects. This
 222 method has been considered as one of the widely used classification algorithms
 223 owing to its simplicity and straightforward implementation. However, it has few
 224 shortcomings affecting its accuracy of classification (Gallego et al., 2022; Gou
 225 et al., 2022) which are discussed in sections 2.4.2 and 2.4.3. This classification
 226 technique has an objective of classification and attribution to a request point q
 227 belonging to a sample of observations Q , the class of the instance of training of
 228 the nearest neighbor based on a metric that define the similarity between ob-
 229 servations and classes of a training dataset A . Moreover, it is useful to consider
 230 more than one neighbor, so the technique is more commonly referred to as K
 231 nearest neighbors (Knn) classification where the K nearest neighbors are used
 232 to determine the class (Cunningham and Delany, 2007). Figure 3 visualizes the
 233 overview scheme for the proposed K nearest neighbors classification of water
 234 masses.

235 We suppose a supervised learning set of data $A = \{(\sigma_s, \pi_s, \lambda_s)\}_{s=1}^N$ as
 236 described previously. In the training step, the dataset A is simply stored with
 237 any explicit learning. In the inference step, for each request instance q belonging
 238 to the dataset $Q = \{(\sigma_j, \pi_j)\}_{j=1}^M$, a Knn search is done to get the K closest
 239 instances $N(\sigma_s, \pi_s) = \{(\sigma_s^{(i)}, \pi_s^{(i)}, \lambda_s^{(i)})\}_{i=1}^k$ which are the nearest to q on the
 240 basis of the metric d . Therefore, the predicted water mass label λ_p is obtained
 241 using a weighted combination of labels $(\lambda^{(i)} |^{(i=1..k)})$ based on the d metric as
 242 follows:

$$243 \quad (2) \quad \lambda_p = f(q, A) = \frac{\sum_{i=1}^k d^{-1}(q, (\sigma^{(i)}, \pi^{(i)})) \cdot \lambda^{(i)}}{\sum_{i=1}^k d^{-1}(q, (\sigma^{(i)}, \pi^{(i)}))}$$

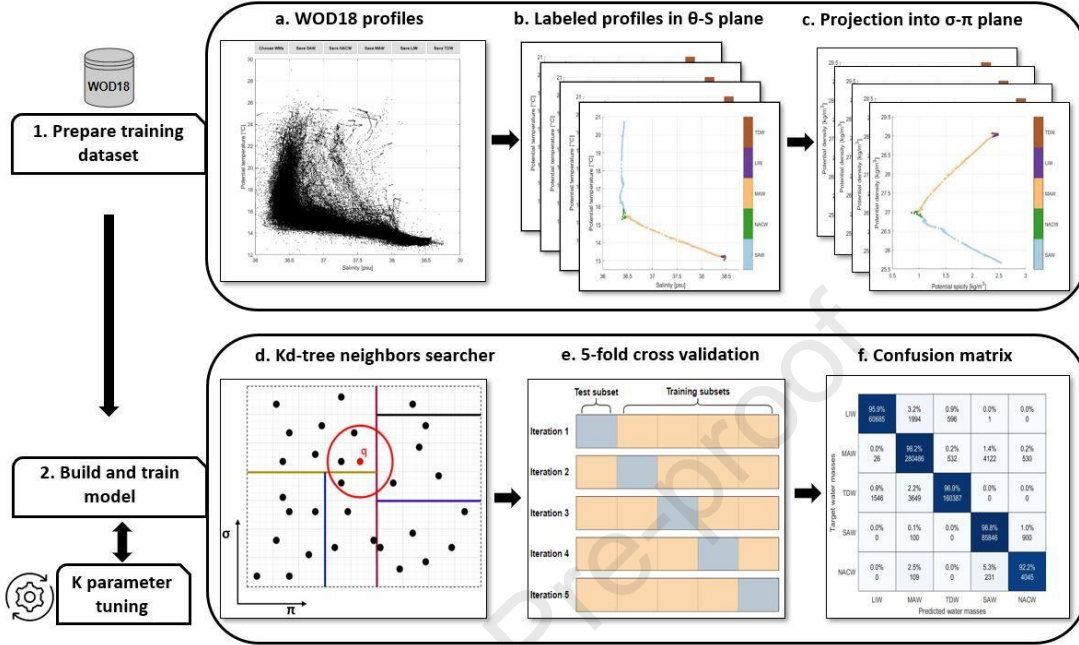


Figure 3: Flowchart of the proposed K nearest neighbors classification of water masses.

244 Thus, a Knn instance with a smaller distance will contribute more to the pre-
 245 diction for the instance.

246 In addition to classifying water masses into different categories, we can quan-
 247 tify the fraction of a given water mass λ_q in a request sample q belonging to Q
 248 as follows :

$$(3) \quad F_{q, \lambda_q} = \frac{\sum_{i=1}^m d^{-1}(q, (\sigma^{(i)}, \pi^{(i)}))}{\sum_{i=1}^k d^{-1}(q, (\sigma^{(i)}, \pi^{(i)}))}$$

249

250 where m is the number of samples representing the water mass λ_q among the K
 251 nearest neighbors on the basis of the metric d . Distances have been normalized
 252 to have all them lying between 0 and 1.

253 Such a quantification is helpful to supplement the information displayed in
 254 figures like 15b or 16b.

255 2.4.2. K nearest neighbors search

256 The simplest solution to the problem stated before remains on computing
 257 the d metric between the request point q and each point of the dataset and
 258 to return the k nearest points on the basis of d . The computing complexity
 259 is $O(N \times L)$, where N is the size of the data ensemble and L is its dimension
 260 (herein $L=2$). This method can be costly due to the huge amount of data.
 261 Within this context, numerous studies have been concerned with finding new
 262 approaches that are efficient with computations through employing fast search
 263 algorithms or using a training dataset size reduction scheme (Ougiaroglou and
 264 Evangelidis, 2016; Hou et al., 2018; Gallego et al., 2022). In our study, this
 265 drawback was overcome by searching for the K nearest neighbors using the
 266 spatial K dimensional tree subdivision structure (Kd-tree) (Chen et al., 2019).
 267 The latter is a well known optimisation for the Knn algorithm convenient for
 268 reduced dimensional spaces. The points ensemble N is divided recursively in
 269 the 2D space (σ - π) into a binary tree with N levels and $\log(N)$ depths.

270 This division continues until reaching at least a well defined number of points
 271 for each node. Therefore, the K nearest neighbors search for a point with a given
 272 request is done following these steps:

- 273 1. The determination of the node to which the query point belongs.
- 274 2. The search of the closest K points within that node on the basis of the
 275 metric d .
- 276 3. The determination of all other nodes having any area that is within the
 277 same metric d , in any direction, from the query point to the K^{th} closest
 278 point on the basis of the metric d .
- 279 4. The search of the closest K points within those nodes on the basis of the
 280 metric d .

281 2.4.3. Parameter definition and performances analysis

282 The Knn performances are known to be sensitive to choices of the metric
 283 and the parameter K which depend on the data characteristics (Jiang et al.,

284 2007). Therefore, they must be chosen appropriately to improve the classifica-
 285 tion performances. The metric selection can affect the form, the volume and
 286 the orientation of classes because some data points can be close for a metric
 287 and distant for another one. A small parameter K can capture a local structure
 288 in the data and therefore the result can be sensitive to noise, however a larger
 289 K permits to capture the global structure of data and suppress the noise effect
 290 but consumes more memory (Ghosh, 2006; Kang, 2021).

291 In this study, the chosen distance metric for the query points categorization
 292 is the Euclidean distance. Therefore, the metric d in equations 2 and 3 has the
 293 following form:

$$(4) \quad d(q, (\sigma^{(i)}, \pi^{(i)})) = \sqrt{(q - (\sigma^{(i)}, \pi^{(i)}))^2}$$

294 As mentioned in section 2.3, we choose the definition of potential spicity
 295 proposed by (Huang et al., 2018) who attempted to rehabilitate in the least
 296 square sense, the (Veronis, 1972) form of orthogonality between this variable
 297 and potential density. Thus, the choice of σ - π coordinates system instead of the
 298 traditional θ - S diagram is justified by the orthogonality and the dimensional
 299 homogeneity of these two pairs (σ and π). This allows a precise and concise
 300 measure of the distance d compared with θ - S diagram (Huang et al., 2018; Gao
 301 et al., 2020). Also, despite the existence of numerous techniques of data scaling,
 302 many authors have shown the impact of these techniques on the stability of ma-
 303 chine learning algorithms performances as it is the case for the Knn (Ambarwari
 304 et al., 2020; Shahriyari, 2017). Furthermore, one of the primary challenges is
 305 selecting the most suitable method for scaling. The latter problem is avoided
 306 here since σ and π share almost the same range.

307 Indeed, the major difference between the use of the two aforementioned
 308 diagrams is to determine boundaries between water masses. In our case, the
 309 previously described labeling method makes it possible to reduce this difference
 310 to 1.2%. However, to show the advantage of σ - π diagram for the computation
 311 of distance on which our method is based; the samples of the training dataset

312 forming the labeled boundaries of water masses have been eliminated in order
 313 to construct separate water masses in the two spaces σ - π and θ -S (figure 4).
 314 This situation represents the case of non-continuity of the training dataset or
 315 the case of difficulty to determine the boundaries characteristics between the
 316 water masses in a subjective way.

317 Taking as reference the classification results of the two transects (section
 318 3.2), we computed the total percentage of samples that changed membership
 319 from one water mass to another for the two diagrams. The results sketched in
 320 figures 5 and 6 represent a total difference of 10.38% for σ - π diagram versus
 321 17.7% for θ -S diagram. Therefore, using our methodology of classification, σ - π
 322 diagram is more appropriate for water masses frontiers determination. Also, the
 323 difference in distance calculation between the two spaces σ - π and θ -S is clearly
 324 visible in the computation of the fraction of a given water mass in a given sample
 325 (figures 7 and 8).

326 Concerning neighborhood size K selection, several methods have been de-
 327 veloped with a view to predict its optimal value and to overcome its sensitivity
 328 (Zhongguo et al., 2017; Zhang et al., 2018; Gou et al., 2019, 2022). In our case,
 329 the choice of the parameter K was based on the traditional L-Fold Cross Valida-
 330 tion method (Paik and Yang, 2004; Ghosh, 2006; Kang, 2021). This validation
 331 technique is based on estimating an accuracy rate for different values of K and
 332 select the one that induces the smallest classification error rate. The latter have
 333 been illustrated using the confusion matrix (Provost and Kohavi, 1998).

334 Indeed, the L-Fold Cross Validation consists in splitting the dataset
 335 $A = \{(\sigma_s, \pi_s, \lambda_s)\}_{s=1}^N$ in L independent subsets randomly selected with quasi-
 336 constant sizes. A subset is used to validate the produced model with the help of
 337 the L-1 remaining subsets. This process is applied L times so that each subset
 338 is used exactly one time for the validation. The classification error rate for all
 339 the partitions L is defined by a set $A' = \{(\sigma_s, \pi_s, \lambda_s)\}_{s=1}^{N'}$ $\in A$, as:

$$\tau = \frac{1}{N'} \sum_{i=1}^{N'} I_i \left(\lambda^{(i)} \neq \lambda_p^{(i)} \right) \quad (5)$$

340 where

$$I_i(\lambda^{(i)} \neq \lambda_p^{(i)}) = 0, \text{ if } \lambda^{(i)} = \lambda_p^{(i)} \\ = 1, \text{ otherwise}$$

341 Two popular choices of L are 5 and 10. In this study, we fix L to 5.

342 In practice, the total classification error rate for the set $A = (\sigma_s, \pi_s, \lambda_s)^{N_s=1}$
 343 is deduced from the confusion matrix. This matrix illustrates not only the
 344 algorithm errors but also how the classification algorithm works for each class
 345 (Markoulidakis et al., 2021). Indeed, the confusion matrix is a cross table where
 346 each column represents the predicted class instances, and each row represents
 347 the real class instances. The classes $\lambda_s \in \{\text{SAW, NACW, MAW, LIW, TDW}\}$,
 348 are listed in the same order in the rows and columns, so the correctly classified
 349 elements are located on the main diagonal. During the cross validation L-fold,
 350 if the predicted class of the test sample is correct, then the diagonal element
 351 of the confusion matrix is incremented by 1. However, if the predicted class
 352 is incorrect, then the element off diagonal is incremented by 1. Once, all the
 353 training samplings are classified, the classification error rate is based on the
 354 ratio of the number of sampling incorrectly classified and the total number of
 355 classified samplings.

356 Numerous evaluations of K between 10 and 100 recorded classification er-
 357 ror rates between 2% and 2.2%. The parameter k=51 seems to be a good
 358 compromise between the complexity and precision of computation. The multi-
 359 classes confusion matrix, a matrix of 5×5 dimension, relative to this value of
 360 K is sketched in Figure 9. This matrix is build from the cross validation 5-fold
 361 applied to a total number of sampling N=604855. The classes SAW; MAW;
 362 TDW record an accuracy beyond the total accuracy of 98%. For the case of the
 363 NACW, the algorithm classifies incorrectly almost 4% of the training samplings
 364 between the surface SAW and subsurface MAW layers. The LIW record the
 365 highest classification error rate. Indeed, 8% of the sampling that are supposed
 366 to belong to this class were confused with the classes MAW; TDW. The reason
 367 behind this is the relatively tight and sinuous relationship between the LIW,

368 MAW and TDW classes in regards to the θ -S and σ - π diagrams.

369 It is mentioned that the experimental environment of model building was
370 performed on a computer with an Intel i7-1165G7 CPU @2.80 GHz with 8
371 GB memory. For a total number of sampling $N=604855$ forming the train-
372 ing dataset, prediction speed was 44000 observations per second and the total
373 training time was 72.03 seconds.

374 2.5. Sensitivity of the method

375 A sensitivity analysis was performed to assess the impact of the spatio-
376 temporal distribution of the hydrological profiles forming the training dataset.
377 This is achieved by computing the percentage of samples that move from one
378 water mass to another, when spatio-temporal variability is reduced in the train-
379 ing dataset. The classification results of the two transects, using all profiles of
380 the database are taken as reference.

381 Regarding the spatial sensitivity and as the distribution of MWs in the Albo-
382 ran Sea mainly depends on latitude (e.g the presence of LIW in the northern 2/3
383 of the basin), the area has been divided into two regions separated by latitude
384 $35^{\circ}45'N$. Hydrological profiles of each region were used separately as a training
385 dataset to examine the impact of database spatial distribution on the classifi-
386 cation results of the proposed method. The results of this analysis is presented
387 in the following section.

388 The temporal sensitivity was examined to evaluate the impact of the tem-
389 poral ranges of the training dataset. This is achieved by dividing the training
390 data into profiles acquired during four periods, from 1950 to 1980, from 1950
391 to 1990, from 1950 to 2000 and from 1950 to 2010. Ocean profiles related to
392 each period were used separately as a training dataset. The confusion matrices
393 computed for these four cases (figure 10) showed that no significant changes oc-
394 cur. Therefore, the classification results are not altered by the temporal ranges
395 of the training data.

396 Also, temporal sensitivity of the seasonal variability of SAW was performed.
397 The training dataset was divided into profiles collected during fall, winter, spring

398 and summer seasons. The confusion matrices computed for these four cases
 399 showed less than 2% of difference between predicted SAW samples using the
 400 whole dataset and those predicted by using the separate seasonal data. Thus,
 401 seasonal variability of SAW does not influence the classification results.

402 3. Results

403 3.1. Water masses labeling in the σ - π plane

404 All the observation of the potential temperature and salinity obtained from
 405 the database used to build the training dataset and reaching a maximum depth
 406 of 700 m are sketched in figure 11. The seven water masses previously described
 407 can be distinguished as follows: the SAW are the lightest and characterized
 408 by a salinity layer quasi-homogeneous subject to intense seasonal variability
 409 and a constant temperature vertical gradient. The NACW is below the SAW
 410 and characterized by a salinity minimum with $\theta - S$ between 14°C-36 psu and
 411 16°C-36.4 psu. Under the Atlantic Waters, the $\theta - S$ diagram shows a linear
 412 stripe limited by the isopycnals $\sigma \approx 27.2 \text{ kg.m}^{-3}$ and $\sigma \approx 28.8 \text{ kg.m}^{-3}$. These
 413 values characterize the MAW resulting from the mixing between the Atlantic
 414 and Mediterranean Waters.

415 Beyond a salinity of 38 psu, the $\theta - S$ diagram is characterized by a tight
 416 and sinuous relationship, representing more than 80% of the total water volume.
 417 During its presence, the WIW is clearly noticed by its local temperature mini-
 418 mum (13°C-13.1°C, 38.25-38.35) linking the AW and the LIW. This water mass
 419 is characterized by a temperature and salinity local maximum shown by the
 420 $\theta - S$ diagram, with salinity values up to 38.58. The TDW is represented by a
 421 curved line linking the LIW and WMDW. This water mass is clearly indicated
 422 by its low temperature (< 12.9°C), its relatively low salinity (≈ 38.4) and its
 423 high density ($\approx 29.09 \text{ kg.m}^{-3}$).

424 Figure 12a sketches a part from the labelled training data in the coordinates
 425 system $\theta - S$. The labels are the water masses SAW ; NACW ; MAW ; LIW ; TDW.
 426 The equivalent result is projected on the $\sigma - \pi$ plan as shown in Figure 12b.

427 The latter shows the water masses characteristics which are clearly identified
428 through the analysis of potential spicity.

429 The general aspect of the water masses in the $\sigma - \pi$ plan are perceived as
430 a rotation transformation of the $\theta - S$ plan around the origin with an angle
431 $\alpha = 45^\circ$ (Figure 12b). Indeed, the Atlantic waters (SAW, NACW and MAW)
432 keep a geometric aspect of a curve as an elbow. The inflexion point of this curve
433 represents the interface between the surface waters (SAW) and those of the
434 subsurface (MAW). These waters are characterized by a linear relation defined
435 by positive and negative coefficients respectively. The NACW reveals a potential
436 spicity minimum $\pi = 0.45 \text{ kg.m}^{-3}$. The Mediterranean waters (LIW and TDW)
437 keep the aspect of broad relationship where the LIW is characterized by a local
438 maximum of potential spicity $\pi = 2.58 \text{ kg.m}^{-3}$, equivalent to a local maximum
439 of salinity $S = 38.52 \text{ psu}$.

440 3.2. Glider transects classification

441 As mentioned in section 2.2, the first days of the glider profiling were dedi-
442 cated to the survey of Moroccan Mediterranean waters offshore and more pre-
443 cisely of the Western Alboran Gyre (WAG), located between the strait of Gibrat-
444 tar and the Tres Forcas cape. This quasi-steady anticyclonic gyre has a typ-
445 ical diameter of approximately 100 km and a depth of 200 m and represents
446 the most intense dynamical structure of the mean circulation in the western
447 Mediterranean sea, with surface currents reaching 1.5 m s^{-1} (Álvaro Viúdez
448 et al., 1996; Vélez-Belchi et al., 2005; Flexas et al., 2006).

449 The vertical profiles of temperature and salinity acquired between 11 and 22
450 February 2021 during the first and second transects are represented in Figures
451 13 and 14 respectively.

452 The warm and fresh anomalies characterizing the WAG appear noticeably
453 in the temperature and salinity fields. Globally, the vertical distribution of
454 T is characterized by decreasing values with depth and by a sharp vertical
455 gradient. The WAG core highlights temperature values higher than 15°C and
456 positive anomaly compared to the ambient environment. The latter results in

457 a deepening of the isothermal layers by several tens of meters inside the WAG
 458 and an upwelling of these layers outside the WAG.

459 Despite the relatively long time period sampling of both transects (~ 7
 460 days for the first transect and ~ 5 days for the second transect), we consider
 461 a quasi-synoptic situation, highlighting the water mass composition during this
 462 period of time in the studied region. Therefore, we only consider the spatial
 463 mixing variability of the water mass.

464 The classification methodology applied on both glider transects (Figures 15b
 465 and 16b) shows that the AW engulfs the top layer (from surface to 200-250m)
 466 of the Alboran Sea, just below the isopycnal $\sigma = 28.9 \text{ kg.m}^{-3}$. This layer
 467 is characterized by a density anomaly which is the result of temperature and
 468 salinity anomalies. The isohaline light layer ($S < 36.6 \text{ psu}$, $\sigma < 27 \text{ kg.m}^{-3}$)
 469 is classified as a SAW in both transects. However, only the second transect
 470 outlines the presence of the NACW in its southern side near the Moroccan
 471 coast as highlighted by a spicity minimum (Figure 16b).

472 Beyond the isopycnal $\sigma = 28.9 \text{ kg.m}^{-3}$, we found the MWs adjoining the
 473 WAG and containing the LIW and the TDW. The LIW layer is absent in the
 474 south, near the Moroccan coasts, and is principally concentrated in the center
 475 and the north of the Alboran Sea; where it thickens. The TDW is mainly present
 476 along the two transects from the south to the north. Nevertheless, we distinguish
 477 an upper TDW which is found just below the LIW and a lower TDW which is
 478 located in the southern side below the AW. The $\sigma - \pi$ diagrams (Figures 15a and
 479 16a) outline the distribution of the TDW : In the southern part of the transect,
 480 the dense MWs (herein the lower TDW) are individually mixed with the AW
 481 leading to a relatively straight shape with some bending in the deep part of the
 482 profiles. However, far away from the African coasts the MWs are overlapped
 483 and slightly mixed leading to a sinuous shape in which the upper TDW tends
 484 to connect the LIW with dense MWs.

485 Thus, the classification of glider transects shows: (i) the formation of the
 486 WAG by the newly flushed AW, (ii) The presence of the LIW in the 2/3 North
 487 of the region, (iii) the presence of an upper TDW below the LIW and (iv) the

488 presence of a lower TDW in the southern side below the AW.

489 3.3. Training dataset sensitivity

490 In the case where only the profiles gathered below $35^{\circ}45'N$ are used as train-
491 ing dataset, the confusion matrix (Figure 17a) shows that the TDW is very well
492 classified while almost 30% of LIW samples move to TDW (26.2%) and MAW
493 (2.5%). The results of the new classification applied to the first glider transect
494 (Figure 18a) show that the algorithm captures the uplift of dense MWs in the
495 southern part of the basin and that the LIW layer is still present in the northern
496 2/3 of the transect but the latter becomes less thick. In the case where only
497 the profiles gathered beyond $35^{\circ}45'N$ are used as training data, the confusion
498 matrix (Figure 17b) shows that the LIW regains about 10% of its samples com-
499 pared to in the previous case, while almost 30% of TDW samples move toward
500 MAW (20.2%) and LIW (8.7%). The results of the new classification applied to
501 the first transect of the glider are shown in Figure 18b. The spatial distribution
502 of the LIW is close to that relating to the use of the entire database. However,
503 the uplift of dense MWs in the south is not well represented. Similar tests were
504 carried out on the data from the second transect with similar results.

505

506 3.4. Clustering analysis

507 To show the advantage of our method compared to those of unsupervised
508 classification, a cluster analysis based on the iterative algorithm k-means (Ap-
509 pendix B), classically used to specify water masses characteristics (Roseli et al.,
510 2015; Molleri et al., 2010), was applied on the $\sigma-\pi$ diagram to classify the water
511 mass in the both transects. As the k-means is also based on distance compu-
512 tation, we choose the $\sigma-\pi$ coordinate system to allow a concise measure of
513 this distance. The similarity between the samples and the centroids of the clus-
514 ters (selected randomly in the first step) is indicated by the euclidean distance
515 defined as in equation 3. The results obtained and the related analyzes being
516 the same for the two transects, we limit ourselves to the presentation of those

517 relating to the first transect. To ensure that the chosen number of clusters,
518 k , is representative of the system, different values of k between 2 and 5 were
519 tested (Figure 19). The silhouette method (Appendix B) is used as the tool to
520 validate the clustering quality and to see how well each sample lies within its
521 cluster. In this test, the positive silhouette value nearest to one, indicate cases
522 where a sample is well clustered. Samples with negative silhouette values are
523 considered as poorly classified.

524

525 The classification results show that the clustering analysis performs well to
526 distinguish the Atlantic and Mediterranean Waters for a k value greater than or
527 equal to 3 (Figures 19b, 19c and 19d). The separating interface between these
528 waters is formed by the 28.8 kg.m^{-3} isopycnal, which is approximately equal
529 to the value found by our method ($\sigma \sim 28.9 \text{ kg.m}^{-3}$). However, when different
530 AWs are considered, only the SAW can be distinguished by the algorithm. The
531 MWs (LIW, upper and lower TDW) are inhomogeneous and this results in a single
532 layer above the 28.8 kg.m^{-3} isopycnal.

533

534 The silhouette values show that when water masses are divided into 2 clusters
535 (Figure not shown) all the samples are correctly clustered, showing a positive
536 and significant silhouette values greater than the mean value (0.96 in this case).
537 Nevertheless, for the other number of clusters ($k=3, 4$ and 5), several samples
538 are wrongly grouped with a negative silhouette values (Figures not shown).
539 however in all this cases, clusters number 1 and 2 still well classified and this
540 explain the fact that (i) the k -means performs well for $k \geq 3$ in distinguishing
541 the Atlantic and Mediterranean Waters and that (ii) in the AWs, only the SAW
542 is well defined.

543 4. Discussion

544 4.1. The hydrographic structure of the WAG

545 The hydrographic structure of the WAG sampled through the glider shows
546 that the gyre vertical extent (~ 180 m) is characterized by a large homogeneous
547 layer in salinity with values lower than 36.6 psu. These results are in agree-
548 ment with other cruises that sampled the WAG at its usual location (as in our
549 study), in this case (Álvaro Viúdez et al., 1996; Nibani et al., 2021). Using data
550 from an intensive field experiments, these authors recorded the same salinity
551 characteristics of the isohaline layer, occupying the upper part of the gyre and
552 reported a typical vertical extension of 180-200m. However, in comparison with
553 oceanographic cruises coinciding with the eastward migration event of the WAG
554 (Vélez-Belchi et al., 2005; Flexas et al., 2006), the salinity within the gyre is
555 higher than that found in our study (up to 0.2 psu). This difference is explained
556 by the fact that in its usual location, the WAG is more exposed to inputs of
557 fresh AW than when it is located further east. Moreover, these authors recorded
558 a reduced vertical extensions of the WAG of 130-150m.

559 4.2. The classification results

560 As shown above, the water in the study region is classified into 5 types
561 via a Knn classification method based on σ - π diagram. The labeling process
562 result shows that the characteristics of the different water masses evolving in
563 the Alboran Sea can be clearly identified through the σ - π coordinate system.
564 Indeed, each water masses represents a physical property and a geometric aspect
565 that correlates with that of the traditional θ - S diagram but which can be
566 studied from a different angle.

567 The classification results obtained for the AW show that the latter is not
568 sensitive to the spatio-temporal variability of the training dataset. The core of
569 the WAG marked with a large vertical thickness of homogeneous salinity layer
570 (<36.6 psu), is principally generated by the SAW. This result is in good agree-
571 ment with previous studies (Álvaro Viúdez et al., 1996; Vélez-Belchi et al., 2005;

572 Flexas et al., 2006). These authors show through a three dimensional descrip-
573 tion of the Western Alboran Sea that the WAG is characterized by recent AW
574 transported from the Strait of Gibraltar into the core of the gyre and occupy-
575 ing a considerable part of it. The 28.9 kg.m^{-3} isopycnal, found as separating
576 interface between AW and MWs, corresponds to that deduced from the $\theta - S$
577 diagram analysis by Gascard and Richez (1985) in their study of water masses
578 and circulation in the western Alboran Sea. The no significant NACW samples
579 detected during the second transect can be interpreted as points being closer to
580 NACW than to any other water mass (in this case SAW and MAW) and not as
581 samples marking the pure NACW.

582 The vertical distributions result of the MWs in both transects is sensitive
583 to the spatial variability of the training dataset. By using the whole labeling
584 data, the obtained result is in agreement with those inferred from the expert
585 analysis of the $\theta - S$ diagram. The spatial distribution of the LIW layer that
586 thickens from the north to the south corroborates with the works of Parrilla
587 et al. (1986); Millot (2014). These authors found that the properties of the LIW
588 is quite recognizable in most of the Alboran Sea, except in the southernmost
589 part near the Moroccan coasts. They showed that the path of the LIW in the
590 Alboran basin did not cross south of $35^{\circ}30'N$. The LIW limits obtained by our
591 classification method are $35^{\circ}30'N$ and $35^{\circ}45'N$ for the first and second transect
592 respectively. The spatial distribution of the TDW (the upper TDW and the
593 lower TDW) is in total concordance with the studies of Millot (2009, 2014) and
594 highlights the direct link between the deep MWs and the AW in the southern
595 side of the two transects. Indeed Millot (2009, 2014), shows through a $\theta - S$
596 diagram analysis of zonal hydrographic transects that, in southern part of the
597 Alboran Sea, dense MWs mixes directly with AW.

598 4.3. Comparison with clustering analysis

599 The comparison between the method adopted in this paper and the k-means
600 algorithm shows that this latter can not distinguish water masses when several
601 AWs and MWs are considered. In fact, the uplift of the dense MWs in the

602 southern part and the presence of the LIW in the $\frac{2}{3}$ parts of the northern basin
603 can not be outlined by the clustering analysis. This comparison corroborates
604 the performed analysis by Cheng et al. (2014); Millot (2019). Indeed, Cheng
605 et al. (2014) shows that in a well-defined range of potential density, water masses
606 having similarities in temperature and salinity are inseparable by the clustering
607 analysis. Millot (2019) shows that the method proposed by Naranjo et al. (2015)
608 is rather a computation of euclidean distances between the samples and a set of
609 centroids representing the water masses than a clustering analysis. He concludes
610 that, in regions of relatively moderate mixing processes such as in the Strait of
611 Gibraltar, a subjective (θ -S) diagram analysis based on a traditional method
612 where boundaries of water masses are defined by experts experience, is much
613 more robust than clustering analysis.

614 Thus, conventional cluster analysis are not always appropriate to discrimi-
615 nate water masses and there is no clear physical meanings of the water masses
616 boundaries. As the labeling process guides the decision of the algorithm to-
617 wards the choice of the water mass representing each sample, our methodology
618 retains a part of this physical meaning through a labeling approach based on a
619 traditional method for defining boundaries between water masses.

620 5. Conclusions

621 The objective of this study was to identify the spatial distribution of water
622 masses in the Western basin of the Alboran Sea. To do this, a novel method-
623 ology based on water masses automatic classification using the Knn search was
624 applied to the T-S data acquired by a glider. These data have been projected
625 on the orthogonal and dimensional homogeneous coordinates system: potential
626 density anomaly-potential spicity (σ - π).The parameters used in this algorithm
627 have been selected in order to get the most accurate classification. The latter
628 have been insured by a supervised machine learning process based on available
629 data from the World Ocean Database 18 and the Global Data Assembly Center
630 (Figure 2). From all the water masses described in section 1, the WIW and

631 the WMDW were not successfully detected by the glider and therefore were ex-
632 cluded from the training dataset. Thus, the water masses in the glider transects
633 were classified in 5 categories: SAW; NACW; NAW; LIW and TDW.

634

635 In comparison to the classic method of classification based on clustering
636 analysis (herein the k-means), the proposed method in this paper permits to
637 ascertain the water masses frontiers with a reasonable and robust approach.
638 In the studied region, the classification results are in good agreement with the
639 circulation schemes established in previous studies and inferred from the tra-
640 ditional method based on the subjective expert analysis of the (θ -S) diagram,
641 showing:

- 642 • The formation of the WAG by the recently advected Atlantic Water (Álvaro
643 Viúdez et al., 1996; Flexas et al., 2006) ;
- 644 • The uplift of the dense Mediterranean Waters (the lower-TDW) near the
645 Moroccan coasts (Millot, 2009, 2014) ;
- 646 • The presence of the LIW in the 2/3 North of the Western basin of the
647 Alboran Sea (Parrilla et al., 1986; Millot, 2009).

648 The application of our approach for ocean water masses classification has
649 many advantages. By combining traditional method based on expert analysis
650 and Machine learning technique, this methodology is useful and appropriate
651 to automatically classify water masses in regions where intense mixing occurs
652 such as the Western Alboran Sea. Although the labeling process requires the
653 knowledge of the water masses characteristic in the study area, the adaptation
654 of this technique to other regions is easy and straightforward. Indeed, this
655 methodology can be applied easily to other sub-basins or marginal seas as long as
656 a sufficient number of in-situ observations describing the whole spatio-temporal
657 variability of the area can be provided as a training dataset.

658 The speed of the proposed method will make it possible on the basis of ba-
659 sic hydrographic data collected during typical research cruises or autonomous

660 systems, to provide classification results in real time. Remarkably, Using the
661 proposed methodology, researchers non-particularly specialists in oceanography,
662 can take advantage of previous knowledge of water masses characteristics val-
663 idated by experts to solve the problem of water masses classification. Within
664 this context, a Graphical User Interface (GUI) is under development in order to
665 enable users performing the entire process described in this manuscript (figure
666 3), within all ocean basins.

667 Acknowledgment

668 The authors gratefully acknowledge Houssini Nibani, President of AGIR As-
669 sociation who has performed this glider mission in the Western Aboran Sea
670 in the framework of the European project ODYSSEA. We are also grateful to
671 Laurent Beguery and Orens Fommervault from Alseamar France For their coop-
672 eration and their answers to the questions concerning the glider data collected in
673 this region. Thanks to the anonymous reviewers for their constructive comments
674 and helpful suggestions.

675 **Appendix A. Examples of using K Nearest Neighbors Classification**
676 **to study the spatial distribution of water mass in the**
677 **Western Alboran Sea.**

678 Other examples of data acquired over different time period have been used to
679 study the spatial distribution of water mass in the Western Alboran Sea (figure
680 A.20). It's about:

- 681 • A CTD transect of a field experiment acquired in September 1992 on board
682 of the R/V Garcia del Cid (Álvaro Viúdez et al., 1996). Being available
683 on WOD18, these data have been removed from the training dataset to
684 assess the results of the classification in a more objective way.
- 685 • A hydrographic (CTD) cast of an intensive oceanographic survey (BIOMEGA)
686 collected on board of the Spanish R/V Garcia del Cid during October
687 2003 (Flexas et al., 2006). Data were provided through SeaDataNet Pan-
688 European infrastructure for ocean and marine data management (<https://www.seadatanet.org>);
689
- 690 • A glider transect (from 11 to 17 November 2020) of the first mission per-
691 formed by the Moroccan association AGIR (Nibani et al., 2021), as part
692 of the European project ODYSSEA (<https://odysseaplatform.eu/fr/home-fr/>).

694 The temperature and salinity fields, and the classification results obtained
695 for the three aforementioned oceanographic cruises, are sketched in figures A.21,
696 A.22 and A.23. This leads to the same interpretation of glider data previously
697 described in section 2.2.

698 **Appendix B. k-means clustering and silhouette method.**

699 k-means is one of the simplest unsupervised learning algorithms that solve
700 the well known clustering problem (Kaufman and Rousseeuw, 1990). It is an
701 iterative, data-partitioning algorithm which aims to partition n observations
702 into k groups, called clusters. The algorithm proceeds as follows :

- 703 1. Select k initial centroids at random after indicating the desired k number
704 of clusters ;
- 705 2. Compute sample-to-cluster-centroid distances of all observations to each
706 centroid and then assign each observation to the cluster with the closest
707 centroid ;
- 708 3. Compute the average of the observations in each cluster to obtain k new
709 centroid locations;
- 710 4. Repeat steps 2 and 3 until cluster assignments do not change, or the
711 maximum number of iterations is reached.

712 k-means aims at minimizing an objective function that depends on the dis-
713 tance of the data points to the cluster centroids. Suppose $D = \{x_1, \dots, x_n\}$ is
714 the dataset to be clustered. K-means problem can be expressed as follows :

$$\min \sum_{k=1}^K \sum_{x \in C_k} f(x, c_k) \quad (\text{B.1})$$

The function 'f' computes the distance between object x and centroid c_k which is defined by:

$$c_k = \sum_{x \in C_k} \frac{x}{n_k} \quad (\text{B.2})$$

715 where n_k is the number of data objects assigned to cluster C_k .

716

717 To evaluate the clustering analysis quality, (Rousseeuw, 1987) introduced the
718 so-called silhouette method. This technique provides a graphical representation
719 which helps the user to select the number of clusters and to see how well each
720 sample lies within its cluster. The silhouette value for each sample is a measure

721 of how similar that sample is to other samples in the same cluster, compared to
 722 samples in other clusters. The silhouette value s_i for the i^{th} sample is defined
 723 as :

$$s_i = \frac{(b_i - a_i)}{\max(a_i, b_i)} \quad (\text{B.3})$$

724
 725 where a_i is the average dissimilarity of the i^{th} sample with all other data within
 726 the same cluster and b_i is the minimum average dissimilarity of the i^{th} sample
 727 to samples in a different cluster. Distance metric is employed to calculate the
 728 dissimilarity between samples. When a cluster contains only a single sample, it
 729 is unclear how a_i should be defined and then s_i is set to 1.

730 Indeed, from the preceding definition, it is clear that $-1 \leq s_i \leq 1$ for each
 731 sample i . A high and positive value indicates that the sample is well matched
 732 to its own cluster, and distant from neighboring clusters. A low or negative
 733 silhouette value, correspond to cases in which samples are assigned to wrong
 734 clusters.

735 References

- 736 Ambarwari, A., Jafar Adrian, Q., Herdiyeni, Y., 2020. Analysis of the effect of
737 data scaling on the performance of the machine learning algorithm for plant
738 identification. *Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi)* 4,
739 117–122. URL: [https://jurnal.iaii.or.id/index.php/RESTI/article/
740 view/1517](https://jurnal.iaii.or.id/index.php/RESTI/article/view/1517), doi:10.29207/resti.v4i1.1517.
- 741 Argo, 2021. Argo float data and metadata from global data assembly
742 centre (argo gdac). URL: <https://www.seanoe.org/data/00311/42182/>,
743 doi:<https://doi.org/10.17882/42182>.
- 744 Bosse, A., Testor, P., Mortier, L., Prieur, L., Taillandier, V., d’Ortenzio,
745 F., Coppola, L., 2015. Spreading of levantine intermediate waters by
746 submesoscale coherent vortices in the northwestern mediterranean sea
747 as observed with gliders. *Journal of Geophysical Research: Oceans* 120,
748 1599–1622. URL: [https://agupubs.onlinelibrary.wiley.com/doi/abs/
749 10.1002/2014JC010263](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2014JC010263), doi:<https://doi.org/10.1002/2014JC010263>,
750 arXiv:<https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1002/2014JC010263>.
- 751 Boyer, T., Baranova, O., Locarnini, R., Mishonov, A., Grodsky, A., Paver,
752 C., Weathers, K., Smolyar, I., Reagan, J., Seidov, D., Zweng, M., 2019.
753 World ocean atlas 2018 product documentation ocean climate laboratory ncei
754 / nesdis / noaa noaa national centers for environmental information. doi:10.
755 13140/RG.2.2.34758.01602.
- 756 Broecker, W., 1991. The great ocean conveyor. *Oceanography* URL: [https:
757 //doi.org/10.5670/oceanog.1991.07](https://doi.org/10.5670/oceanog.1991.07).
- 758 Bryden, H.L., HL, B., HM, S., 1982. Origin of the mediterranean outflow .
- 759 Chen, Y., Zhou, L., Tang, Y., Singh, J.P., Bouguila, N., Wang, C., Wang, H.,
760 Du, J., 2019. Fast neighbor search by using revised k-d tree. *Information
761 Sciences* 472, 145–162. URL: <https://www.sciencedirect.com/science/>

- 762 article/pii/S0020025518307126, doi:<https://doi.org/10.1016/j.ins.>
763 2018.09.012.
- 764 Cheng, G.S., Sun, J.D., Zu, T.T., Chen, J., Wang, D.X., 2014. Analysis of water
765 masses in the northern south china sea in summer 2011. *Journal of Tropical*
766 *Oceanography* 33, 10–16.
- 767 Clark, P., Pisias, N., Stocker, T., Weaver, A., 2002. The role of the thermoha-
768 line circulation in abrupt climate change. *Nature* 415, 863–9. doi:[10.1038/](https://doi.org/10.1038/415863a)
769 415863a.
- 770 Cochrane, J.D., 1958. The frequency distribution of water characteristics in
771 the pacific ocean. *Deep Sea Research* (1953) 5, 111–127. URL: <https://www.sciencedirect.com/science/article/pii/0146631358900029>,
772 doi:[https://doi.org/10.1016/0146-6313\(58\)90002-9](https://doi.org/10.1016/0146-6313(58)90002-9).
- 774 Cover, T., Hart, P., 1967. Nearest neighbor pattern classification. *IEEE Trans-*
775 *actions on Information Theory* 13, 21–27. doi:[10.1109/TIT.1967.1053964](https://doi.org/10.1109/TIT.1967.1053964).
- 776 Cunningham, P., Delany, S.J., 2007. k-nearest neighbour classifiers.
- 777 Fix, E., Hodges, J.L., 1989. Discriminatory analysis. nonparametric discrim-
778 ination: Consistency properties. *International Statistical Review / Revue*
779 *Internationale de Statistique* 57, 238–247. URL: [http://www.jstor.org/](http://www.jstor.org/stable/1403797)
780 [stable/1403797](http://www.jstor.org/stable/1403797).
- 781 Flament, P., 2002. A state variable for characterizing water masses and their
782 diffusive stability: Spiciness. *Progress in Oceanography* 54, 493–501.
- 783 Flexas, M., Gomis, D., Ruiz, S., Pascual, A., Leon, P., 2006.
784 In situ and satellite observations of the eastward migration of the
785 western alboran sea gyre. *Progress in Oceanography* 70, 486–
786 509. URL: [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0079661106000590)
787 [S0079661106000590](https://www.sciencedirect.com/science/article/pii/S0079661106000590), doi:[https://doi.org/10.1016/j.pocean.2006.03.](https://doi.org/10.1016/j.pocean.2006.03.017)
788 017. gabriel T. Csanady: *Understanding the Physics of the Ocean*.

- 789 Gallego, A.J., Rico-Juan, J.R., Valero-Mas, J.J., 2022. Efficient k-nearest neighbor
790 search based on clustering and adaptive k values. *Pattern Recognition*
791 122, 108356. URL: [https://www.sciencedirect.com/science/article/
792 pii/S0031320321005367](https://www.sciencedirect.com/science/article/pii/S0031320321005367), doi:[https://doi.org/10.1016/j.patcog.2021.
793 108356](https://doi.org/10.1016/j.patcog.2021.108356).
- 794 Gao, Y., Huang, R.X., Zhu, J., Huang, Y., Hu, J., 2020. Using the sigma-
795 pi diagram to analyze water masses in the northern south china sea in
796 spring. *Journal of Geophysical Research: Oceans* 125, e2019JC015676.
797 URL: [https://agupubs.onlinelibrary.wiley.com/doi/abs/10.
798 1029/2019JC015676](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2019JC015676), doi:<https://doi.org/10.1029/2019JC015676>,
799 arXiv:<https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2019JC015676>.
800 e2019JC015676 2019JC015676.
- 801 Gascard, J., Richez, C., 1985. Water masses and circulation in
802 the western alboran sea and in the straits of gibraltar. *Progress*
803 *in Oceanography* 15, 157–216. URL: [https://www.sciencedirect.
804 com/science/article/pii/007966118590031X](https://www.sciencedirect.com/science/article/pii/007966118590031X), doi:[https://doi.org/10.
805 1016/0079-6611\(85\)90031-X](https://doi.org/10.1016/0079-6611(85)90031-X).
- 806 Ghosh, A., 2006. On optimum choice of k in nearest neighbor classification.
807 *Computational Statistics Data Analysis* 50, 3113–3123. doi:[10.1016/j.
808 csda.2005.06.007](https://doi.org/10.1016/j.csda.2005.06.007).
- 809 Gou, J., Ma, H., Ou, W., Zeng, S., Rao, Y., Yang, H., 2019. A
810 generalized mean distance-based k-nearest neighbor classifier. *Ex-
811 pert Systems with Applications* 115, 356–372. URL: [https:
812 //www.sciencedirect.com/science/article/pii/S0957417418305293](https://www.sciencedirect.com/science/article/pii/S0957417418305293),
813 doi:<https://doi.org/10.1016/j.eswa.2018.08.021>.
- 814 Gou, J., Sun, L., Du, L., Ma, H., Xiong, T., Ou, W., Zhan, Y., 2022.
815 A representation coefficient-based k-nearest centroid neighbor classi-
816 fier. *Expert Systems with Applications* 194, 116529. URL: [https:
817 //www.sciencedirect.com/science/article/pii/S0957417422005293](https://www.sciencedirect.com/science/article/pii/S0957417422005293).

- 817 //www.sciencedirect.com/science/article/pii/S095741742200288,
818 doi:<https://doi.org/10.1016/j.eswa.2022.116529>.
- 819 Hansen, B., 1916. Nogen hydrografiske metoder. Forh Skand Naturf Mote 16,
820 357–359.
- 821 Hansen, B.H., Nansen, F., . The norwegian sea : its physical oceanography
822 based upon the norwegian researches 1900-1904.
- 823 Hou, W., Li, D., Xu, C., Zhang, H., Li, T., 2018. An advanced k nearest
824 neighbor classification algorithm based on kd-tree, in: 2018 IEEE Interna-
825 tional Conference of Safety Produce Informatization (IICSPI), pp. 902–905.
826 doi:[10.1109/IICSPI.2018.8690508](https://doi.org/10.1109/IICSPI.2018.8690508).
- 827 Huang, R.X., Yu, L.S., Zhou, S.Q., 2018. New definition of potential spicity
828 by the least square method. Journal of Geophysical Research: Oceans 123,
829 7351–7365. URL: [https://agupubs.onlinelibrary.wiley.com/doi/abs/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2018JC014306)
830 [10.1029/2018JC014306](https://doi.org/10.1029/2018JC014306), doi:<https://doi.org/10.1029/2018JC014306>,
831 arXiv:<https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2018JC014306>.
- 832 Hur, H.B., Jacobs, G.A., Teague, W.J., 1999. Monthly variations of water
833 masses in the yellow and east china seas, november 6, 1998. Journal of
834 Oceanography 55, 171–184.
- 835 Jackett, D.R., McDougall, T.J., 1985. An oceanographic variable for the charac-
836 terization of intrusions and water masses. Deep Sea Research Part A. Oceano-
837 graphic Research Papers 32, 1195–1207.
- 838 Jiang, L., Cai, Z., Wang, D., Jiang, S., 2007. Survey of improving k-nearest-
839 neighbor for classification. Fourth International Conference on Fuzzy Systems
840 and Knowledge Discovery (FSKD 2007) 1, 679–683.
- 841 Kang, S., 2021. k-nearest neighbor learning with graph neural networks. Mathe-
842 matics 9. URL: <https://www.mdpi.com/2227-7390/9/8/830>, doi:[10.3390/](https://doi.org/10.3390/math9080830)
843 [math9080830](https://doi.org/10.3390/math9080830).

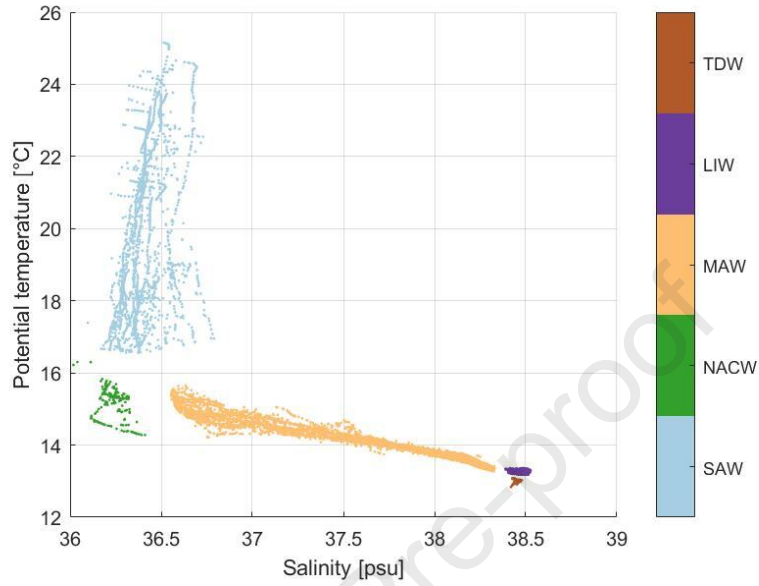
- 844 Kaufman, L., Rousseeuw, P.J., 1990. Finding Groups in Data: An Introduction
845 to Cluster Analysis. John Wiley.
- 846 Kim, K., Kim, K.R., Rhee, T., Rho, H., Limeburner, R., Beardsley, R.,
847 1991. Identification of water masses in the yellow sea and the east
848 china sea by cluster analysis, in: Takano, K. (Ed.), Oceanography of
849 Asian Marginal Seas. Elsevier. volume 54 of Elsevier Oceanography Series,
850 pp. 253–267. URL: [https://www.sciencedirect.com/science/article/
851 pii/S0422989408701004](https://www.sciencedirect.com/science/article/pii/S0422989408701004), doi:[https://doi.org/10.1016/S0422-9894\(08\)
852 70100-4](https://doi.org/10.1016/S0422-9894(08)70100-4).
- 853 Markoulidakis, I., Rallis, I., Georgoulas, I., Kopsiaftis, G., Doulamis, A.,
854 Doulamis, N., 2021. Multiclass confusion matrix reduction method and
855 its application on net promoter score classification problem. Technolo-
856 gies 9. URL: <https://www.mdpi.com/2227-7080/9/4/81>, doi:10.3390/
857 technologies9040081.
- 858 McDougall, T.J., Krzysik, O.A., 2015. Spiciness. Journal of Marine Research
859 73, 141–152.
- 860 Miller, A., Stanley, R., 1961. Volumetric ts diagrams for the mediterranean sea.
861 ATLANTIS 263, 6009–6146.
- 862 Millot, C., 2009. Another description of the mediterranean sea outflow.
863 Progress in Oceanography - PROG OCEANOGR 82, 101–124. doi:10.1016/
864 j.pocean.2009.04.016.
- 865 Millot, C., 2014. Heterogeneities of in- and out-flows in the mediter-
866 ranean sea. Progress in Oceanography 120, 254–278. URL: [https:
867 //www.sciencedirect.com/science/article/pii/S0079661113001857](https://www.sciencedirect.com/science/article/pii/S0079661113001857),
868 doi:<https://doi.org/10.1016/j.pocean.2013.09.007>.
- 869 Millot, C., 2019. Comments on computations about the mediterranean outflow
870 composition. Bollettino di Geofisica Teorica ed Applicata 60.

- 871 Millot, C., Candela, J., Fuda, J.L., Tber, Y., 2006. Large warming and
872 salinification of the mediterranean outflow due to changes in its compo-
873 sition. *Deep Sea Research Part I: Oceanographic Research Papers* 53,
874 656–666. URL: [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0967063706000112)
875 [S0967063706000112](https://www.sciencedirect.com/science/article/pii/S0967063706000112), doi:<https://doi.org/10.1016/j.dsr.2005.12.017>.
- 876 Molleri, G.S.F., Kampel, M., de Moraes Novo, E.M.L., 2010. Spectral classifi-
877 cation of water masses under the influence of the amazon river plume. *Acta*
878 *Oceanologica Sinica*, URL: [http://www.aosocean.com/en/article/doi/](http://www.aosocean.com/en/article/doi/10.1007/s13131-010-0031-1)
879 [10.1007/s13131-010-0031-1](http://www.aosocean.com/en/article/doi/10.1007/s13131-010-0031-1), doi:[10.1007/s13131-010-0031-1](https://doi.org/10.1007/s13131-010-0031-1).
- 880 Montgomery, R., 1958. Water characteristics of atlantic ocean and of
881 world ocean. *Deep Sea Research (1953)* 5, 134–148. URL: [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/0146631358900042)
882 [0146631358900042](https://www.sciencedirect.com/science/article/pii/0146631358900042),
883 doi:[https://doi.org/10.1016/0146-6313\(58\)90004-2](https://doi.org/10.1016/0146-6313(58)90004-2).
- 884 Naranjo, C., Sammartino, S., García-Lafuente, J., Bellanco, M.J., Taupier-
885 Letage, I., 2015. Mediterranean waters along and across the strait
886 of gibraltar, characterization and zonal modification. *Deep Sea Re-*
887 *search Part I: Oceanographic Research Papers* 105, 41–52. URL: [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0967063715001399)
888 [S0967063715001399](https://www.sciencedirect.com/science/article/pii/S0967063715001399),
889 doi:<https://doi.org/10.1016/j.dsr.2015.08.003>.
- 890 Nibani, H., Hilmi, K., Damghi, A., Beguery, L., Fommervault, O., Amhaoach,
891 Z., 2021. AL HOCEIMA LAUNCHES ITS FIRST FUNCTIONAL MARINE
892 OBSERVATORY IN NORTH AFRICA, in: 9th EuroGOOS International
893 conference, Shom and Ifremer and EuroGOOS AISBL, Brest, France. URL:
894 <https://hal.archives-ouvertes.fr/hal-03331067>.
- 895 Ougiaroglou, S., Evangelidis, G., 2016. RHC: a non-parametric cluster-based
896 data reduction for efficient k-nn classification. *Pattern Anal. Appl.* 19,
897 93–109. URL: <http://dx.doi.org/10.1007/s10044-014-0393-7>, doi:[10.1007/s10044-014-0393-7](https://doi.org/10.1007/s10044-014-0393-7).
898

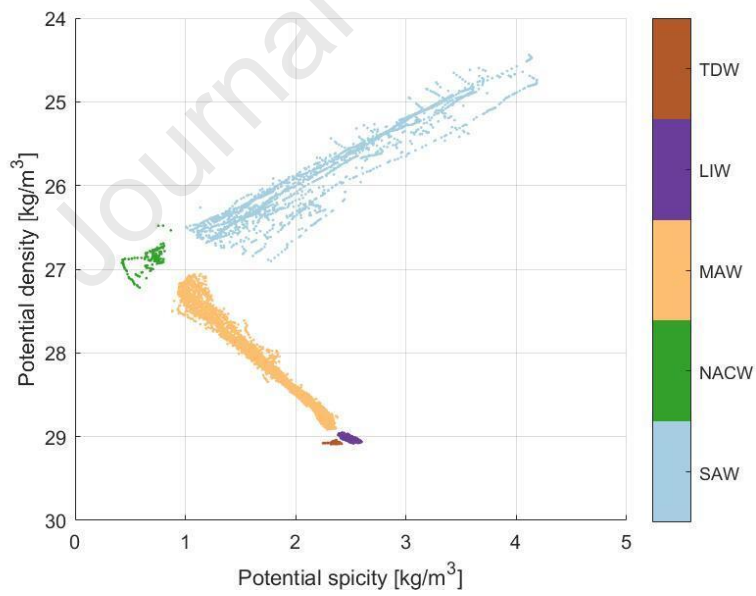
- 899 Paik, M., Yang, Y., 2004. Combining nearest neighbor classifiers versus
900 cross-validation selection. *Statistical Applications in Genetics and Molecu-*
901 *lar Biology* 3. URL: <https://doi.org/10.2202/1544-6115.1054>, doi:doi:
902 10.2202/1544-6115.1054.
- 903 Pantiulin, A.N., 2002. Water masses: birth of the idea URL:
904 [https://ices-library.figshare.com/articles/report/Theme_](https://ices-library.figshare.com/articles/report/Theme_Session_on_ACFM_and_Assessment_Working_Group_Reports/19258325)
905 [Session_on_ACFM_and_Assessment_Working_Group_Reports/19258325](https://ices-library.figshare.com/articles/report/Theme_Session_on_ACFM_and_Assessment_Working_Group_Reports/19258325),
906 doi:10.17895/ices.pub.8512.
- 907 Parrilla, G., Kinder, T., 1987. The physical oceanography of the alboran sea
908 184, 31.
- 909 Parrilla, G., Kinder, T.H., Preller, R.H., 1986. Deep and intermediate mediter-
910 ranean water in the western alboran sea. *Deep Sea Research Part A. Oceanographic Research Papers* 33, 55–88. URL: <https://www.sciencedirect.com/science/article/pii/0198014986901081>, doi:[https://doi.org/10.1016/0198-0149\(86\)90108-1](https://doi.org/10.1016/0198-0149(86)90108-1).
- 914 Pistek, P., De Strobel, F., Montanari, C., 1985. Deep-sea circulation in the
915 alboran sea. *Journal of Geophysical Research: Oceans* 90, 4969–4976.
- 916 Pollak, M., 1958. Frequency distribution of potential temperatures and salini-
917 ties in the indian ocean. *Deep Sea Research (1953)* 5, 128–133. URL: <https://www.sciencedirect.com/science/article/pii/0146631358900030>,
918 doi:[https://doi.org/10.1016/0146-6313\(58\)90003-0](https://doi.org/10.1016/0146-6313(58)90003-0).
- 920 Provost, F., Kohavi, R., 1998. Guest editors' introduction: On applied research
921 in machine learning. *Machine learning* 30, 127–132.
- 922 Renault, L., Oguz, T., Pascual, A., Vizoso, G., Tintore, J., 2012. Sur-
923 face circulation in the alborán sea (western mediterranean) inferred
924 from remotely sensed data. *Journal of Geophysical Research: Oceans*
925 117. URL: <https://agupubs.onlinelibrary.wiley.com/doi/abs/>

- 926 10.1029/2011JC007659, doi:<https://doi.org/10.1029/2011JC007659>,
927 arXiv:<https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1029/2011JC007659>.
- 928 Roseli, N.H., Akhir, M.F., Husain, M., Tangang, F., Ali, A., 2015. Water mass
929 characteristics and stratification at the shallow sunda shelf of southern south
930 china sea. *Open Journal of Marine Science* 05, 455–467. doi:10.4236/ojms.
931 2015.54036.
- 932 Rousseeuw, P., 1987. Silhouettes: a graphical aid to the interpretation and val-
933 idation of cluster analysis. *J. Comput. Appl. Math.* 20, 53–65. URL: <http://portal.acm.org/citation.cfm?id=38772>, doi:[http://dx.doi.org/10.](http://dx.doi.org/10.1016/0377-0427(87)90125-7)
934 [1016/0377-0427\(87\)90125-7](http://dx.doi.org/10.1016/0377-0427(87)90125-7).
935
- 936 Shahriyari, L., 2017. Effect of normalization methods on the performance of
937 supervised learning algorithms applied to htseq-fpkm-uq data sets: 7sk rna
938 expression as a predictor of survival in patients with colon adenocarcinoma.
939 *Briefings in bioinformatics* 20. doi:10.1093/bib/bbx153.
- 940 Veronis, G., 1972. On properties of seawater defined by temperature, salinity,
941 and pressure. *Journal of Marine Research* 30, 227–255.
- 942 Álvaro Viúdez, Tintoré, J., Haney, R.L., 1996. Circulation in the alboran sea
943 as determined by quasi-synoptic hydrographic observations. part i: Three-
944 dimensional structure of the two anticyclonic gyres. *Journal of Physical*
945 *Oceanography* 26, 684 – 705. URL: [https://journals.ametsoc.org/view/](https://journals.ametsoc.org/view/journals/phoc/26/5/1520-0485_1996_026_0684_citasa_2_0_co_2.xml)
946 [journals/phoc/26/5/1520-0485_1996_026_0684_citasa_2_0_co_2.xml](https://journals.ametsoc.org/view/journals/phoc/26/5/1520-0485_1996_026_0684_citasa_2_0_co_2.xml),
947 doi:10.1175/1520-0485(1996)026<0684:CITASA>2.0.CO;2.
- 948 Vélez-Belchi, P., Vargas-Yáñez, M., Tintoré, J., 2005. Observation of a western
949 alboran gyre migration event. *Progress in Oceanography* 66, 190–210. doi:10.
950 1016/j.pocean.2004.09.006.
- 951 Worthington, L.V., 1981. The water masses of the world ocean : Some results
952 of a fine-scale census.

- 953 Wright, W., Worthington, L., 1970. Serial atlas of the marine environment.
954 folio 19. the water masses of the north atlantic ocean: a volumetric census
955 of temperature and salinity. Technical Report. Woods hole oceanographic
956 institution ma.
- 957 Zhang, S., Li, X., Zong, M., Zhu, X., Wang, R., 2018. Efficient knn classification
958 with different numbers of nearest neighbors. IEEE Transactions on Neural
959 Networks and Learning Systems 29, 1774–1785. doi:10.1109/TNNLS.2017.
960 2673241.
- 961 Zhongguo, Y., Hongqi, L., Liping, Z., Qiang, L., Ali, S., 2017. A case
962 based method to predict optimal k value for k-nn algorithm. J. Intell.
963 Fuzzy Syst. 33, 55–65. URL: <https://doi.org/10.3233/JIFS-161062>,
964 doi:10.3233/JIFS-161062.

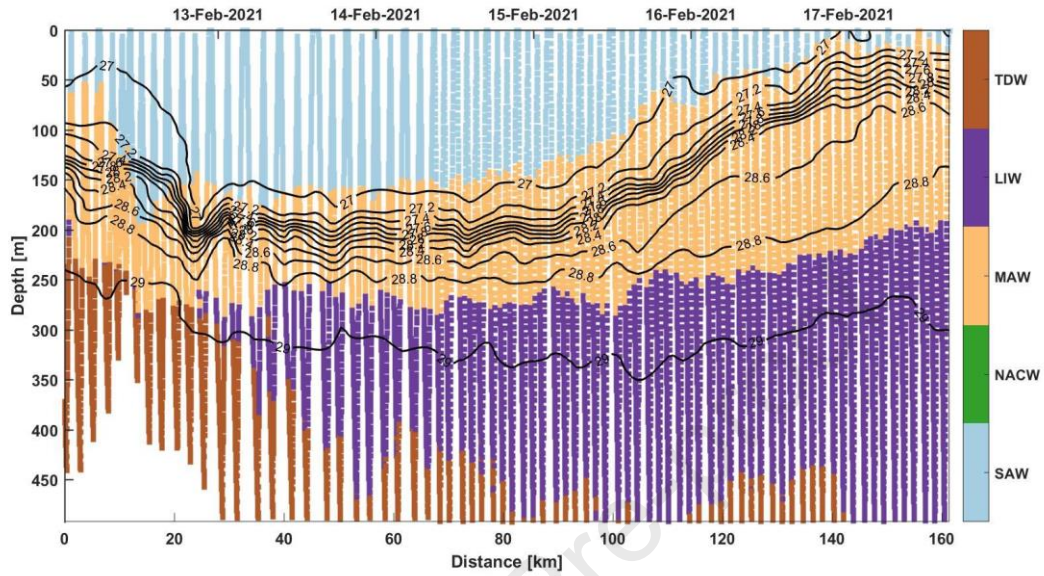


(a)

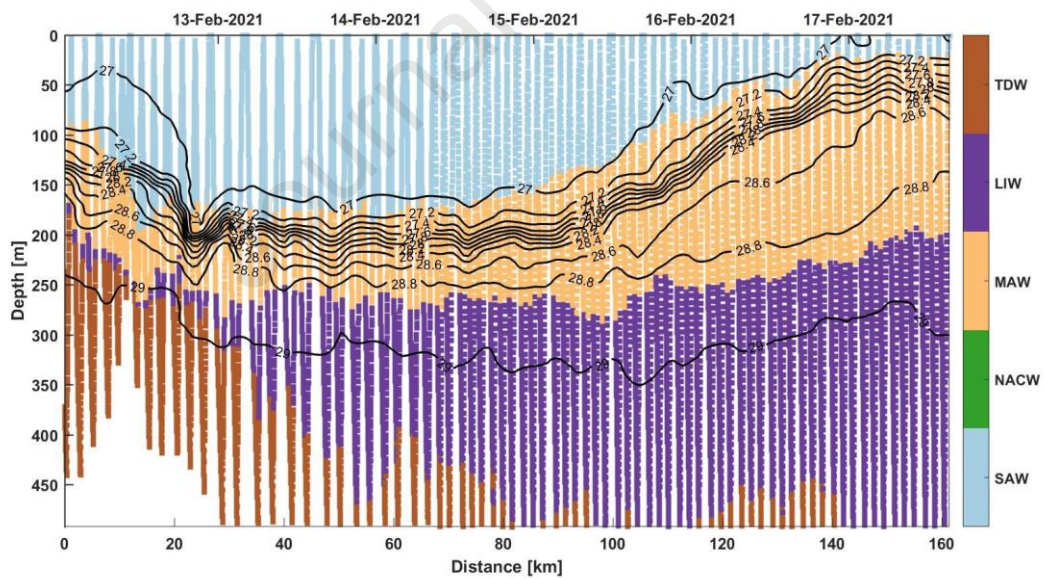


(b)

Figure 4: Example of a $\theta - S$ diagram (a) labelled by the Atlantic (SAW, NACW, MAW) and Mediterranean (LIW, TDW) water masses, with frontiers removed, and its equivalent $\sigma - \pi$ diagram (b).

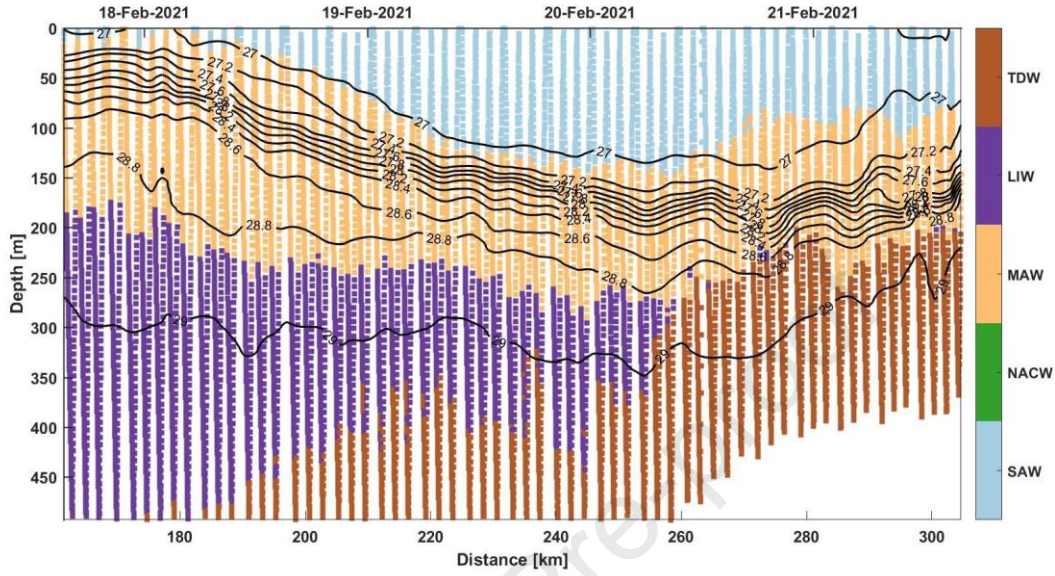


(a)

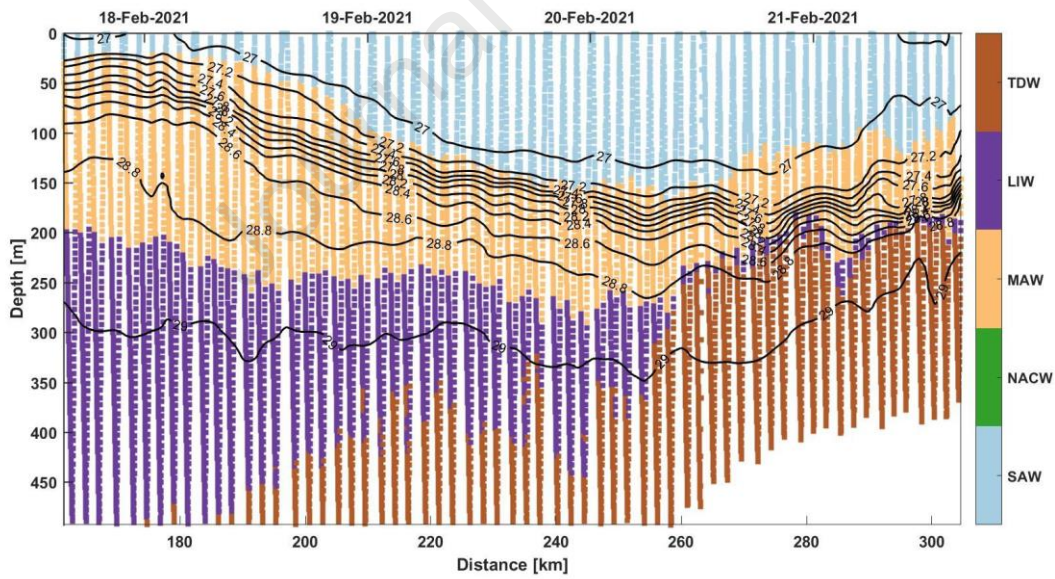


(b)

Figure 5: Classification of the water masses in the first transect using training dataset with frontiers removed for $\sigma - \pi$ (a) and $\theta - S$ (b) diagrams.

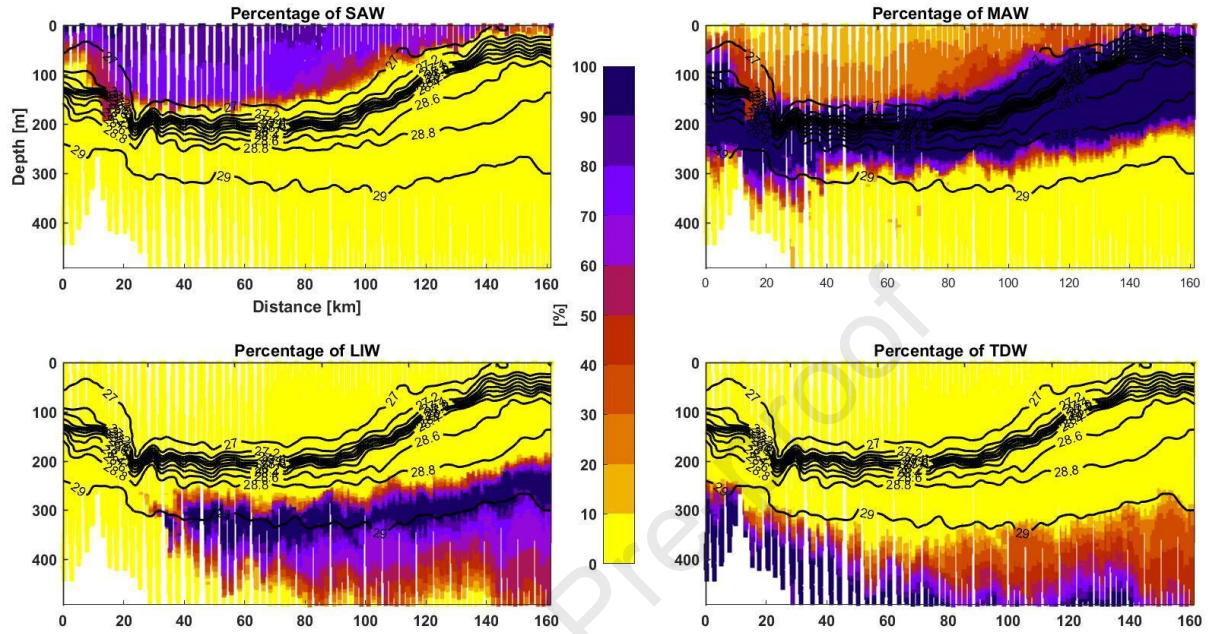


(a)

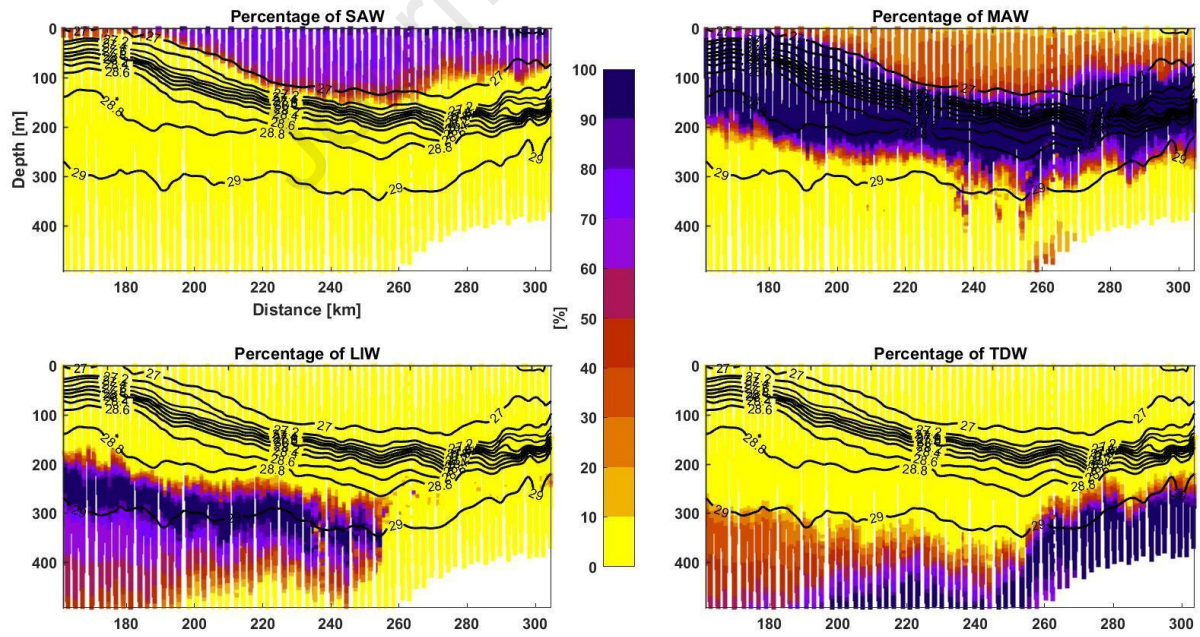


(b)

Figure 6: Classification of the water masses in the second transect using training dataset with frontiers removed for $\sigma - \pi$ (a) and $\theta - S$ (b) diagrams.

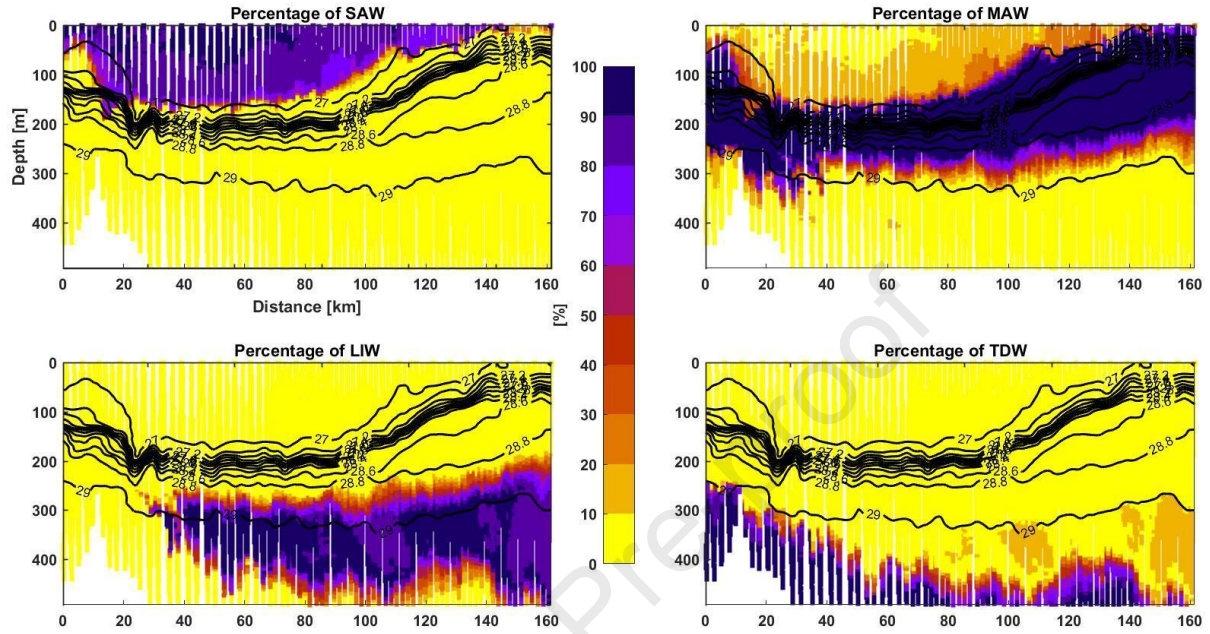


(a)

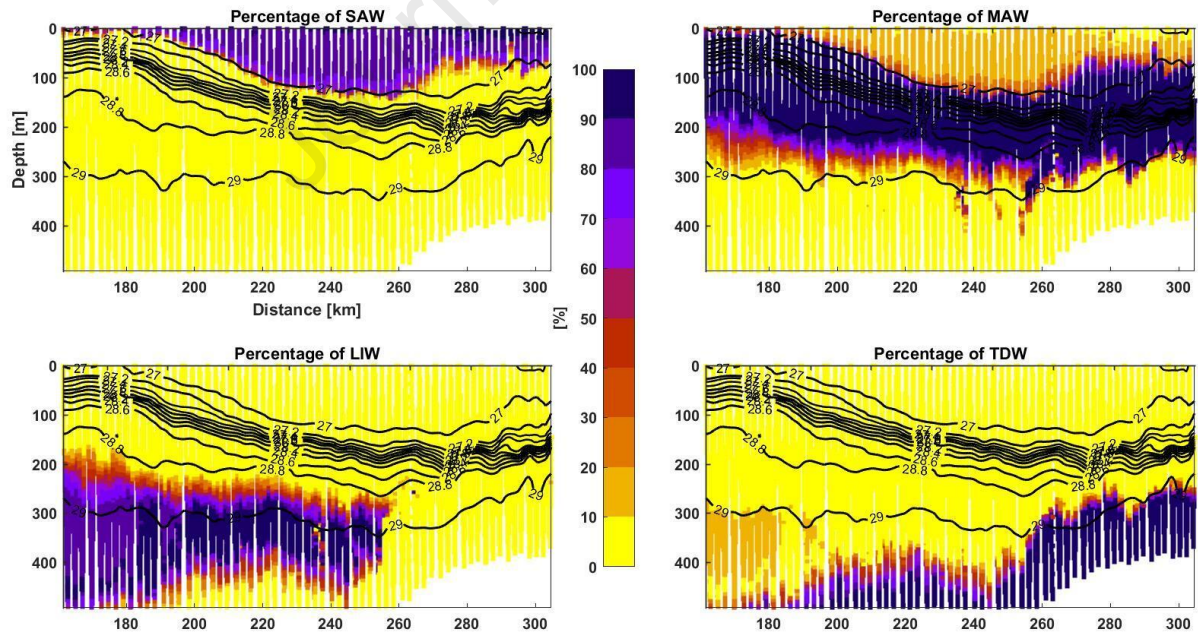


(b)

Figure 7: Percentage of the AWs and MWs along the first (a) and second (b) transects using $(\sigma-\pi)$ diagram. The sum of the four contributions leads to 100% in the Atlantic and Mediterranean layers.



(a)



(b)

Figure 8: Percentage of the AWs and MWs along the first (a) and second (b) transects using $(\theta - S)$ diagram. The sum of the four contributions leads to 100% in the Atlantic and Mediterranean layers.

Target water masses	LIW	91.9% 58156	4.6% 2880	3.5% 2240	0.0% 0	0.0% 0
	MAW	0.9% 2483	98.6% 281630	0.5% 1340	0.0% 117	0.0% 126
	TDW	1.1% 1858	0.7% 1175	98.2% 161549	0.0% 0	0.0% 0
	SAW	0.0% 0	0.2% 214	0.0% 0	99.7% 86673	0.0% 29
	NACW	0.0% 0	2.2% 95	0.0% 0	2.0% 86	95.9% 4204
		LIW	MAW	TDW	SAW	NACW
		Predicted water masses				

Figure 9: Confusion matrix for the 5 water masses deduced from the classification during the training stage.

Predicted water masses for the whole data	LIW	96.9% 6520	0.5% 32	2.6% 176	0.0% 0
	MAW	5.1% 531	94.7% 9807	0.0% 2	0.2% 16
	TDW	2.4% 114	4.1% 196	93.4% 4419	0.0% 0
	SAW	0.0% 0	0.0% 2	0.0% 0	100.0% 5868
		LIW	MAW	TDW	SAW
Predicted water masses for data from 1950 to 1980					

(a)

Predicted water masses for the whole data	LIW	98.9% 6653	0.1% 7	1.0% 68	0.0% 0
	MAW	3.7% 380	96.1% 9956	0.1% 6	0.1% 14
	TDW	0.5% 23	1.9% 92	97.6% 4614	0.0% 0
	SAW	0.0% 0	0.1% 8	0.0% 0	99.9% 5861
		LIW	MAW	TDW	SAW
Predicted water masses for data from 1950 to 1990					

(b)

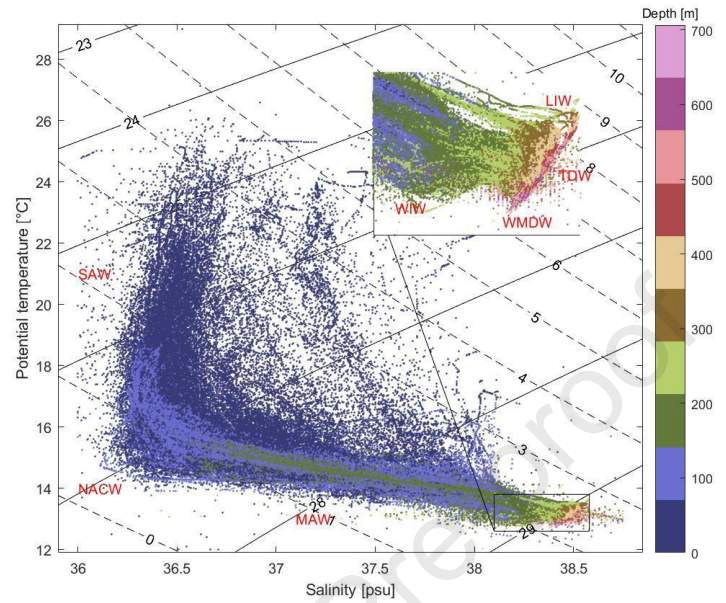
Predicted water masses for the whole data	LIW	99.6% 6703	0.1% 5	0.3% 20	0.0% 0
	MAW	2.2% 230	97.5% 10095	0.2% 22	0.1% 9
	TDW	1.7% 82	0.9% 42	97.4% 4605	0.0% 0
	SAW	0.0% 0	0.1% 6	0.0% 0	99.9% 5864
		LIW	MAW	TDW	SAW
Predicted water masses for data from 1950 to 2000					

(c)

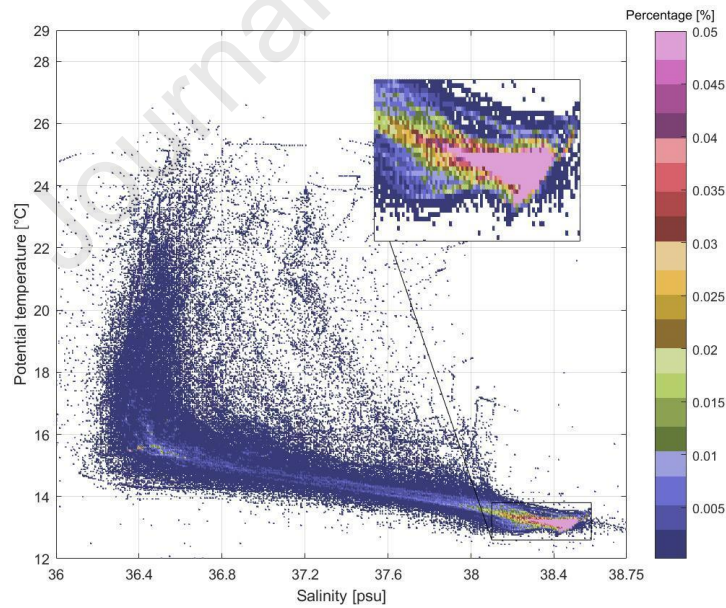
Predicted water masses for the whole data	LIW	99.6% 6703	0.1% 5	0.3% 20	0.0% 0
	MAW	1.2% 124	98.5% 10201	0.2% 22	0.1% 9
	TDW	1.7% 82	0.9% 43	97.4% 4604	0.0% 0
	SAW	0.0% 0	0.1% 6	0.0% 0	99.9% 5864
		LIW	MAW	TDW	SAW
Predicted water masses for data from 1950 to 2010					

(d)

Figure 10: Confusion matrices of training dataset profiles gathered from 1950 to 1980 (a), from 1950 to 1990 (b), from 1950 to 2000 (c) and from 1950 to 2010 (d).

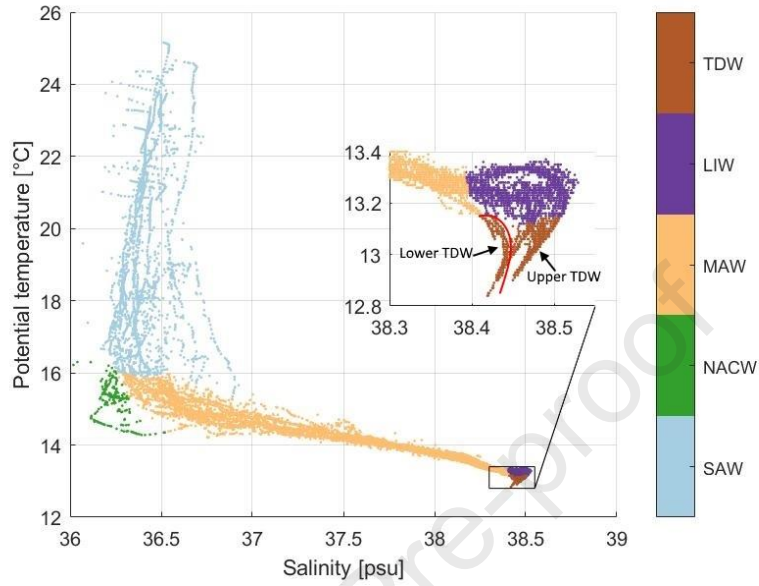


(a)

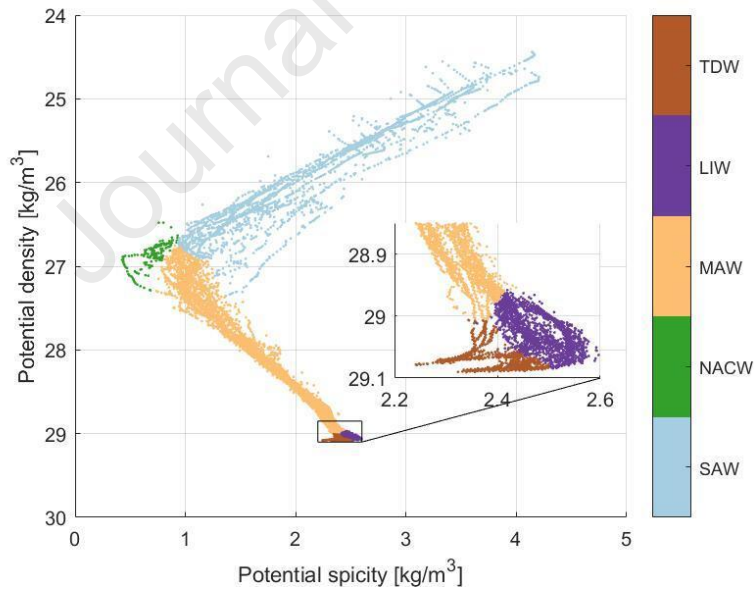


(b)

Figure 11: (a) $\theta - S$ diagram for all the database. Depths between the surface and 700 m are illustrated in different colors. isopycnals (solid lines) and spicity isopleths (dotted lines) are plotted $1 \text{ kg}\cdot\text{m}^{-3}$ apart. (b) Occurrence of water types as a function of temperature and salinity over temporal range of the WOD18 (1951-2020). Bin is scaled to represent percentage of total points. The color scales go from 0 to 0.05 for the sake of clarity.



(a)



(b)

Figure 12: Example of a θ -S diagram (a) labelled by the Atlantic (SAW, NACW, MAW) and Mediterranean (LIW, TDW) water masses and its equivalent σ - π diagram (b). Upper-TDW and lower-TDW are separated by the red curve plotted in the inset (a).

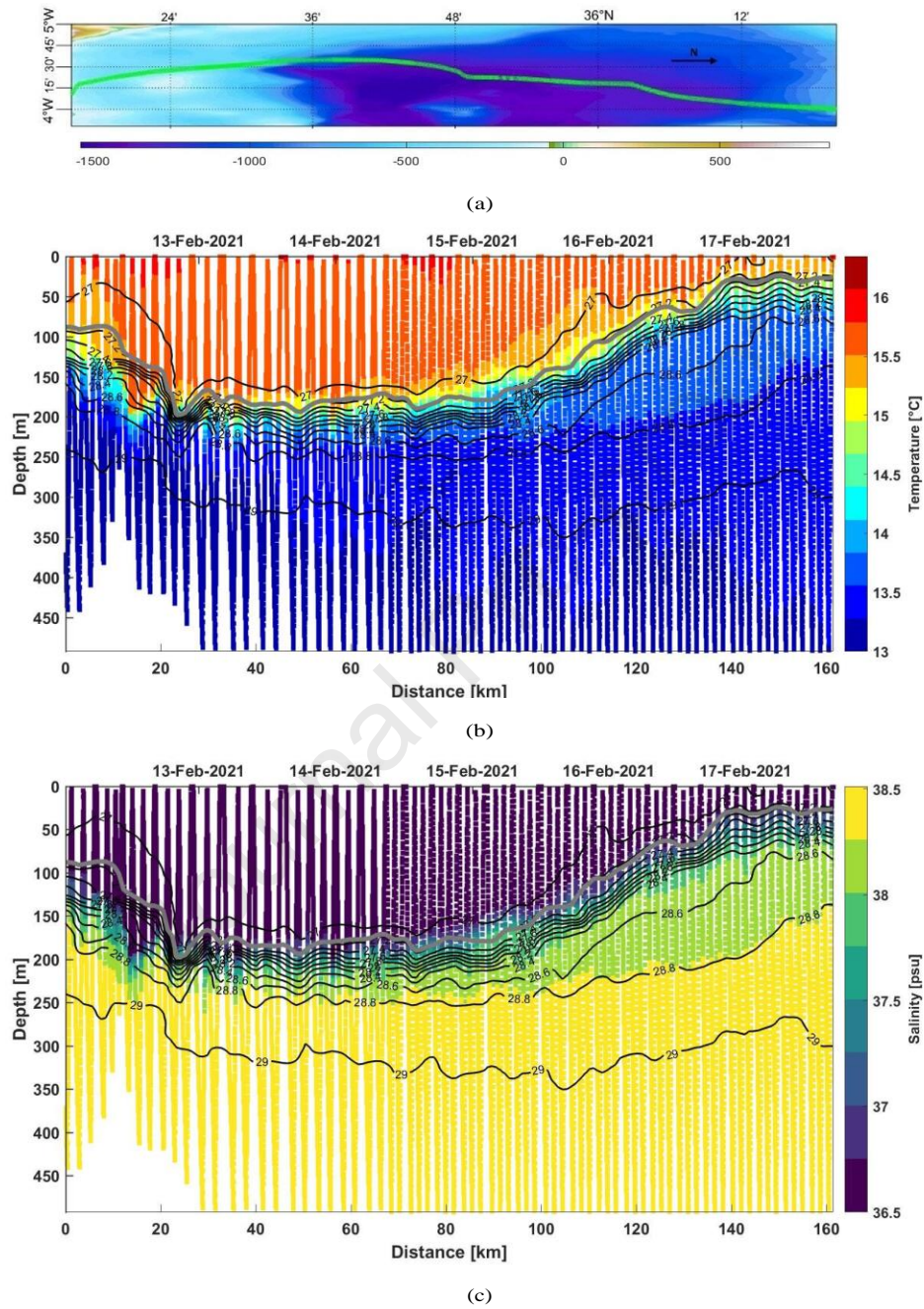


Figure 13: Temperature (b) and salinity (c) along the first glider transect (a). The black lines are the isopycnal levels and the gray line is the Mixed Layer Depth, defined using the threshold method with a finite difference criterion (density criterion of $0.03 \text{ kg}\cdot\text{m}^{-3}$). The black arrow in (a) points in the North direction.

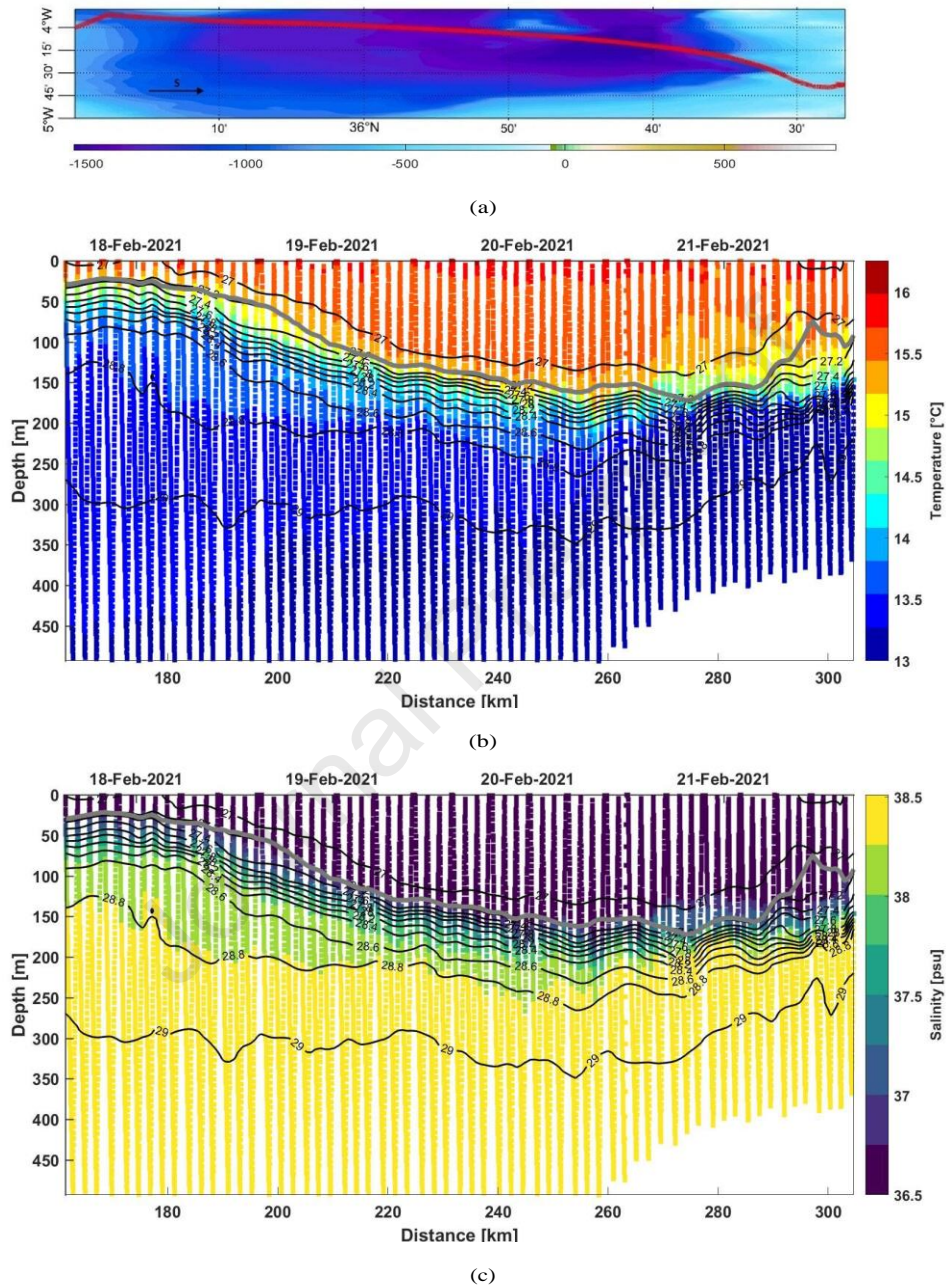


Figure 14: Temperature (b) and salinity (c) along the second glider transect (a). The black lines are the isopycnal levels and the gray line is the Mixed Layer Depth, defined using the threshold method with a finite difference criterion (density criterion of $0.03 \text{ kg}\cdot\text{m}^{-3}$). The black arrow in (a) points in the South direction.

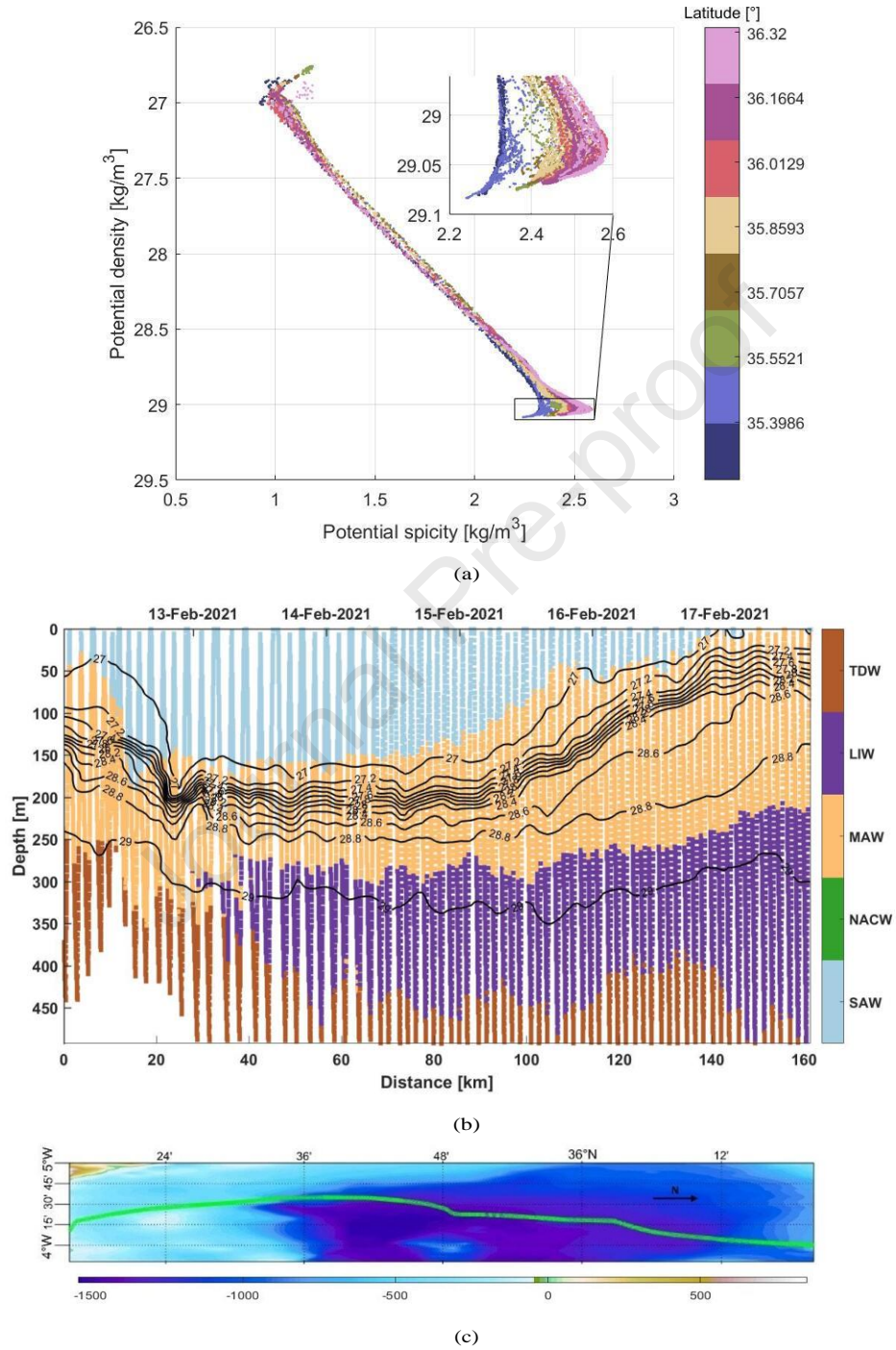


Figure 15: The $\sigma - \pi$ diagram (a) and the classification (b) of the water masses in the first glider transect (c). The black arrow in (c) points in the North direction.

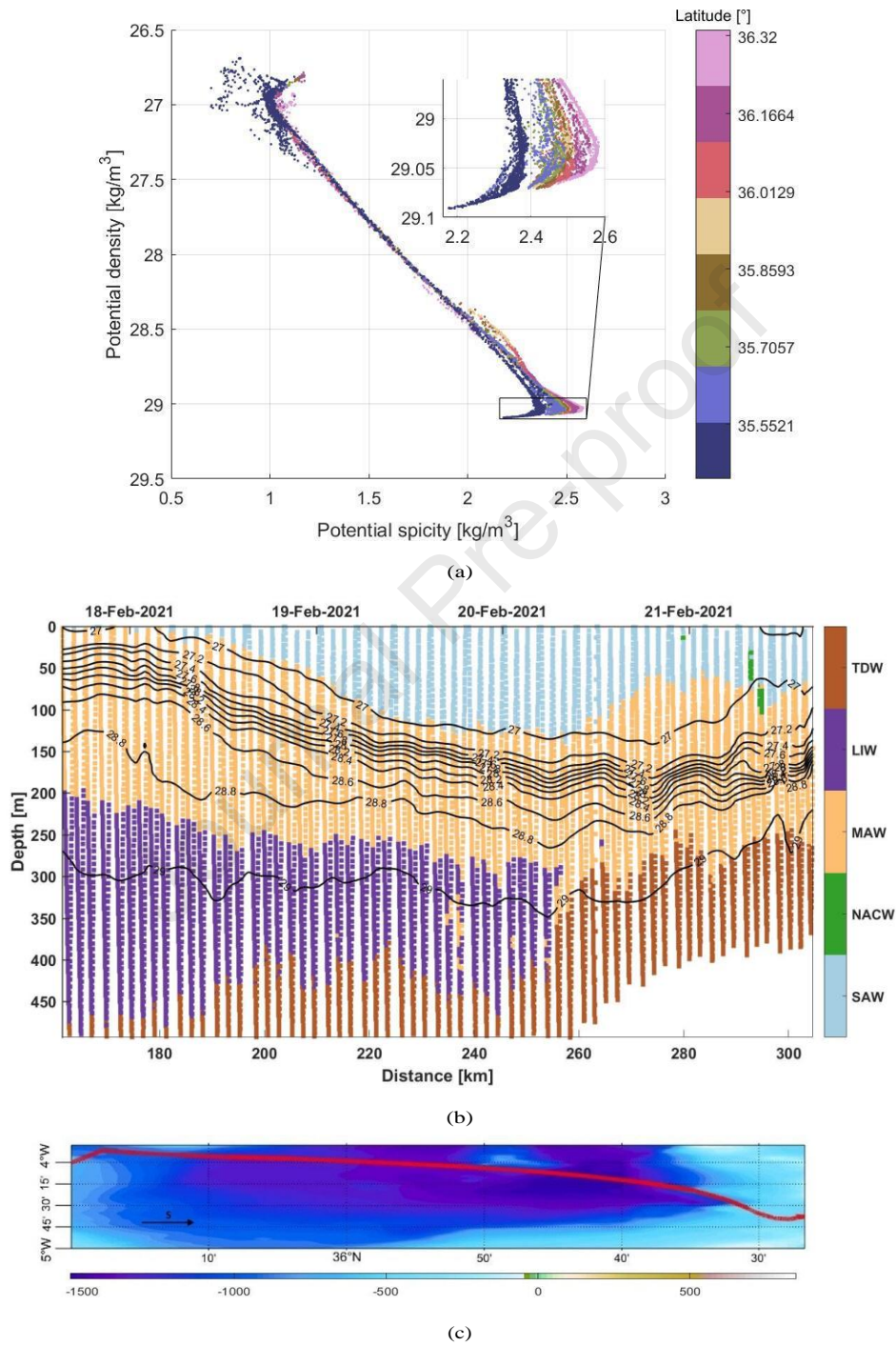


Figure 16: The $\sigma - \pi$ diagram (a) and the classification (b) of the water masses in the second glider transect (c). The black arrow in (c) points in the South direction.

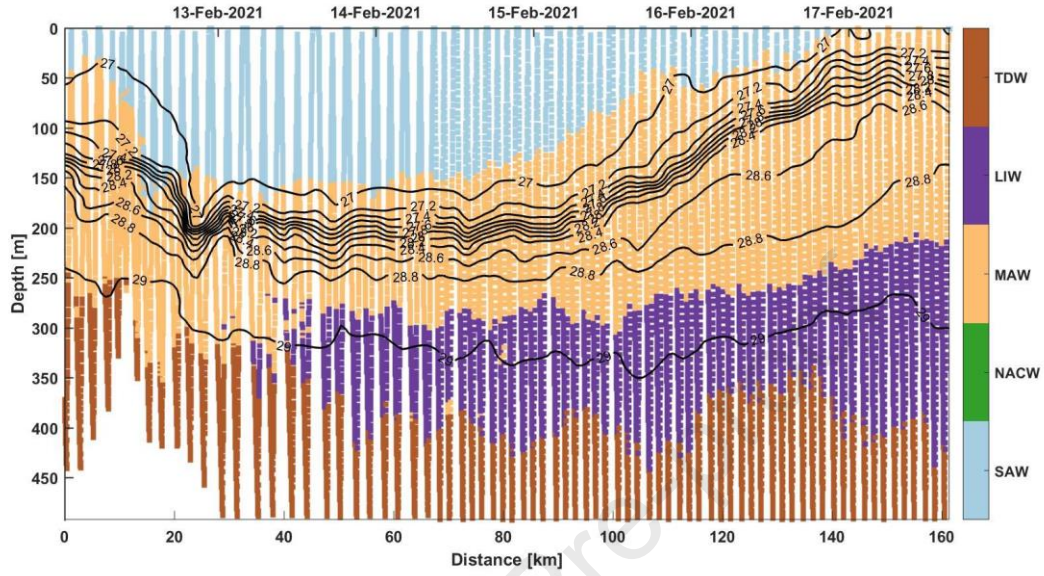
Predicted water masses for the whole data	LIW	71.3% 4797	2.5% 168	26.2% 1763	0.0% 0
	MAW	0.0% 5	97.2% 10065	2.8% 286	0.0% 0
	TDW	0.0% 0	0.4% 21	99.6% 4708	0.0% 0
	SAW	0.0% 0	3.6% 214	0.0% 0	96.4% 5650
		LIW	MAW	TDW	SAW
Predicted water masses for data with latitude < 35°45'					

(a)

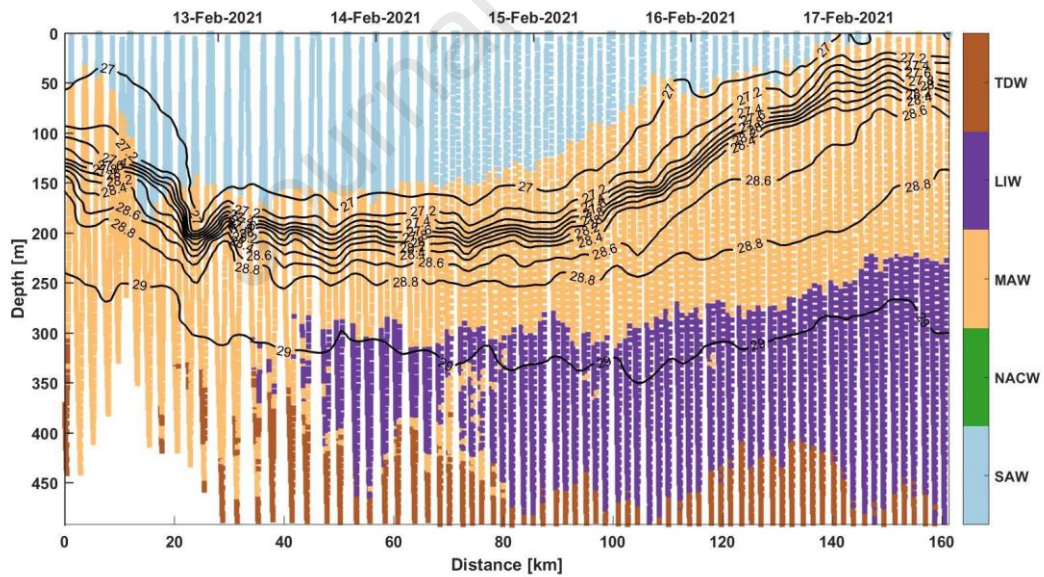
Predicted water masses for the whole data	LIW	81.0% 5449	16.4% 1105	2.6% 174	0.0% 0
	MAW	0.0% 0	98.9% 10241	0.7% 69	0.4% 45
	TDW	8.7% 412	20.2% 954	71.1% 3363	0.0% 0
	SAW	0.0% 0	0.0% 0	0.0% 0	100.0% 5861
		LIW	MAW	TDW	SAW
Predicted water masses for data with latitude > 35°45'					

(b)

Figure 17: Confusion matrices of training dataset profiles gathered below (a) and beyond (b) 35°45'N.

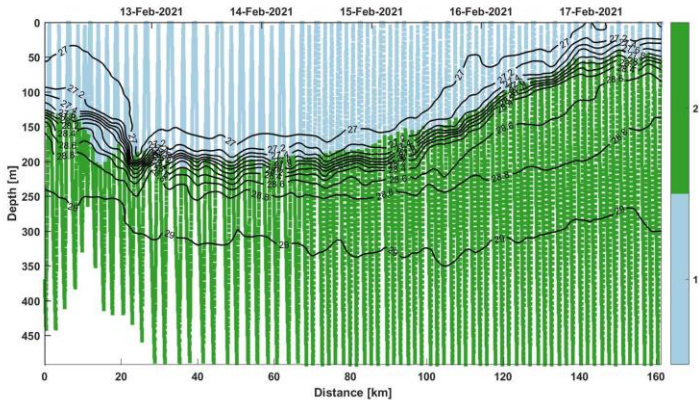


(a)

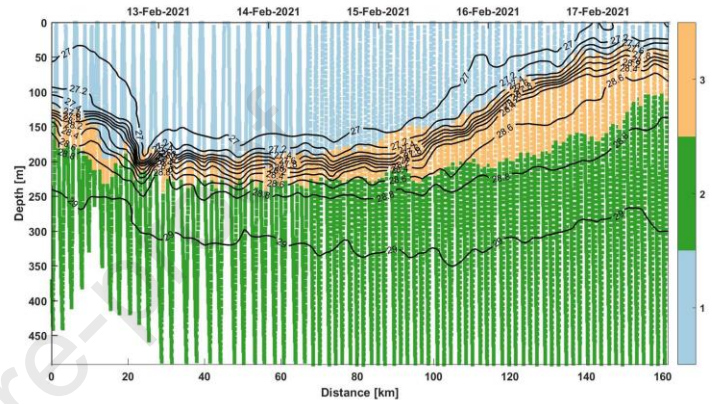


(b)

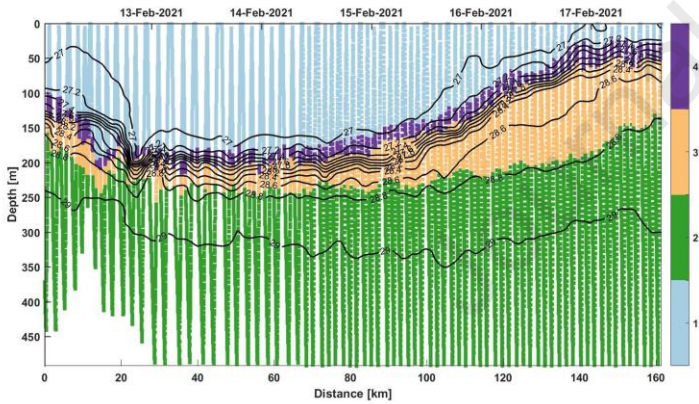
Figure 18: Classification of the water masses in the first transect using training dataset profiles gathered below (a) and beyond (b) $35^{\circ}45'N$.



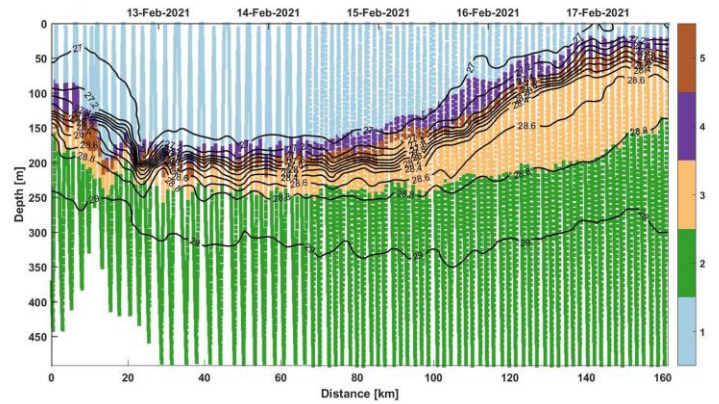
(a)



(b)



(c)



(d)

Figure 19: Classification of water masses in the first transect (Figure 15c) provided by the k-means clustering. The number of clusters used in: (a) $k = 2$, (b) $k = 3$, (c) $k = 4$ and (d) $k = 5$.

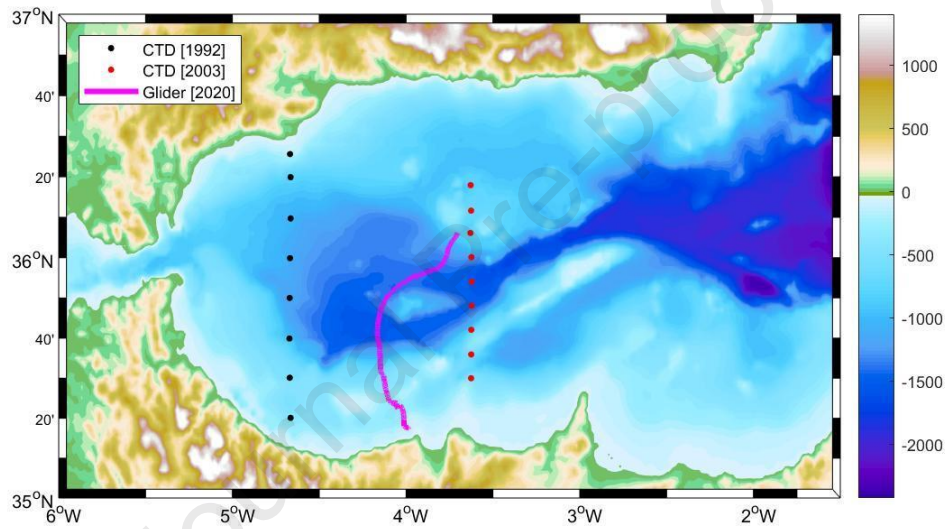
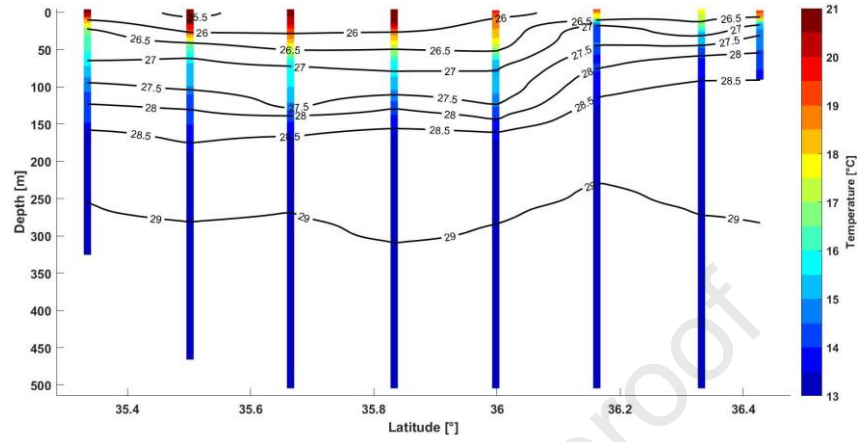
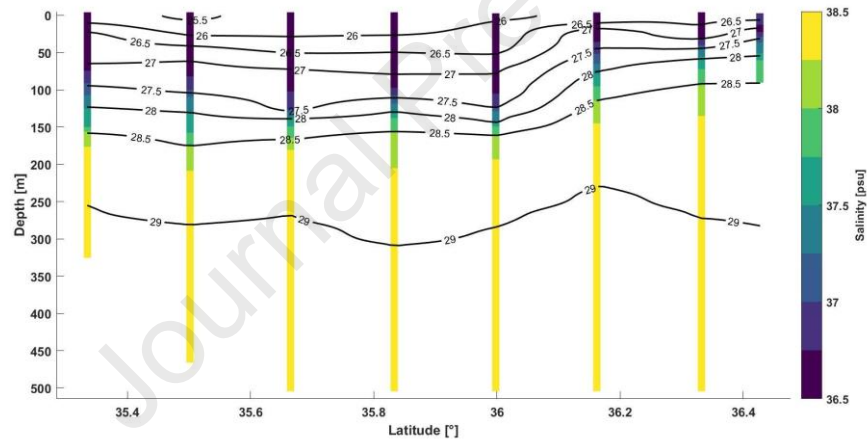


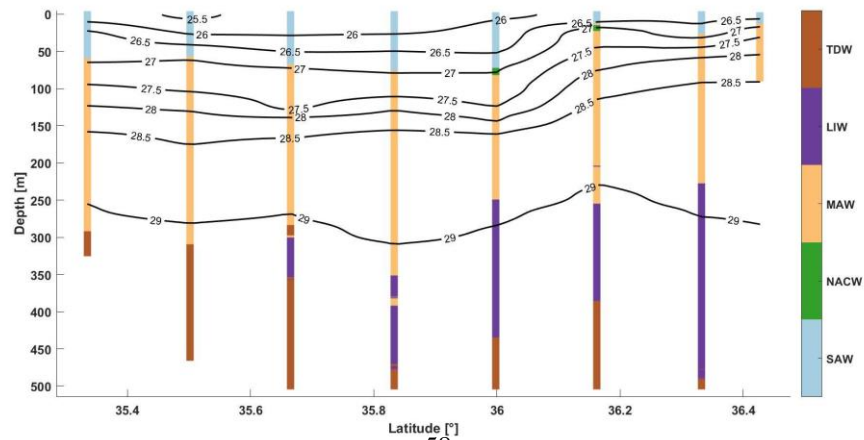
Figure A.20: Map of the Western Alboran Sea sketching the bathymetric depths and topographic elevations in meters (m) relative to the mean sea level. Black dots indicate the position of the CTD data gathered in 1992. Red dots represent the localization of the CTD data collected in 2003. The glider transect is sketched in magenta.



(a)

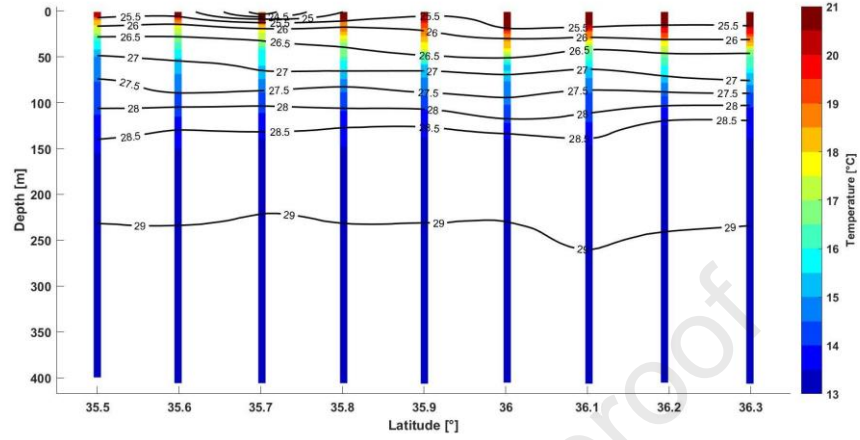


(b)

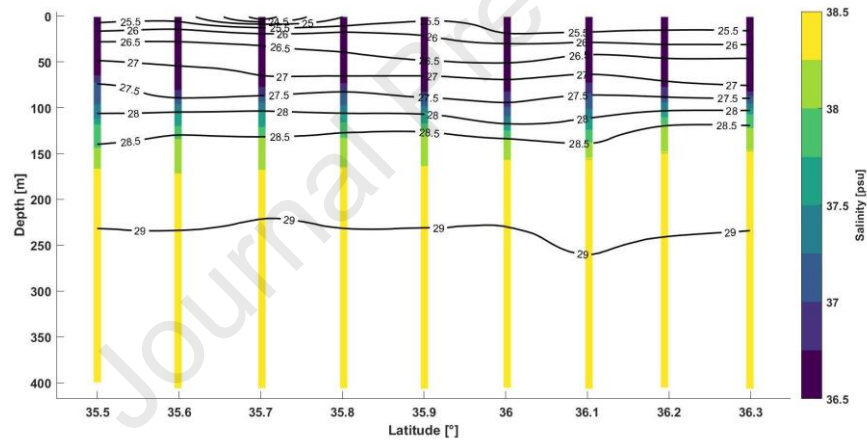


(c)

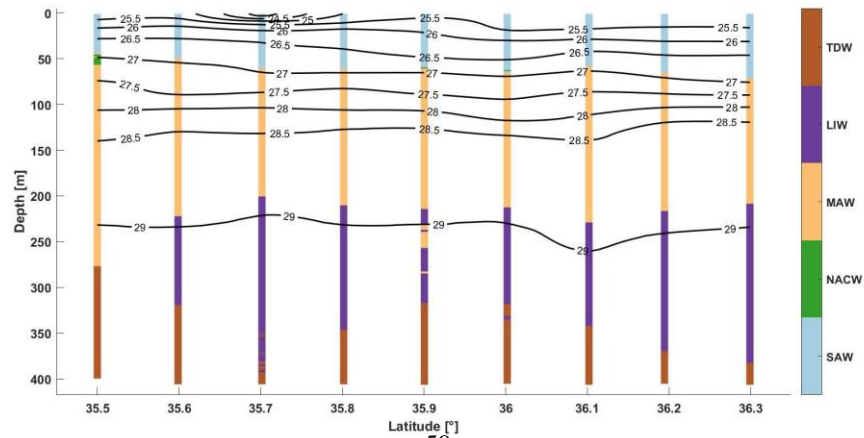
Figure A.21: Temperature (a) and salinity (b) along the westernmost transect (figure A.20). (c) represent the classification results. The black lines are the isopycnal levels



(a)

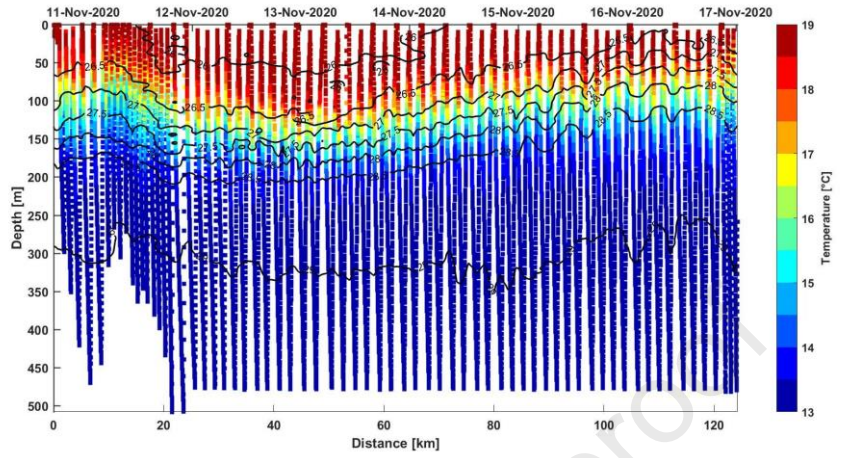


(b)

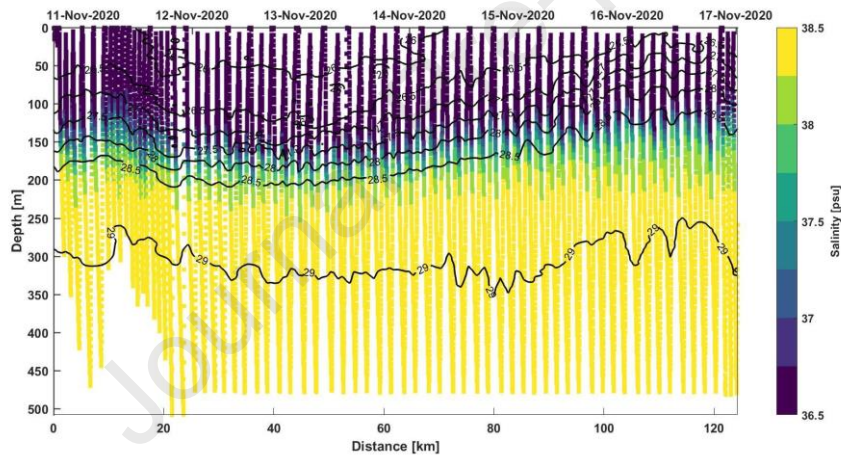


(c)

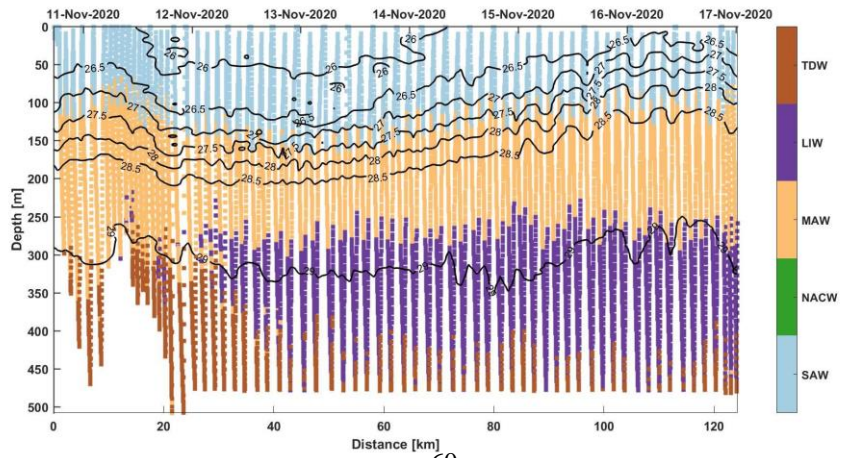
Figure A.22: Temperature (a) and salinity (b) along the easternmost transect (figure A.20). (c) represent the classification results. The black lines are the isopycnal levels.



(a)



(b)



(c)

Figure A.23: Temperature (a) and salinity (b) along the glider transect (figure A.20). (c) represent the classification results. The black lines are the isopycnal levels.

- High spatial resolution glider profiles of θ -S in the western Alboran sea ;
- Water masses derived on a $(\sigma$ - π) diagram using Knn algorithm ;
- Classification results confirm earlier derived circulation schemes ;
- The proposed method outperforms classical clustering analysis in delineating water mass boundaries.

Journal Pre-proof

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:

Journal Pre-proof