

# Time series segmentation for state-model generation of autonomous aquatic drones: a systematic framework

Alberto Castellini\*, Manuele Bicego, Francesco Masillo,  
Maddalena Zuccotto, Alessandro Farinelli

*Department of Computer Science, University of Verona, Verona, Italy,  
Email: name.surname@univr.it, \*corresponding author*

---

## Abstract

Autonomous surface vessels are becoming increasingly important for water monitoring. Their aim is to navigate rivers and lakes with limited intervention of human operators, to collect real-time data about water parameters. To reach this goal, these intelligent systems must interact with the environment and act according to the situations they face. In this work we propose a framework based on the integration of recent time-series clustering/segmentation methods and cluster validity indices, for detecting, modeling and evaluating aquatic drone states. The approach is completely data-driven and unsupervised. It takes unlabeled multivariate time series of sensor traces and returns both a set of statistically significant state-models (generated by different mathematical approaches) and a related segmentation of the dataset. We test the approach on a real dataset containing data of six campaigns, two in rivers and four in lakes, in different countries for about 5.6 hours of navigation. Results show that the methodology is able to recognize known states and to discover unknown states, enabling novelty detection. The approach is therefore an easy-to-use tool for discovering and interpreting significant states in sensor data, that enables improved data analysis and drone autonomy.

*Keywords:* Time series segmentation, situation assessment, state-model generation, autonomous surface vessels, activity recognition, water monitoring, model interpretation/explanation, sensor data analysis

---

## 1. Introduction

Autonomous robots have recently had a strong impact in the transition from manual (passive) to autonomous (active) water monitoring. These intelligent systems, used also in several other application domains, such as surveillance and monitoring (Farinelli et al. (2012)), are able to autonomously collect large amounts of data, providing crucial support to human operations. Aquatic drones involved in autonomous monitoring of catchments navigate rivers and lakes acquiring real-time data about water parameters, such as pH and dissolved oxygen. While human operators are usually involved in such data collection activities, direct tele-operation of the drones is often not an option for an entire mission, hence autonomous navigation is required. Navigation strategies usually aim at maximizing the information content of acquired data (Bottarelli et al. (2016, 2019)), while adapting to the conditions of the environment. Although data are very noisy in this context, applications require minimal number of sensors to reduce the costs.

A key factor for the success of autonomous data acquisition campaigns is *mission awareness* (Endsley (1995)), which is composed of three main elements: knowledge of mission objectives, internal self-situational awareness, and external self-situational awareness. In this work we specifically focus on the problem of detecting, modeling and interpreting aquatic drone states with data-driven methods, an aspect of self-situational awareness. By state we mean an abstract, compact and informative descriptor of key properties of the drone-environment system. In particular, we aim at developing *interpretable models of drone states* from traces of sensor data acquired during water-monitoring campaigns, by means of machine learning and artificial intelligence methods (Hastie et al. (2001); Bishop (2006); Russell and Norvig (2009)). Generating such a set of drone state-models is important for two reasons, namely, it supports *offline data analysis* by improving the extraction of knowledge from large sensor traces, and it enhances the autonomy of the drone by providing key information for *online decision making* (Kaelbling and Lozano-Perez (2013); Asperti et al. (2019)).

Automatic detection of aquatic drone states from sensor data can be performed by supervised or unsupervised methods. Supervised methods are typically more accurate than unsupervised methods but they need labeled datasets, usually hard, expensive and sometimes impossible to collect in real monitoring campaigns. Ad-hoc experiments could be performed to generate labelings, but they usually consider only subsets of situations that the

38 drone faces during real campaigns. On the other hand, many data is usually  
39 available from past campaigns that can be mined by unsupervised methods.

40 This work focuses on unsupervised approaches, namely *clustering* and  
41 *time series segmentation*, able to split multivariate time series into groups  
42 of observations corresponding to system states and having common proper-  
43 ties that can be compactly represented by mathematical *models*. The goal is  
44 to discover these states (and models) using data-driven methods from sensor  
45 data of past campaigns. The literature (see Section 2) proposes several meth-  
46 ods for this purpose, characterized by different assumptions and extracting  
47 different types of patterns. The main difference between the works in the  
48 literature and our work is that we propose a *systematic framework* for gener-  
49 ating and evaluating statistically significant state-models for aquatic drones,  
50 while the literature mainly proposes novel clustering methods or it compares  
51 standard methods in different application domains.

52 We first investigated clustering and subspace clustering methods for de-  
53 tecting aquatic drone states in (Castellini et al. (2018b, 2019c)). Here, we ex-  
54 tend those works using both classic (Bishop (2006)) and very recent methods,  
55 including SubCMedians (Peignier et al. (2018)), Toeplitz Inverse Covariance-  
56 based Clustering (TICC) (Hallac et al. (2017)) and Inertial Hidden Markov  
57 Models (IHMM) (Montanez et al. (2015)). The proposed framework is tested  
58 on a large datasets with observations from many campaigns. State-models  
59 are analyzed and interpreted in terms of situations faced by the drones. The  
60 statistical significance of state-models is computed by comparing their prop-  
61 erties with those of random clusters. Since different aspects of state-model  
62 performance must be evaluated, we select a set of validity indices (Arbelaitz  
63 et al. (2013)) satisfying the requirements of our domain.

64 The main contributions of this paper are summarized in the following:

- 65 • we propose an easy-to-use framework for systematically generating and  
66 evaluating significant state-models in multivariate time series;
- 67 • we successfully apply the proposed framework to a real dataset of sensor  
68 data collected by aquatic drones involved in water monitoring;
- 69 • we present, analyze and interpret, with high level of detail, both the dis-  
70 covered state-models and the application procedures used to generate  
71 these models, which makes this manuscript a valuable reference also  
72 for practitioners interested in analyzing similar data and performing  
73 extensive cross-comparison of methodologies;

74 • we present and make available the dataset used in this analysis<sup>1</sup>.

75 The rest of the manuscript is organized as follows. Section 2 provides an  
76 overview of the state-of-the-art on this research topic. Section 3 introduces  
77 the aquatic drone architecture and the proposed framework for state-model  
78 generation. In Section 4 we describe the dataset and the labelings. Section 5  
79 introduces clustering and segmentation methods, and the procedures for the  
80 generation of random clusterings and segmentations. Section 6 defines some  
81 clustering validity indices and performance measures. Section 7 illustrates  
82 the results and some state-models generated by the proposed framework.  
83 Conclusions and future directions are drawn in Section 8.

## 84 2. Related work

85 From the *application* point of view, strong similarities are present with  
86 sensor-based human activity recognition (Chen et al. (2012); Dhiman and  
87 Vishwakarma (2019)), where sensors are used to acquire data about human  
88 movements and machine learning methods are employed to generate activity  
89 models and to predict human activities in novel contexts. The main difference  
90 between our problem and human activity recognition is that data collected  
91 by aquatic drones are very noisy, since they come from several sources (not  
92 only accelerometers as in applications of human activity recognition) and  
93 are strongly influenced by unstructured and diversified environments (e.g.,  
94 rivers and lakes in different parts of the world have disparate environmental  
95 properties). Moreover, aquatic drones collect two kinds of data, some relating  
96 to movement, others to water properties, and both sources of information can  
97 be used to assess the drone state.

98 From a *methodological* viewpoint, the main theoretical connections with  
99 our work concern clustering (Bishop (2006)) and time series segmentation (Fu  
100 (2011); Castellini et al. (2015)). K-means, Gaussian mixture models (GMM)  
101 and hierarchical clustering, have been recently used to identify activities of  
102 both humans (Abdallah et al. (2012); Trabelsi et al. (2013); Kwon et al.  
103 (2014); Barták and Vomlelová (2017)) and flying drones (Barták and Vomlelová  
104 (2017)) from sensor data. Hidden Markov models (HMMs) have been  
105 applied (Kim et al. (2010); Trabelsi et al. (2013); Barták and Vomlelová  
106 (2017)) and also extended (Fox et al. (2008); Montanez et al. (2015)) in

---

<sup>1</sup>The dataset is available here <https://data.mendeley.com/datasets/gtt7stf5x8/1>

107 the same context. Time series segmentation (Hallac et al. (2016a, 2017);  
108 Chiu et al. (2003)), change point detection (Barnett and Onnela (2016)) and  
109 motif discovery methods, have been employed to identify homogeneous in-  
110 tervals in sequential time-dependent data. The last techniques have been  
111 very recently applied also to problems related to driver identification (Hal-  
112 lac et al. (2016b)) and state representation of modern automobiles (Hallac  
113 et al. (2018)).

114 In previous works we tested standard clustering methods on single cam-  
115 paigns (Castellini et al. (2018a,b)) and introduced the usage of subspace clus-  
116 tering for generating sparse state-models (Castellini et al. (2019c,a)). What  
117 differentiates this paper from our previous work and the approaches in the lit-  
118 erature mentioned above is that here we propose a systematic framework for  
119 generating statistically significant state-models using very recent techniques  
120 and, most important, for evaluating them by several internal and external  
121 validity indices. Moreover, we test the proposed framework on a large real  
122 dataset in the application domain of autonomous water monitoring and we  
123 analyze the statistical properties of detected states. Furthermore, we select  
124 some validity indices (Arbelaitz et al. (2013); Moshtaghi et al. (2019)) and  
125 used them to evaluate and rank the state-models generated by five clustering  
126 techniques.

### 127 **3. System overview**

128 In this section we describe the two main elements of our system, namely  
129 the aquatic drone architecture and the framework for state-model generation.

#### 130 *3.1. Data acquisition system: autonomous aquatic drones*

131 Data acquisition campaigns are performed by Lutra mono hull boats (see  
132 Figure 1) produced by Platypus<sup>2</sup> and customized in the EU Horizon 2020  
133 INTCATCH project<sup>3</sup> to accomplish water monitoring of catchments. Lo-  
134 calization and orientation are provided by an on-board smartphone which  
135 gathers information from GPS, compass and gyroscope. Sensor manage-  
136 ment and sensor data transmission to the cloud is performed by a Go-Sys

---

<sup>2</sup><http://senseplatypus.com>

<sup>3</sup><http://www.intcatch.eu>

137 BlueBox<sup>4</sup> control unit connected to an arduino e-board. Operators can de-  
 138 fine desired paths by setting waypoints in a map on a tablet, to perform  
 139 autonomous navigation, or they can manually drive the drone using an RC  
 140 controller. Drones are equipped with sensors for GPS position, water temper-  
 141 ature, dissolved oxygen and electrical conductivity, commands to propellers  
 142 and battery voltage. Sensor traces are stored in log files on the smartphone  
 143 or transmitted to the cloud by a Go-Sys BlueBox. Log files are preprocessed  
 144 using Platypus Python libraries to obtain a matrix of time series having one  
 145 sensor signal in each row and time instants in columns. Since different sen-  
 146 sors have different sampling frequencies the alignment of sensor traces was  
 147 obtained via interpolation and re-sampling, with sampling frequency of 1Hz.

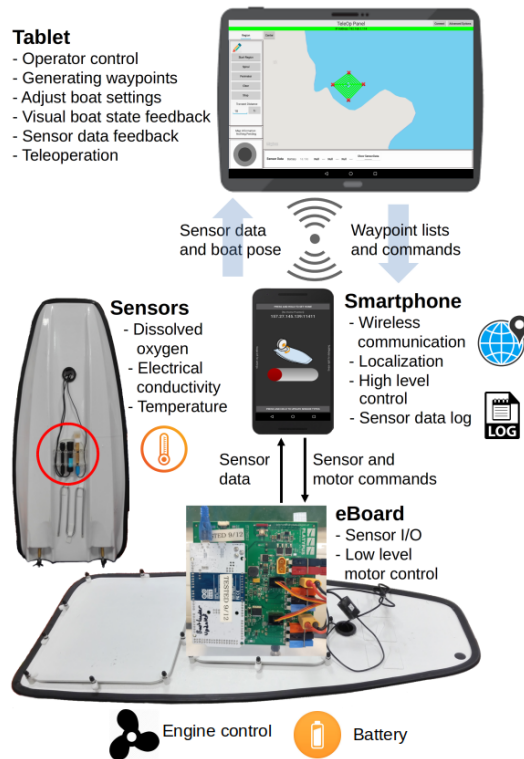


Figure 1: Overview of the drone architecture.

<sup>4</sup><https://www.go-sys.de/en/bluebox/>

148 3.2. Framework for state-model generation and evaluation

149 The framework proposed in this work is outlined in Figure 2. The input  
 150 *dataset* is a matrix of multivariate time series with engineered features (see  
 151 Section 4), which contains sensor readings from multiple campaigns. Data  
 152 are processed by five *clustering and segmentation methods*, namely, k-means  
 153 (KM), Toeplitz Inverse Covariance-based Clustering (TICC), Hidden Markov  
 154 Models (HMM), Inertial Hidden Markov Models (IHMM), and SubCMedians  
 155 (SCM). They generate clusterings depending on parameter settings. Multiple  
 156 instances of random clustering (RC) and random segmentation (RS) are also  
 157 generated. They are used as baselines to evaluate the significance of the  
 158 state-models generated by real clustering algorithms (see Section 5).

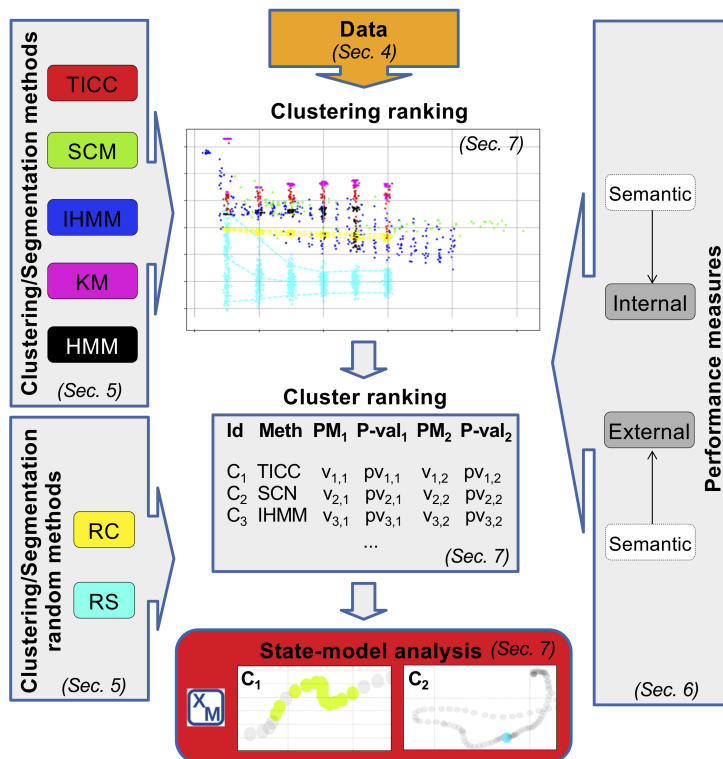


Figure 2: Overview of the proposed framework for state-model generation and evaluation.

159 Clusterings and related clusters are then evaluated by means of *performance*  
 160 *measures* (see Section 6). They have different semantics and can  
 161 favour different kinds of patterns (i.e., states) in the data (e.g., the *silhouette*  
 162 is maximized if clusters are both compact and distant from each other, while

163 *spread* considers only the cluster compactness). Performance measures enable to rank clusterings and clusters, and to identify the best state-models. 164 After computing performance, we also determine cluster (clustering) p-values 165 using random partitioning as baselines. Only clusters (clusterings) with low 166 p-values are considered statistically significant. The last step of the proposed 167 framework involves the *analysis and interpretation of significant state-models* 168 (performed in Section 7). Since each state-model is generated by a cluster- 169 ing method, evaluated by some performance measures, and interpreted as a 170 situation, the framework enables different kinds of analyses involving combi- 171 nations of these properties. For instance, we analyze the statistical properties 172 of significant state-models, compare the capability of different methods to dis- 173 cover specific situations, and compare the capability of different performance 174 measures to rank situations. State-model analysis is supported by a Python 175 tool called eXplainable Modeling<sup>5</sup> (Castellini et al. (2019d)) that integrates 176 several data visualization and statistical tools. 177

#### 178 4. Dataset

179 We analyze sensor traces generated in six independent campaigns (also 180 called experiments in the following). Table 1 shows the name, number of sam- 181 ples, duration and type of catchment (i.e., river or lake) of each campaign. 182 Since our goal is to generate a unique set of state-models, we concatenated 183 the traces of all the campaigns, obtaining a single dataset (called *CON-* 184 *CAT*) with 20187 observations and about 5.6 hours of navigation, since the 185 sampling frequency is 1Hz. Variables available in the raw dataset are time, 186 latitude, longitude, altitude, speed, electrical conductivity, dissolved oxygen, 187 temperature, battery voltage, heading, acceleration, command to propeller 188 0 and command to propeller 1 (the boat has two propellers). Using only 189 these variables we obtain experiment-dependent state-models because of the 190 strong differences in environmental parameters among different campaigns. 191 To avoid this problem we generate new variables by feature extraction. In 192 particular, we compute *moving means* and *standard deviations* over a slid- 193 ing windows of 10 seconds, and *variations* between couples of consecutive 194 observations. The list of 27 variables in the final dataset is reported in Ta- 195 ble 2. Z-score standardization was performed on each variable to improve 196 the performance of clustering and segmentation methods.

---

<sup>5</sup><https://github.com/XModeling/XM>



197 **Mathematical notation.** In the following, we use notation  $X = \{x_1, x_2,$   
198  $\dots x_n\}$  to represent the dataset, where  $n$  is the number of observations (i.e.,  
199  $n = 20187$  in our dataset), each observation  $x_i \in X$  has  $D$  variables (i.e.,  
200  $D = 27$  in our dataset). Each variable is represented by a number ranging  
201 from 1 to  $D$ , and the set of all variables is denoted  $\mathcal{D} = \{1, \dots, D\}$ .

<b>Id</b>	<b>Campaign name</b>	<b>Samples</b>	<b>Duration</b>	<b>Lake/River</b>
1	ESP2	2814	47'	R
2	ESP5	3601	60'	R
3	ESP4	2374	39'	L
4	GARDA3	2451	40'	L
5	ITA1	7243	121'	L
6	ITA6	1704	28'	L
-	CONCAT	20187	335'	-

Table 1: List of data acquisition campaigns in the dataset.

<b>Symbol</b>	<b>Description</b>
$s, v, a$	Instantaneous speed, voltage, acceleration
$m_0, m_1$	Instantaneous signal to propeller 0 and 1
$\bar{s}, \bar{v}, \bar{a}$	Moving average mean of speed, voltage, acceleration
$\bar{m}_0, \bar{m}_1$	Moving average mean of signal to propeller 0 and 1
$\hat{s}, \hat{v}, \hat{a}$	Moving average std of speed, voltage, acceleration
$\hat{e}c, \hat{d}o, \hat{T}$	Moving average std of electrical conductivity, dissolved oxygen, temperature
$\hat{m}_0, \hat{m}_1$	Moving average std of signal to propeller 0 and 1
$\hat{h}$	Moving average std of heading
$\tilde{s}, \tilde{a}, \tilde{v}$	Variation of speed, voltage, acceleration
$\tilde{m}_0, \tilde{m}_1$	Variation of signal to propeller 0 and 1
$\tilde{e}c, \tilde{d}o, \tilde{h}$	Variation of electrical conductivity, dissolved oxygen, temperature

Table 2: List of variables extracted from the dataset and used for clustering/segmentation.

#### 202 4.1. Known drone states

203 Some drone states are easy to identify by observing the drone paths in  
204 geographical maps but hard to detect from sensor traces, hence recognizing

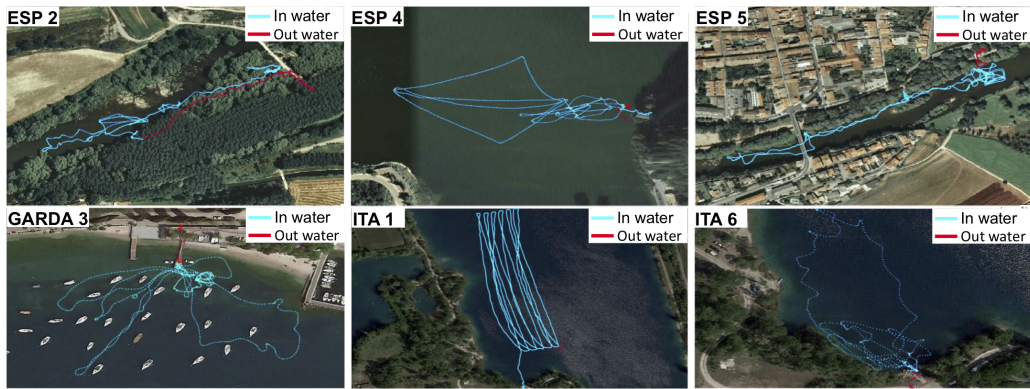


Figure 3: Geo-localization of monitoring campaigns and manual labelling of situations “drone into the water” (blue) and “drone out of the water” (red) (best viewed in color).

205 them is not a trivial task for clustering methods. We use these states to test  
 206 the ability of different methods to detect real situations. The states that we  
 207 manually label are: drone into the water (IW), drone out of the water (OW),  
 208 upstream navigation (US), downstream navigation (DS), no water stream  
 209 (NS), manual drive (MD), autonomous drive (AD), and turning (T). Figure  
 210 3 shows the labelled paths of states IW (cyan) and OW (red).

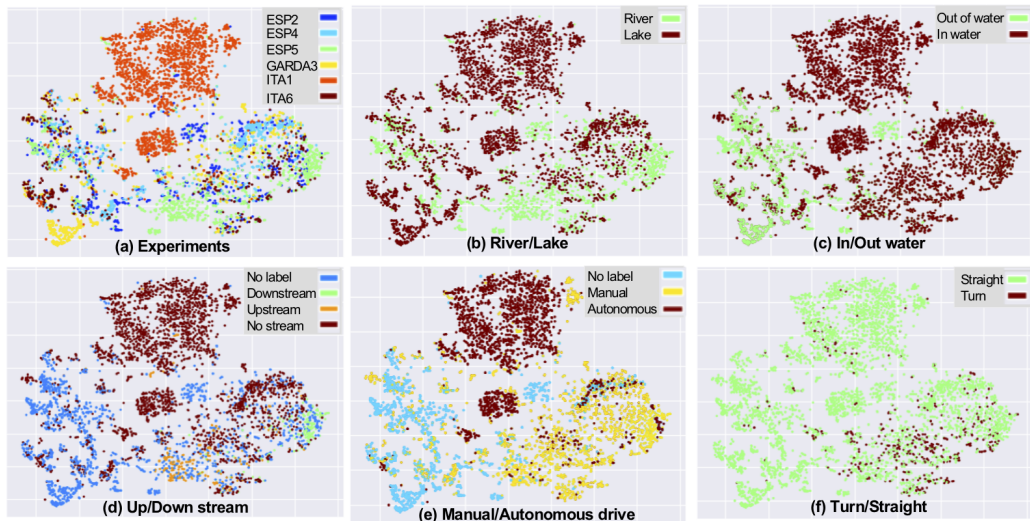


Figure 4: t-SNE projections. Points represent data observations and colors correspond to known situations (best viewed in color).

211 *4.2. Dimensionality reduction analysis*

212 We use t-Distributed Stochastic Neighbor Embedding (t-SNE) (van der  
213 Maaten and Hinton (2008)) to see if known situations correspond to implicit  
214 structures in the data. t-SNE allows the implicit structure in the data to  
215 influence the way in which subset of data points are gathered, hence it reveals  
216 structures at different scales. In Figure 4.a, for instance, colors represent  
217 experiments (e.g., ESP2) and in Figure 4.c they represent situations in/out  
218 water. Projections are informative, they show grouping of observations and  
219 correspondence between groups and situations (colors). For instance, the  
220 coloring related to in/out water (Figure 4.c) identifies well separated clusters,  
221 as expected, although more than one dense region is present for each label.

222 **5. Clustering and time series segmentation methods**

223 We generate our state-models by five clustering or time series segmenta-  
224 tion methods, namely, k-means, SubCMedians, TICC, HMMs and IHMMs.  
225 The main difference between clustering and time series segmentation is that  
226 clustering does not consider time proximity between observations, while time  
227 series segmentation considers it, generating groups of *adjacent* observations  
228 (called segments) having common properties. Here we briefly introduce the  
229 methodologies and their peculiarities. The sets of parameters used in the  
230 training phase, for each method, are also described (see Table 3). Since all  
231 methods are unsupervised, the real number of clusters is unknown, hence we  
232 test several combinations of methods and parameters and leave the selection  
233 of the best state-models to subsequent statistical analysis. Finally, we de-  
234 scribe the procedures for generating random clusterings and segmentations.

235 *5.1. K-means (KM)*

236 K-means<sup>6</sup> is an iterative descent clustering method (Bishop (2006)) which  
237 aims at minimizing the objective function  $J = \sum_{i=1}^n \sum_{c=1}^k r_{ic} \|x_i - \mu_c\|^2$ ,  
238 where  $r_{ic} \in \{0, 1\}$  is a binary indicator of point-cluster membership,  $x_i$  is a  
239 data point,  $\mu_c$  is the centroid of cluster  $c$ ,  $n$  is the number of data points and  $k$   
240 the number of clusters. Each clustering is a set of centroids that minimizes  $J$ .  
241 We use Euclidean distance  $\|\cdot\|^2$ , number of clusters  $k$  listed in Table 3, and  
242 for each clustering, we re-initialized the algorithm 100 times and selected the

---

<sup>6</sup><https://scikit-learn.org/>

Method	Parameter	Values
KM	$k$	{5, 10, 15, 20, 25, 30}
	# repeats	50
SCM	NbExtClust	{2, 3, 4, 5, 6, 10, 15, 20, 25, 30}
	# repeats	10
TICC	$k$	{5, 10, 15, 20, 25, 30}
	$\lambda$	{0.1, 0.5, 0.7, 1.0}
	$\beta$	{0, 50, 100, 150, 200}
	$w$	{ 1, 3 }
	# repeats	1
HMM	$k$	{5, 10, 15, 20, 25}
	# repeats	50
IHMM	$k$	{2, 4, 6, ..., 38, 40}
	$\zeta$	{0, 5, 10, ..., 65, 70}
	# repeats	1
RC	$k$	{5, 10, 15, 20, 25, 30}
	# repeats	200
RS	$k$	{5, 10, 15, 20, 25, 30}
	# repeats	200

Table 3: Learning parameters of all clustering methods tested.

243 best clustering, since initial conditions influence the solution. We compute  
244 50 clusterings (# repeats in Table 3) for each  $k$ .

### 245 5.2. SubCMedians (SCM)

246 SubCMedians is a recent center-based subspace clustering technique (Peignier  
247 et al. (2018)). This algorithm is based on a K-medians paradigm and it aims  
248 at clustering data points around suitable candidate centers  $m_i \in \mathcal{M}$ , where  
249 centers are defined in different subspaces (i.e., subsets of variables)  $\mathcal{D}_i \subseteq \mathcal{D}$ .  
250 In our work, each subspace cluster represents a putative state of the aquatic  
251 drone. Formally, the goal of SCM is to build a set of centers  $\mathcal{M}$ , so as to  
252 minimize the Sum of Absolute Errors between the dataset and the centers  
253  $SAE(X, \mathcal{M}) = \sum_{x \in X} AE(x, \mathcal{M})$ , and such that  $Size(\mathcal{M}) \leq SD_{max}$ , where  
254  $Size(\mathcal{M}) = \sum_i |\mathcal{D}_i|$ , and  $SD_{max}$  is a parameter denoting the maximum Sum  
255 of Dimensions used in  $\mathcal{M}$  to describe all its centers. The Absolute Error  
256  $AE(x, \mathcal{M})$  represents the distance between each point  $x \in X$  and its closest  
257 center  $m_i \in \mathcal{M}$ , and it is computed as  $AE(x, \mathcal{M}) = \min_{m_i \in \mathcal{M}} dist(x, m_i)$ ,

258 where  $dist(x, m_i) = \sum_{d \in \mathcal{D}_i} |x_d - m_{i,d}| + \sum_{d \in \mathcal{D} \setminus \mathcal{D}_i} |x_d - \mu_d|$  is an extension of  
 259 the Manhattan distance, with  $m_{i,d}$  the coordinate of  $m_i$  along variable  $d$ , and  
 260  $\mu_d$  the mean of the coordinates of all points in  $X$  along  $d$ .

261 The algorithm<sup>7</sup> has three main parameters, namely  $SD_{max}$  (described  
 262 above), the sample size  $N$  (the algorithm considers only  $N$  randomly chosen  
 263 observations at each iteration) and the number of iterations  $NbIter$  of the  
 264 training process. The number of centers is not fixed in advance. In (Peignier  
 265 et al. (2018)), guidelines are provided to compute all parameters from a single  
 266 meta-parameter called  $NbExpClust$  and representing the expected number  
 267 of clusters. The actual number of clusters is then computed during training.  
 268 Table 3 shows the values of  $NbExpClust$  that we test and the number of  
 269 repetitions of each test. The algorithm needs less than one minute to compute  
 270 a clustering on an Intel(R) Core(TM) i7-6700HQ CPU @ 2.60GHz with 8GB  
 271 of RAM.

### 272 5.3. Toeplitz Inverse Covariance-Based Clustering (TICC)

273 TICC clusters are modeled as sparse Gaussian inverse covariance (Toeplitz)  
 274 matrices representing dependencies between variables. In particular, off-  
 275 diagonal elements represent partial correlations and on-diagonal elements the  
 276 inverse of variable variances (i.e., variable compactness) inside the cluster.  
 277 Formally, TICC computes a set of  $k$  Toeplitz matrices  $\Theta = \{\Theta_1, \dots, \Theta_k\}$  and  
 278 a clustering (i.e., assignment of observations to clusters)  $P = \{P_1, \dots, P_k\}$   
 279 that solve the following optimization problem (Hallac et al. (2017)):

$$\operatorname{argmin}_{\Theta \in \mathcal{T}, P} \sum_{j=1}^k \left[ \overbrace{\|\lambda \circ \Theta_j\|_1}^{\text{sparsity}} + \sum_{Y_i \in P_j} \left( \overbrace{-\ell\ell(Y_i, \Theta_j)}^{\text{log likelihood}} + \overbrace{\beta \mathbb{1}\{Y_{i-1} \notin P_j\}}^{\text{temporal consistency}} \right) \right]$$

280 where  $\mathcal{T}$  is the set of symmetric block Toeplitz matrices,  $\|\lambda \circ \Theta_j\|_1$  is an  
 281  $\ell_1$ -norm penalty of the Hadamard product aiming to sparsify the inverse  
 282 covariance matrices,  $\lambda$  is a matrix of regularization parameters that we set  
 283 to a single value  $\lambda \in \mathbb{R}$  to simplify parameter setting,  $Y_i$  is a concatenation  
 284 of observations  $x_{i-w+1}, \dots, x_i$ ,  $w \in \mathbb{R}$ ,  $\ell\ell(Y_i, \Theta_j)$  is the log-likelihood that  
 285 observation  $Y_i$  belongs to cluster  $\Theta_j$ ,  $\beta$  is a regularization parameter for  
 286 temporal consistency, and  $\mathbb{1}\{Y_{i-1} \notin P_j\}$  is an indicator function checking if  
 287 neighbouring observations are assigned to same cluster.

---

<sup>7</sup><https://sergiopeignier.github.io/>

288 The algorithm<sup>8</sup> uses four parameters, namely,  $\lambda$  that controls Toeplitz  
289 matrix sparsity,  $\beta$  that controls temporal consistency in clusters, the windows  
290 size  $w$  used to generate matrix  $Y$  from the dataset  $X$ , and the number of  
291 clusters  $k$ . The parameter values and the number of repetitions we test are  
292 displayed in Table 3. We set the maximum number of iterations to 100. For  
293 time reasons, tests using  $w = 3$  are performed only with  $\lambda = 1.0$  and  $\beta = 0.0$ .  
294 On an Intel(R) Core(TM) i7-6700HQ CPU @ 2.60GHz with 8GB of RAM  
295 the algorithm takes from 1 to 30 minutes to compute a clustering with  $w = 1$   
296 (longer time is taken with smaller  $\lambda$ s and  $\beta$ s) and between 40 minutes and  
297 1.5 hours with  $w = 3$ .

#### 298 5.4. Hidden Markov Models (HMM)

299 Hidden Markov models (Rabiner (1989); Bishop (2006)) are probabilistic  
300 models which describe Markovian stochastic processes. Observation models  
301 are set to single component multivariate Gaussian distributions (with one  
302 dimension for each observed variable). The initial state distribution is set  
303 to uniform over the set of hidden states, the initial transition matrix is set  
304 to a random stochastic matrix, initial means are computed by k-means and  
305 initial covariance matrices are set according to the obtained k-means clusters.  
306 The maximum number of iterations for the EM algorithm<sup>9</sup> is set to 100.  
307 The Viterbi algorithm (Bishop (2006)) is used to generate the most likely  
308 sequence of hidden states (i.e., drone states) given the observed sequence of  
309 sensor readings. We generated models having number of hidden states (i.e.,  
310 clusters) listed in Table 3. The learning algorithm was not able to generate  
311 clusterings with 30 or more clusters which are instead available for all other  
312 methods.

#### 313 5.5. Inertial Hidden Markov Models (IHMM)

314 IHMMs (Montanez et al. (2015)) are a regularization-based extension of  
315 HMMs in which the transition matrix is biased towards the inertial property,  
316 namely, it has increased self-transition (i.e., on-diagonal) values to better  
317 adapt to naturally “long lasting” activities observed in several contexts, such  
318 as human activity recognition. The basic idea is to introduce prior knowledge,  
319 in the form of a supplementary learning parameter  $\zeta$ , related to the expected

---

<sup>8</sup><https://github.com/davidhallac/TICC>

<sup>9</sup><https://scikit-learn.org/>

320 duration of activities, so that the HMM tends to reduce state transitions and,  
321 consequently, to generate long segments along the time axis instead of frag-  
322 menting adjacent observations in several states. The observation model of  
323 each state is represented by the parameters of a multivariate Gaussian distri-  
324 bution. IHMMs are trained by standard EM algorithm, where the transition  
325 matrix update is modified to consider parameter  $\zeta$ . In our tests we set pa-  
326 rameters  $k$  and  $\zeta$  as shown in Table 3. The algorithm<sup>10</sup> needs between 30  
327 seconds and 100 minutes (longer time is needed when more hidden states are  
328 used) to compute a single clustering on an Intel(R) Core(TM) i7-6700 CPU  
329 @ 3.40GHz with 16GB of RAM.

### 330 5.6. Random clustering (RC)

331 Random clusterings are generated by assigning to each observation in the  
332 dataset a uniformly random number from 1 to  $k$  (the number of clusters). The  
333 obtained vector of labels (i.e., numbers from 1 to  $k$ ) is used as a clustering,  
334 hence observations assigned to the same label are put together in the same  
335 group. We generate 200 random clusterings for each  $k \in \{5, 10, 15, 20, 25, 30\}$   
336 (see Table 3) and use them to compute the statistical significance of cluster-  
337 ings and clusters generated by standard methods.

### 338 5.7. Random segmentation (RS)

339 Random segmentations are generated by selecting  $k - 1$  different random  
340 splitting points between 2 and  $n - 1$ , and then assigning label 1 to the  
341 observations before the first splitting point, label 2 to observations between  
342 the first and the second splitting point, and so on, until the last interval of  
343 observations (between the last splitting point and the last observation) which  
344 was assigned to label  $k$ . In this way we generate  $k$  segments of random length,  
345 in which each segment is related to a single cluster. As for RC we generate  
346 200 random segmentations for each  $k \in \{5, 10, 15, 20, 25, 30\}$  (see Table 3).

## 347 6. Performance measures

348 A key element for evaluating state-models generated by different clus-  
349 tering methods are performance measures. Since different aspects of the  
350 performance must be evaluated, here we propose an ensemble of indices that

---

<sup>10</sup><https://github.com/george-montanez/InertialRegularizedHMM>

351 satisfy the requirements of our and possibly other application domains. Se-  
 352 lected indices can be split into three categories, namely, measures for eval-  
 353 uating *clusterings*, measures for evaluating single *clusters* (i.e., state-models  
 354 in our context), and measures for evaluating state-model *variables*. The first  
 355 and second categories can be further divided into *external* and *internal*. The  
 356 former uses a ground truth to evaluate the clustering/cluster, while the lat-  
 357 ter does not require any labeling. Since the goal of the proposed framework  
 358 is to provide quality state-models from unlabeled data, we focus our analy-  
 359 sis on internal performance measures, however, some external measures are  
 360 presented to assess the capability of clustering methods to detect known situ-  
 361 ations. For each internal and external measure we specify if it can be applied  
 362 at clustering level, at cluster level or both. The measures are then used in  
 363 Section 7 to evaluate, rank, select and interpret state-models generated by  
 364 different methods. Symbol  $\uparrow$  ( $\downarrow$ ) is used to identify measures that must be  
 365 maximized (minimized). In all indices below the notation  $d_e(x_i, x_j)$  is used  
 366 to represent the Euclidean distance between observations  $x_i$  and  $x_j$ . We no-  
 367 tice that the performance indices here used focus on cluster and clustering  
 368 goodness, not on their prediction capabilities. We do not split our dataset in  
 369 training and test set, compute models on training set and evaluate them on  
 370 test set (a way to evaluate prediction capabilities of state-models). The prob-  
 371 lem we tackle here comes before the prediction problem, in fact we generate  
 372 state-models that could be eventually processed to learn prediction models.  
 373 An advantage of this approach is a lower time complexity (computing predic-  
 374 tion performance on test sets needs time consuming cross-validation) which  
 375 allows us to select optimal state-model among a large set of clusters generated  
 376 by several combinations of clustering methods and parameter settings.

### 377 6.1. Internal measures

378 **Silhouette** ( $\mathcal{S}, \uparrow$ ). The *silhouette* (Rousseeuw (1987); Arbelaitz et al.  
 379 (2013)) is an internal measure that contrasts the average distance to elements  
 380 in the same cluster with the average distance to elements in other clusters.  
 381 Cluster cohesion is measured based on the distance between all the points  
 382 in the same cluster, the separation between clusters is based on the nearest  
 383 neighbour distance. The silhouette of a single observation  $x_i$  assigned to a  
 384 cluster  $z_c$  is defined as:

$$\mathcal{S}(x_i^c) = \frac{b(x_i, z_c) - a(x_i, z_c)}{\max\{a(x_i, z_c), b(x_i, z_c)\}}$$



385 where  $a(x_i, z_c)$  is the average distance of  $x_i$  from the other observations in  
386 cluster  $z_c$  and  $b(x_i, z_c)$  is the minimum average distance between  $x_i$  and the  
387 observations in clusters  $z_l \neq z_c$ . Silhouette can be computed for a specific  
388 cluster  $z_c$ , as  $\mathcal{S}(z_c) = 1/|z_c| \sum_{x_i \in z_c} \mathcal{S}(x_i)$ , or for an entire clustering  $Z$ , as  
389  $\mathcal{S}(Z) = 1/n \sum_{z_c \in Z} \sum_{x_i \in z_c} \mathcal{S}(x_i)$ . Its values range from -1 to 1 where high  
390 values indicate points belonging to perfectly compact and separated clusters  
391 and low values indicate clustering with mixed clusters.

392 **Davies-Bouldin index** ( $\mathcal{DB}, \downarrow$ ). Davies-Bouldin index (Davies and  
393 Bouldin (1979); Arbelaitz et al. (2013)) estimates the cohesion as the distance  
394 from the observations in a cluster to its centroid (computationally faster than  
395 computing distances between all pairs of observations in the cluster, as in sil-  
396 houette) and the separation based on the distance between centroids (also  
397 faster than silhouette). The cohesion is divided by the separation, hence the  
398 index must be minimized. The index formula is

$$\mathcal{DB}(Z) = 1/k \sum_{z_c \in Z} \max_{z_l \neq z_c} \left\{ \frac{C(z_c) + C(z_l)}{d_e(\bar{z}_c, \bar{z}_l)} \right\},$$

399 where  $\bar{z}_c$  is the centroid of cluster  $z_c$  and  $C(z_c)$  is the estimated cohesion of  
400 cluster  $z_c$ ,  $C(z_c) = 1/|z_c| \cdot \sum_{x_i \in z_c} d_e(x_i, \bar{z}_c)$ .

**Calinski-Harabasz index** ( $\mathcal{CH}, \uparrow$ ). Calinski-Harabasz index (Caliński  
and Harabasz (1974); Arbelaitz et al. (2013)) estimates cluster cohesion from  
the distances between cluster points and related cluster centroids. The sep-  
aration is estimated from the distance between the centroids and the global  
centroid of the dataset  $\bar{X}$ . The separation term is finally divided by the cohe-  
sion term, hence this index is ratio-based and must be maximized. Formally,

$$\mathcal{CH}(Z) = \frac{n - k}{k - 1} \frac{\sum_{z_c \in Z} |z_c| d_e(\bar{z}_c, \bar{X})}{\sum_{z_c \in Z} \sum_{x_i \in z_c} d_e(x_i, \bar{z}_c)}$$

401 where  $\bar{z}_c$  is the number of observations in cluster  $z_c$ ,  $\bar{z}_c$  is the centroid of  $z_c$ .

**Spread** ( $\mathcal{Q}, \downarrow$ ). The spread of a cluster is a measure of cluster cohesion  
(Kelley et al. (1996)). Given a cluster  $z_c$  containing  $|z_c|$  observations the  
spread is given by

$$\mathcal{Q}(z_c) = \frac{(\sum_{x_i \in z_c} \sum_{x_j \in z_c, j > i} d_e(x_i, x_j))}{|z_c|(|z_c| - 1)/2}.$$

402 The measure can be extended to clusterings by averaging cluster spreads as

403 
$$\mathcal{Q}(Z) = \frac{\sum_{c=1}^k \mathcal{Q}(z_c)}{k}.$$

**Weighted spread** ( $\mathcal{R}, \downarrow$ ). Since clusters with small number of observations are more likely to be more compact, and consequently to have smaller spread than large clusters, we computed a weighted version of the cluster spread, in which the spread is divided by the percentage of observations in the cluster, namely,

$$\mathcal{R}(z_c) = (\mathcal{Q}(z_c)/|z_c|) \cdot n.$$

404 The extension to clusterings is obtained as a sum of weighted cluster spread,  
405 that is  $\mathcal{R}(Z) = \sum_{z_c \in Z} \mathcal{R}(z_c)$ .

**NMRCLUST penalty** ( $\mathcal{P}, \downarrow$ ). In (Kelley et al. (1996)) an internal measure is proposed to compare clusterings having different number of clusters and possibly being generated by different methods. The index is computed for a clustering  $Z$  as  $\mathcal{P}(Z) = \mathcal{N}\mathcal{Q}(Z) + k$ , where the first term is the sum of the normalized average spread of the clustering

$$\mathcal{N}\mathcal{Q}(Z) = \left( \frac{n - 2}{\max_i(\mathcal{Q}(Z_i)) - \min_i(\mathcal{Q}(Z_i))} \right) (\mathcal{Q}(Z) - \min_i(\mathcal{Q}(Z_i))) + 1,$$

406 where  $\max_i(\mathcal{Q}(Z_i))$  and  $\min_i(\mathcal{Q}(Z_i))$  are the maximum and minimum values  
407 of the average spread of all available clusterings, and the second term is  
408 the number of clusters in  $Z$ , which is used to compensate the change of  
409 normalized average spread among clusterings having different numbers of  
410 clusters.

## 411 6.2. External measures

412 **Purity** ( $\mathcal{U}, \uparrow$ ). The purity of a clustering  $Z$  with respect to a labeling  
413  $L$  is a measure of the extent to which clusters contain a single class. It is  
414 computed by the formula  $\mathcal{U}(Z) = \frac{1}{n} \sum_{c=1}^k \max_{l \in L} |z_c \cap l|$ , where  $Z$  is a clustering,  
415  $n$  is the total number of observations,  $k$  is the number of clusters,  $z_c$  is the  
416  $c$ -th cluster,  $L$  is the set of classes (i.e., observations with specific labels).  
417 Purity close to  $1/|L|$  represents fragmented clusterings, while purities close  
418 to 1 identify clusterings with almost only one label for each cluster.

419 **Precision** ( $\mathcal{P}, \uparrow$ ). The precision of a cluster  $z_c$  with respect to a label  
420 class  $l$  is a measure of the extent to which the cluster contains the label class.  
421 It is computed as  $\mathcal{P}_l(z_c) = \frac{|z_c \cap l|}{|z_c|}$ , where  $|z_c \cap l|$  is the number of observations  
422 in the intersection between cluster  $z_c$  and label class  $l$ , and  $|z_c|$  is the number  
423 of observations in the cluster  $z_c$ . Values close to 1 are obtained when all the  
424 observations in the cluster correspond to label class  $l$ , values close to 0 are  
425 obtained when no observation in  $z_c$  corresponds to class label  $l$ . We use this

426 measure to find clusters having good match with known states. For instance,  
 427 to find clusters corresponding to drone turning we search clusters  $z_c$  having  
 428  $\mathcal{P}_T(z_c) \geq 0.5$ , where  $\mathcal{P}_T$  is the precision for drone turning.

### 429 6.3. Measures for model variables

430 **Symmetrical uncertainty** ( $\mathcal{SU}, \uparrow$ ). Symmetrical uncertainty (Hong  
 431 et al. (2008)) is a measure of relevance of a variable  $v_d, d \in \{1, \dots, D\}$  with  
 432 respect to a clustering  $Z$  and can be computed as

$$\mathcal{SU}(v_d, Z) = 2 \left( \frac{IG(v_d | Z)}{H(v_d) + H(Z)} \right)$$

433 where  $H(Z)$  is the entropy of the clustering labels and  $IG(v_d | Z)$  is the  
 434 information gain that is computed as  $IG(v_d | Z) = H(v_d) - H(v_d | Z)$ ,  
 435 and  $H(v_d)$  is the entropy of variable  $v_d$  and  $H(v_d | Z)$  is the conditional  
 436 entropy of  $v_d$  given  $Z$ . A value 1 of  $\mathcal{SU}$  indicates that the variable  $v_d$  is  
 437 completely related to clustering  $Z$  while a value 0 means that the variable  $v_d$   
 438 is absolutely irrelevant since it does not share any information with clustering  
 439  $Z$ . It happens for instance, if  $v_d$  is a uniformly distributed random variable.

### 440 6.4. Statistical significance of clusterings and clusters

441 For each internal and external measure defined above it is possible to  
 442 compute the statistical significance, based on p-value, of a clustering  $Z$  with  
 443 respect to the random clustering RC and the random segmentation RS de-  
 444 scribed in Subsections 5.6 and 5.7, respectively. The p-value of a clustering  
 445  $Z$  with respect to a performance measure  $I$  is computed as the percentage of  
 446 random clusterings (random segmentations) that outperform clustering  $Z$  in  
 447 terms performance measure  $I$ . The same approach can be used to compute  
 448 the statistical significance of single clusters. Only clusters/clusterings with  
 449 percentage less than 0.05 are considered statistically significant.

## 450 7. Results and discussion

451 We generate 1076 clusterings of our dataset using the five clustering meth-  
 452 ods described in Section 5 with different parameter settings for each method  
 453 (see Table 3): 126 clusterings are generated by TICC, 300 by IHMM, 100 by  
 454 SCM, 300 by KM and 250 by HMM. The total number of clusters generated  
 455 in this way is 19320 (i.e., 2205 clusters produced by TICC, 5739 by IHMM,  
 456 2376 by SCM, 5250 by KM and 3750 by HMM). To evaluate the statistical

457 significance of clusterings and clusters we compute 200 random clusterings  
 458 (RC) and 200 random segmentations (RS) for each  $k \in \{10, 15, 20, 25, 30\}$ ,  
 459 a total of 1200 random segmentations (21000 random segments) and 1200  
 460 random clusterings (21000 random clusters), and we use them to compute  
 461 clustering and cluster p-values with respect to different performance mea-  
 462 sures. We rank both single clusters and entire clusterings according to their  
 463 performance, and compute their statistical significance with respect to the  
 464 random clusterings/segmentations. In this way, we select a subset of cluster-  
 465 ings and clusters having clear evidence of being non-random and to represent  
 466 drone states. In the following, we first perform an analysis of single cluster  
 467 and then of entire clusterings. We always compare clusters (clusterings)  
 468 having the same parameter  $k$  since all performance measures considered are  
 469 influenced by this parameter. Specific focus is put on  $k = 10$  and  $k = 20$ ,  
 470 two levels of granularity (i.e., abstraction) of interest to discover macroscopic  
 471 states (e.g., in water) and microscopic states (e.g., turning). We notice that  
 472 the extraction of statistically significant state-models is often better achieved  
 473 using cluster validity indices than clustering performance indices, because  
 474 good (e.g., compact and separated) clusters are sometimes present also in  
 475 clusterings having average/low performance, which would not be selected us-  
 476 ing only clustering performance indices. This happens, for instance, when  
 477 a high number of clusters is used, which favours the identification of small  
 478 patterns but also generates non-significant clusters that reduce the overall  
 479 performance of the clustering, even in the presence of good clusters. This  
 480 motivates our choice to analyze deeper single clusters than complete cluster-  
 481 ings, although the analysis of clusterings is an important tool for identifying,  
 482 for instance, the number of clusters in the dataset.

### 483 *7.1. Analysis of single clusters*

484 Clusters are first ranked according to performance measures of Section 6.  
 485 We consider only statistically significant clusters, having p-value less than  
 486 0.05 for at least one performance measure. A summary of properties and  
 487 performance of investigated clusters is reported in Table 4. Figure 5 shows  
 488 the results for two internal measures, i.e., silhouette ( $\mathcal{S}$ ) and weighted spread  
 489 ( $\mathcal{R}$ ), and one external measure, i.e., precision in detecting drone turns ( $\mathcal{P}_T$ ).  
 490 For each performance measure, we show on the left a scatter plot displaying  
 491 all the 61320 clusters (19320 generated by clustering methods, 21000 by RC  
 492 and 21000 by RS) where each point is a cluster, the x-axis is the number of  
 493 states  $k$  in the clustering, and the y-axis is the performance of the cluster.

494 On the right, we display clusters having a specific range of  $k$  and p-value  
495 less than 0.05 for RS. Below, we propose an analysis of few of these clusters,  
496 showing that they have a clear interpretation in terms of drone states.  
497 Further analysis is reported in supplementary material.

498 **Ranking by cluster silhouette.** Figure 5.a shows cluster silhouette  
499 and the ranking by silhouette of clusters with  $k$  between 9 and 11. The  
500 cyan and yellow dashed lines, on the left, characterize the 5th and the 95th  
501 percentile with respect to RS and RC, respectively. Clusters located above  
502 these lines are statistically significant. Focusing on  $k$  between 9 and 11 (see  
503 the blue box on the left of Figure 5.a) we find 249 clusters, of which 27  
504 generated by TICC, 21 by IHMM, 9 by SCM, 100 by KM and 92 by HMM.  
505 These clusters are ranked by silhouette on the right of Figure 5.a where the  
506 point color depends on clustering techniques and point size on cluster size.

507 Clusters  $C_1$  and  $C_2$  have the highest silhouette, respectively 0.76 and 0.68,  
508 and are generated by IHMM. As displayed in Table 4, they have a very small  
509 number of observations, namely three per cluster (see column  $\mathcal{O}$ ), they do not  
510 correspond to a turn ( $\mathcal{P}_T = 0.00$ ), but they correspond to locations in which  
511 the drone was into the water ( $\mathcal{P}_{IW} = 1.00$ ), manually driven ( $\mathcal{P}_{MD} = 1.00$ )  
512 and navigating outside strong streams ( $\mathcal{P}_{NS} = 0.00$ ). Note that information  
513 about precision comes from manual labeling. It is used for result validation  
514 and not provided to the (unsupervised) clustering learning process.

515 We discovered that these clusters identify a real pattern in experiment  
516 ESP4 which can be traced back to a specific (possibly *anomalous*) situation.  
517 The boxplot of variable  $\hat{ec}$  in Figure 6.a shows that clusters  $C_1$  and  $C_2$  have  
518 much higher standard deviation of electrical conductivity than other clusters.  
519 Then, the boxplot of variable  $\tilde{ec}$ , in the same figure, points out that in  $C_1$  the  
520 variation of  $ec$  is positive (increment) and in  $C_2$  it is negative (decrement).  
521 The third and fourth boxplots instead say the two clusters have also high  
522 standard deviation of temperature and voltage. The geolocalization in Figure  
523 6.b shows that cluster  $C_2$  precedes cluster  $C_1$ . All these information, together,  
524 suggest that this pair of clusters could be associated to a location where the  
525 drone was suddenly extracted from and put back into the water. The location  
526 of the clusters is in the middle of a lake, hence the situation could be due to  
527 manual intervention of an operator from a boat, anomalous conditions (e.g.,  
528 obstacles or waves), or sensor faults. It is important to detect such situations  
529 to improve data analysis and avoid misinterpretations of sensor readings.

530 Other key information about this state is provided by the parameters of  
531 the IHMM representing the state-models. Figure 6.c shows the heatmaps of

<b>Id</b>	<b>Clustering method</b>	<b>Selection method</b>	<b>Parameters</b>	<b><math>\mathcal{O}</math></b>	<b><math>\mathcal{S}</math></b>	<b><math>\mathcal{R}</math></b>	<b><math>\mathcal{P}_T</math></b>	<b><math>\mathcal{P}_{IW}</math></b>	<b><math>\mathcal{P}_{MD}</math></b>	<b><math>\mathcal{P}_{US}</math></b>	<b><math>\mathcal{P}_{DS}</math></b>	<b><math>\mathcal{P}_{NS}</math></b>	<b>p-val</b>
$C_1$	IHMM	$\mathcal{S}$ (1st)	$k = 10, \zeta = 30$	3	<b>0.76</b>	21816.9	0.00	<b>1.00</b>	<b>1.00</b>	0.00	0.00	<b>1.00</b>	0.00
$C_2$	IHMM	$\mathcal{S}$ (2nd)	$k = 10, \zeta = 30$	3	<b>0.68</b>	29516.8	0.00	<b>1.00</b>	<b>1.00</b>	0.00	0.00	<b>1.00</b>	0.001
$C_3$	KM	$\mathcal{S}$ (3rd)	$k = 10$	33	<b>0.57</b>	5143.7	0.00	<b>0.64</b>	<b>1.00</b>	0.00	0.00	<b>1.00</b>	0.0025
$C_4$	HMM	$\mathcal{S}$ (53th)	$k = 10$	33	<b>0.57</b>	5143.7	0.00	<b>0.64</b>	<b>1.00</b>	0.00	0.00	<b>1.00</b>	0.0025
$C_5$	SCM	$\mathcal{S}$ (86th)	$NbExpClust = 3$	6774	<b>0.49</b>	17.8	0.02	<b>0.99</b>	0.08	0.00	0.01	<b>0.99</b>	0.005
$C_6$	TICC	$\mathcal{S}$ (246th)	$k = 10, \lambda = 1.0,$ $\beta = 0.0, w = 3.0$	1007	<b>0.21</b>	132.9	0.12	<b>1.00</b>	<b>0.86</b>	<b>0.77</b>	0.00	0.23	0.047
$C_7$	TICC	$\mathcal{R}$ (3th)	$k = 20, \lambda = 1.0,$ $\beta = 50.0, w = 1.0$	8111	0.35	<b>5.32</b>	0.02	<b>0.98</b>	0.16	0.00	0.03	<b>0.97</b>	0.0005
$C_8$	TICC	$\mathcal{R}$ (160th)	$k = 20, \lambda = 0.1,$ $\beta = 200.0, w = 1.0$	4172	-0.024	<b>18.59</b>	0.01	0.22	<b>0.89</b>	0.00	0.00	<b>1.00</b>	0.0305
$C_9$	TICC	$\mathcal{P}_T$ (13th)	$k = 20, \lambda = 0.5,$ $\beta = 100.0, w = 1.0$	317	-0.174	766.20	<b>0.75</b>	<b>1.00</b>	<b>1.00</b>	0.41	0.00	<b>0.59</b>	0.0045
$C_{10}$	SCM	$\mathcal{P}_T$ (287th)	$NbExpClust = 6$	1905	-0.19	68.87	<b>0.39</b>	<b>1.00</b>	<b>0.92</b>	0.10	0.17	<b>0.73</b>	0.027

Table 4: Performance measures and main properties of ten selected clusters.  $Id$  is the cluster identifier, *clustering method* the technique by which the cluster was generated, *selection method* the performance measure by which it was selected (only clusters having p-value less than 0.05 for that measure were considered), *parameters* are the clustering parameters used to generate the cluster,  $\mathcal{O}$  is the number of observations in the cluster,  $\mathcal{S}$  is the cluster silhouette,  $\mathcal{R}$  its weighted spread,  $\mathcal{P}_T$  the cluster precision for drone turns,  $\mathcal{P}_{IW}$  the precision for state “in water” (notice that the precision for the state “out of water” can be calculable as  $1 - \mathcal{P}_{IW}$ ),  $\mathcal{P}_{MD}$  the precision for state “manual drive” (precision of autonomous drive is  $1 - \mathcal{P}_{MD}$ ),  $\mathcal{P}_{US}$  is the precision for state “upstream navigation”,  $\mathcal{P}_{DS}$  is the precision for state “downstream navigation”,  $\mathcal{P}_{NS}$  is the precision for state “no-stream”, and *p-val* is the p-value for RS related to the index in the selection method.

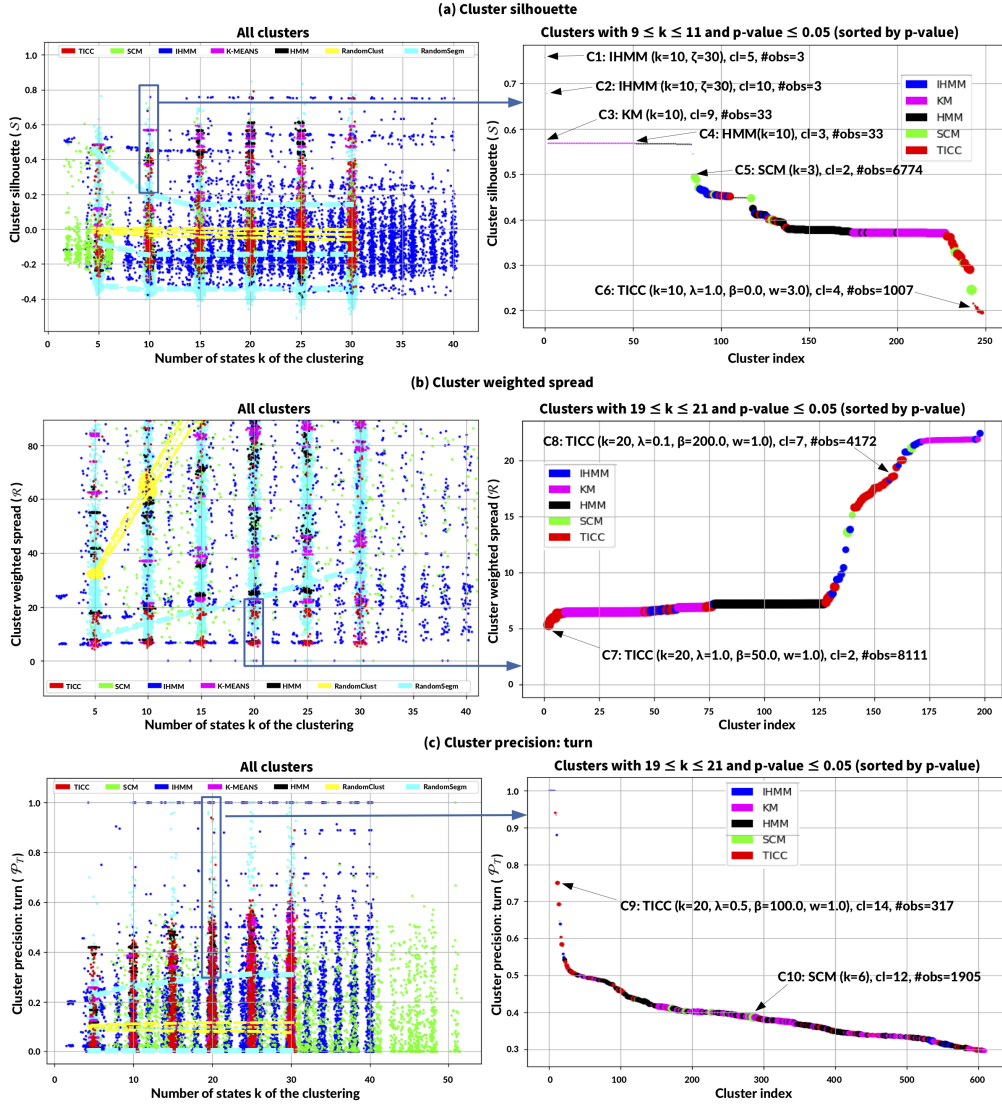


Figure 5: Performance of single clusters (best viewed in colors). Left: X-axes are number of states  $k$  in the clustering, y-axes are values of cluster performance, colors are clustering methods, light blue dashed lines represent 5-th and 95-th percentiles for RS, yellow dashed lines 5-th and 95-th percentiles for RC. Right: statistically significant clusters sorted by performance. (a) Cluster silhouette: significant if above the upper dashed lines. (b) Cluster weighted spread: significant if below the lower dashed lines; only the 5-th percentile line is visible for RS because the figure is zoomed on the lower part of the y-axis. (c) Cluster precision for drone turns: significant if above the upper dashed lines.

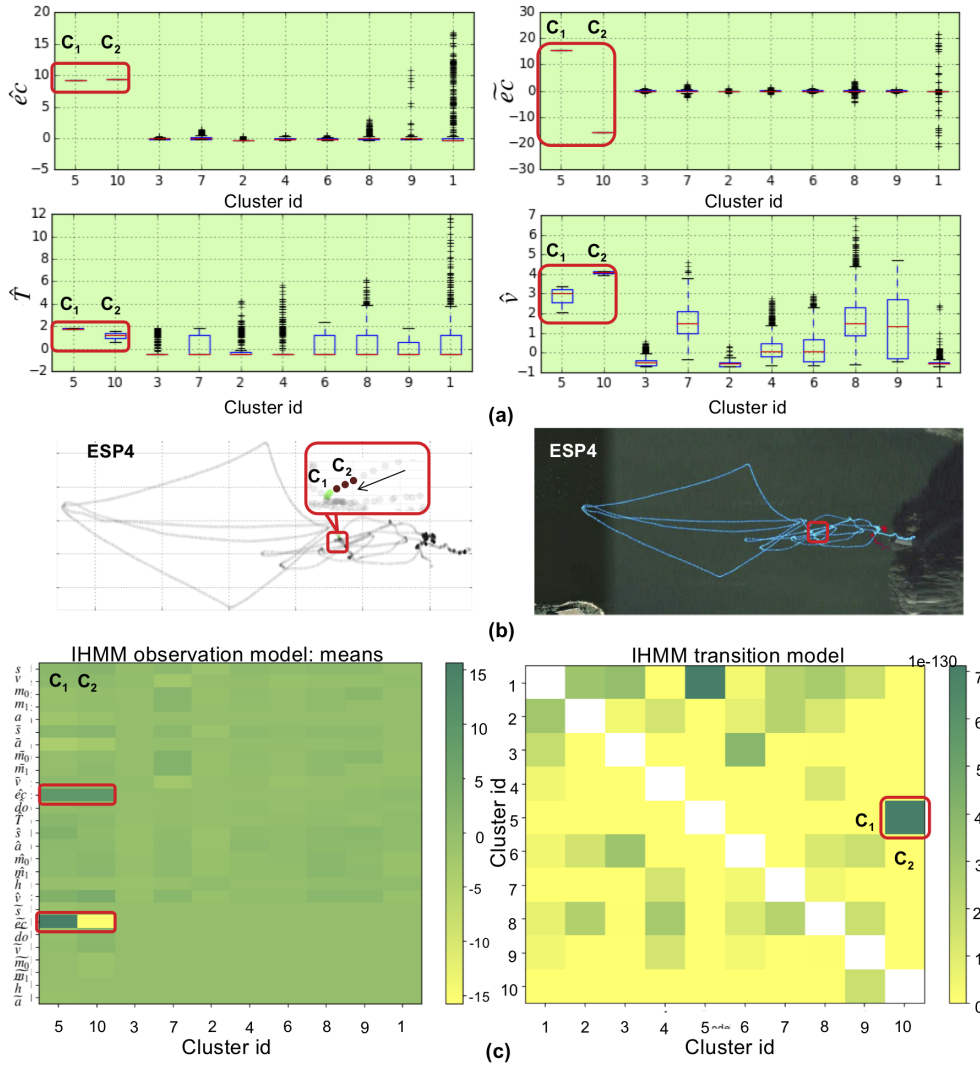


Figure 6: Clusters  $C_1$  and  $C_2$ . (a) Box plots of variables  $\hat{e}c$ ,  $\tilde{e}c$ ,  $\hat{T}$ ,  $\hat{v}$ . (b) Maps of cluster locations. (c) State-model parameters (variable means and transition matrix).

532 variable means for each cluster (on the left) and the transition matrix (on the  
 533 right). Cluster  $C_1$  has strongly positive means for  $\hat{e}c$  and  $\tilde{e}c$  (see dark green  
 534 cells in the first column of the means matrix) and cluster  $C_2$  has strongly  
 535 positive mean for  $\hat{e}c$  and strongly negative mean for  $\tilde{e}c$  (second column of  
 536 the means matrix). Moreover, the switch between cluster  $C_2$  and cluster  $C_1$   
 537 is represented by the high parameter in the highlighted cell of the transition



538 matrix (on the right). We reported other analysis on clusters  $C_3$  to  $C_6$  in the  
539 supplementary material.

540 **Ranking by cluster weighted spread.** This ranking of clusters is dis-  
541 played in Figure 5.b. On the right we show the significant clusters with  $k$   
542 between 19 and 21. We found 199 significant clusters, of which 42 generated  
543 by TICC, 29 by IHMM, 3 by SCM, 75 by KM and 50 by HMM. Cluster  $C_7$   
544 has almost the best performance in the ranking (two other clusters perform  
545 better but they contain only one observation). It was generated by TICC,  
546 contains 8111 observations, has weighted spread 5.32 and silhouette 0.35.  
547 This cluster corresponds to observations in which the drone was into the wa-  
548 ter (i.e.,  $\mathcal{P}_{IW} = 0.98$ ), autonomously driven (i.e.,  $\mathcal{P}_{MD} = 0.16$ ), not in strong  
549 streams (i.e.,  $\mathcal{P}_{NS} = 0.97$ ) and not turning (i.e.,  $\mathcal{P}_T = 0.02$ ). Interestingly  
550 enough, this cluster contains almost the same points of cluster  $C_5$ , which was  
551 generated by SubCMedians and selected from the silhouette ranking. This  
552 shows that different clustering methods (i.e., SubCMedians and TICC in this  
553 case) were able to discover the same state of the drone although using differ-  
554 ent state representations (i.e., centroids and Toeplitz matrices). Cluster  $C_8$   
555 is analyzed in the supplementary material.

556 **Ranking by cluster precision for drone turning.** The third ranking  
557 we analyze is based on the precision to detect drone turns. A scatter plot  
558 of clusters arranged by  $k$  (x-axis) and precision to detect drone turns  $\mathcal{P}_T$   
559 (y-axis) is displayed on the left of Figure 5.c. We focus, in particular, on  $k$   
560 between 19 and 21. These clusters are 609 in total, of which 101 generated  
561 by TICC, 36 by IHMM, 17 by SCM, 212 by KM and 243 by HMM. The best  
562 15 clusters, having  $\mathcal{P}_T \geq 0.69$ , are all generated by TICC or IHMM that  
563 seem to have the best capability to detect drone turns.

564 Cluster  $C_9$  is the first “large” cluster in the ranking (317 observations) and  
565 it is generated by TICC. Its precision on drone turns  $\mathcal{P}_T$  is 0.75, meaning that  
566 the 75% of its observations in the cluster correspond to real turn, according  
567 to our manual labeling. According to Table 4 this cluster corresponds to  
568 observations taken into the water (i.e.,  $\mathcal{P}_{IW} = 1.00$ ) during manual drive (i.e.,  
569  $\mathcal{P}_{MD} = 1.00$ ), partially in upstream navigation and partially with no stream  
570 (i.e.,  $\mathcal{P}_{US} = 0.41$  and  $\mathcal{P}_{NS} = 0.59$ ). Among the main statistical properties  
571 of variables characterizing this clusters there are high standard deviation of  
572 signal to propellers  $\hat{m}_0$  (and  $\hat{m}_1$ ), and high standard deviation of voltage  
573  $\hat{v}$ , as shown in the two boxplots of Figure 7.a. The geolocalization of this  
574 cluster confirms its correspondence to curves in the drone path, as shown  
575 in Figure 7.b that displays five locations belonging to three campaigns (i.e.,

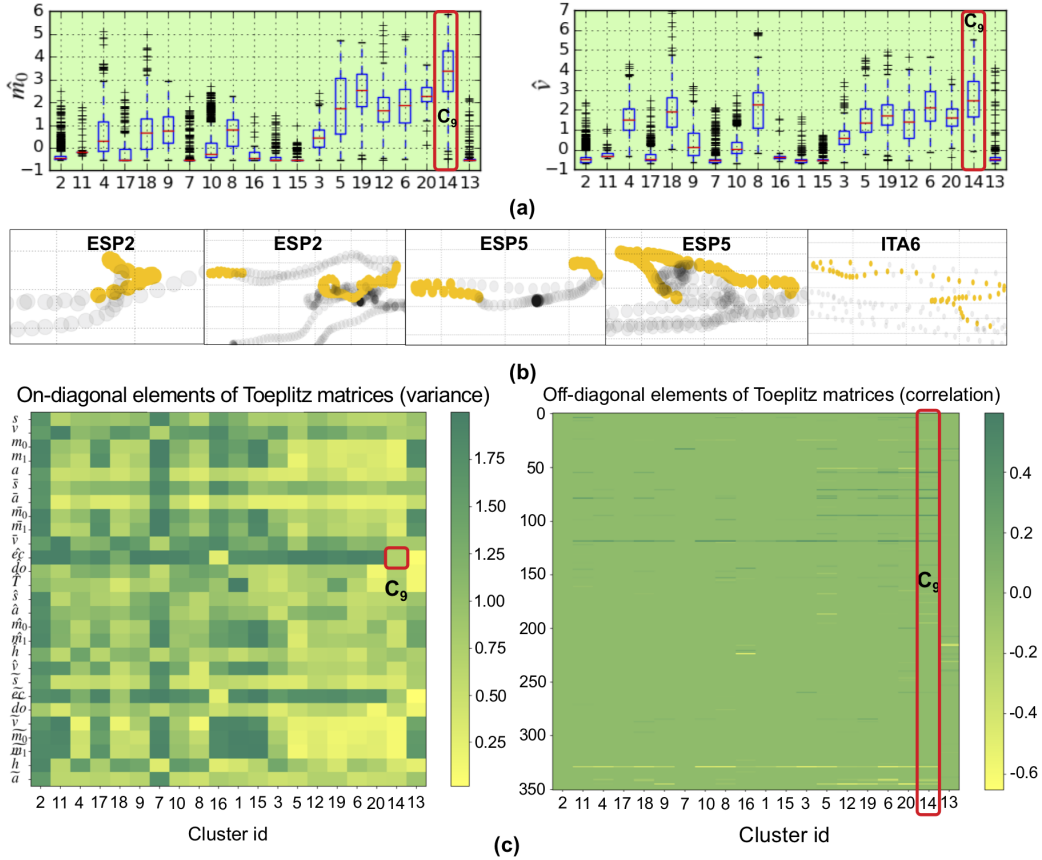


Figure 7: Clusters  $C_9$ . (a) Box plot of variables  $\hat{m}_0, \hat{v}$ . (b) Maps showing cluster locations. (c) State-model parameters (on-diagonal and off-diagonal elements of Toeplitz matrices).

576 ESP2, ESP5 and ITA6). We observe that the cluster really characterizes  
 577 the turning pattern in the data. Figure 7.c shows the on-diagonal elements  
 578 (on the left) and the off-diagonal elements (on the right) of the Toeplitz  
 579 matrix representing this state. Cluster  $C_{10}$  is analyzed in the supplementary  
 580 material.

## 581 7.2. Analysis of clusterings

582 Here we perform a second kind of analysis based on clustering significance  
 583 (the previous one was on cluster significance). We evaluate our clusterings,  
 584 computed by different methods and different parameter settings, according  
 585 to four internal measures, namely silhouette ( $\mathcal{S}$ ), Davis-Bouldin index ( $\mathcal{DB}$ ),

586 weighted spread ( $\mathcal{R}$ ), and Calinski-Harabaz index ( $\mathcal{CH}$ ). Results are sum-  
 587 marized in Figure 8, which has a similar structure to Figure 5. Scatter plots,  
 588 on the left, contain one point for each clustering. The x-axis represents the  
 589 number of clusters  $k$  in the clustering and the y-axis the performance mea-  
 590 sure of interest. Point colors correspond to different clustering methods. On  
 591 the right hand side some selections of significant clusterings, with specific  $k$   
 592 and p-value less than or equal to 0.05, are displayed by ascending/descending  
 593 performance.

594 Clustering silhouette is displayed in Figure 8.a. As expected the best sil-  
 595 houette is achieved by clustering with small number of clusters (e.g.,  $k = 2$  for  
 596 IHMM,  $k = 5$  for k-means and TICC,  $k = 6$  for SCM). The average clustering  
 597 silhouette however increases from  $k = 10$  to  $k = 25$  and then it decreases for  
 598  $k > 25$ , showing a peak around  $k = 25$  for all methodologies. This is interest-  
 599 ing because it suggests a best number of clusters (around 25) for this dataset.  
 600 Moreover, silhouette of SCM and IHMM with  $k > 30$  sharply degrades to  
 601 zero or less than zero. Surprisingly, the best silhouette is achieved by k-means  
 602 for all  $k$  (see pink points in the chart). Then TICC reaches the second best  
 603 silhouette performance, followed by SubCMedians and IHMM that has simi-  
 604 lar average performance to HMM but better performance considering the  
 605 best parameter settings. The silhouette of non-random clusterings is almost  
 606 always higher than silhouette of random segmentations. This behavior is  
 607 very different from that observed for clusters, wherein several superpositions  
 608 were present. Ranking by silhouette of clusterings with  $k$  between 9 and 11  
 609 (on the right of Figure 8.a) show that the best clustering was generated by  
 610 SCM and has a silhouette of 0.17. It is followed by k-means (about 0.15)  
 611 and TICC (about 0.14), then there is a big jump to reach the best IHMM  
 612 clustering, having silhouette 0.08, and HMM with silhouette 0.07.

613 The Davis-Bouldin index, in Figure 8.b, is again dominated by k-means  
 614 (see the pink points in the chart) that shows, as for silhouette, an optimum  
 615 (i.e., a minimum for Davis-Bouldin index) in  $k$  between 20 and 25. The  
 616 performance of the other methods (considering the best models for each  
 617 technique while  $k$  varies between 5 and 30) are quite constants over  $k$ , with  
 618 best performance achieved mainly by TICC, SCM and IHMM depending on  
 619  $k$ . Not considering small  $k$ , TICC has its best performance in  $k = 25$ , IHMM  
 620 and HMM in  $k = 20$ , SCM in  $k = 39$  (with small differences with other  
 621  $k$ ). All points are below the cyan and yellow points of RS and RC (yellow  
 622 points are not displayed because of too high values). Weighted spread and  
 623 Calinski-Harabaz indices are analyzed in supplementary material.

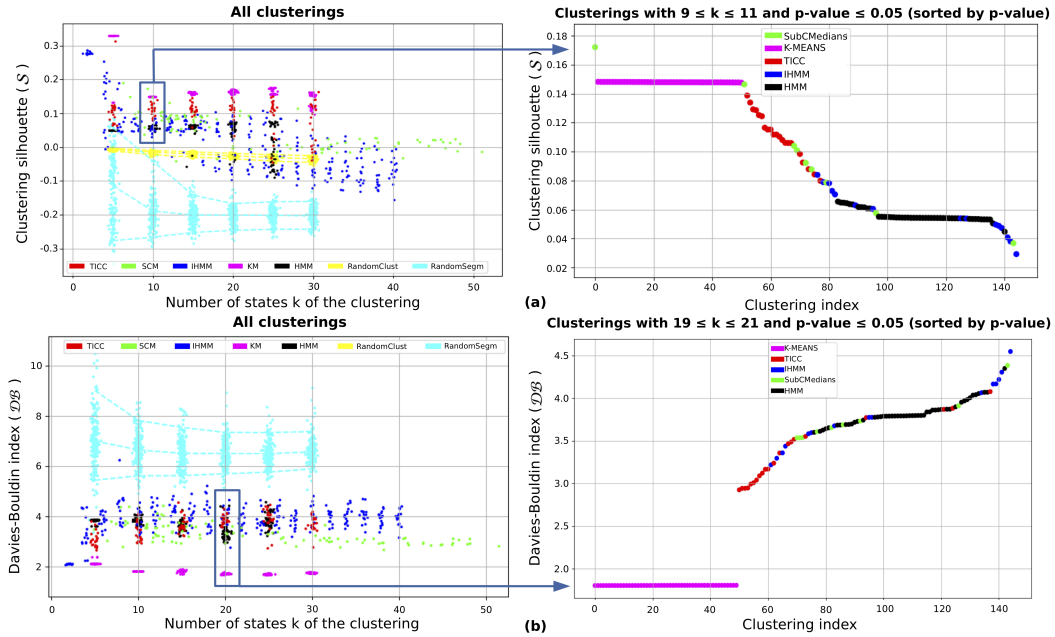


Figure 8: Performance of clusterings. Left: x-axis is the number of states  $k$ , y-axis is the performance value, colors are clustering methods. Each point is a clustering. Right: significant clusterings sorted by performance. (a) silhouette, (b) Davis-Bouldin index.

624 A final comment is focused on clustering p-values. Differently from cluster-  
625 ters, clusterings are almost all statistically significant with respect to RC  
626 and RS. This holds for all the four internal performance measures analyzed  
627 in this section, as displayed in Figure 8, where the points related to non-  
628 random clusterings are almost always out of the areas delimited by the 5th  
629 and 95th percentile lines (yellow and cyan dashed lines). This is possibly  
630 due to the fact that randomly generate clusterings with performance similar  
631 to that of state-of-the-art clustering algorithms is more difficult than ran-  
632 domly generate single clusters with performance similar to that generated by  
633 state-of-the-art methods.

## 634 8. Conclusions and future work

635 The framework proposed in this work allows to identify significant states  
636 of aquatic drones involved in water monitoring by means of diverse unsu-  
637 pervised clustering and segmentation methodologies. The analysis of the  
638 models of these states, namely, centroids, Toeplitz matrices, and multivari-

639 ate Gaussian distributions (depending on the methodology that generated  
640 them), allows us to discover the statistical properties that characterize some  
641 of these states and, consequently, to provide interpretations for the related  
642 models. This result has direct consequences on the analysis of the data ac-  
643 quired by the drones since we can now label the dataset by discovered states,  
644 obtaining a compact semantic-based way to represent each campaign. This  
645 could have strong impact on water monitoring projects involving the citi-  
646 zenship in collecting evidence about water healthiness (following the citizen  
647 science approach), since unskilled people need support in data interpretation.

648 From a more general point of view, the proposed framework represents an  
649 easy-to-use tool for discovering significant states in multivariate time series  
650 datasets and for comparing the capabilities of different clustering techniques.  
651 It only needs a dataset and a set of parameter settings for each methodology,  
652 and produces several rankings of clusterings/clusters with associated signif-  
653 icance levels, allowing to compare the performance of different methods to  
654 identify states in specific application domains (and related datasets). The  
655 choice of a clustering/segmentation method for real datasets is a challenging  
656 activity and our approach could provide valuable support in this direction.

657 Future activities will aim to release an easy-to-use software for supporting  
658 the proposed framework. Then we want to merge the clusters discovered  
659 by different methods using different levels of granularity (i.e., parameter  $k$ )  
660 into a hierarchical (voting) structure, so that each observation could be part  
661 of several clusters of different abstraction levels (e.g., drone into the water,  
662 turning and moving upstream). Another goal is to focus on specific situations  
663 of interest, such as anomalies and dangerous states (e.g., high waves). We are  
664 planning specific field tests to this purpose. Finally, we want to integrate our  
665 state recognition method into online sequential decision making algorithms,  
666 such as those based on Partially Observed Markov Decision Processes (known  
667 as POMDPs) that we started to develop in (Castellini et al. (2019b)). This  
668 direction could improve drone autonomy by supporting the generation of  
669 policies based on improved system states.

## 670 **9. Acknowledgements**

671 This work is partially funded by the European Union’s Horizon 2020  
672 research and innovation programme under grant agreement No 689341. This  
673 work reflects only the authors’ view and the EASME is not responsible for  
674 any use that may be made of the information it contains.

675 **10. References**

- 676 Abdallah, Z. S., Gaber, M. M., Srinivasan, B., Krishnaswamy, S., 2012.  
677 CBARS: Cluster Based Classification for Activity Recognition Systems.  
678 Springer Berlin Heidelberg, pp. 82–91.
- 679 Arbelaitz, O., Gurrutxaga, I., Muguerza, J., Prez, J. M., Perona, I., 2013. An  
680 extensive comparative study of cluster validity indices. *Pattern Recognition*  
681 46 (1), 243 – 256.
- 682 Asperti, A., Cortesi, D., Sovrano, F., 2019. Crawling in Rogue’s dungeons  
683 with (partitioned) A3C. In: *The 4th Int. Conf. Machine Learning, Opti-*  
684 *mization and Data science (LOD 2018)*, Volterra, Italy. Springer.
- 685 Barnett, I., Onnela, J.-P., 2016. Change point detection in correlation net-  
686 works. *Scientific Reports* 6, 18893.
- 687 Barták, R., Vomlelová, M., 2017. Using machine learning to identify activities  
688 of a flying drone from sensor readings. In: *Proceedings of Florida Artificial*  
689 *Intelligence Research Society Conference, FLAIRS 2017*. pp. 436–441.
- 690 Bishop, C. M., 2006. *Pattern Recognition and Machine Learning (Informa-*  
691 *tion Science and Statistics)*. Springer-Verlag New York, USA.
- 692 Bottarelli, L., Bicego, M., Blum, J., Farinelli, A., 2016. Skeleton-based orien-  
693 teering for level set estimation. In: *ECAI 2016 - 22nd European Conference*  
694 *on Artificial Intelligence*. pp. 1256–1264.
- 695 Bottarelli, L., Bicego, M., Blum, J., Farinelli, A., 2019. Orienteering-based  
696 informative path planning for environmental monitoring. *Engineering Ap-*  
697 *plications of Artificial Intelligence* 77, 46–58.
- 698 Caliński, T., Harabasz, J., 1974. A dendrite method for cluster analysis.  
699 *Communications in Statistics-Simulation and Computation* 3 (1), 1–27.
- 700 Castellini, A., Beltrame, G., Bicego, M., Bloisi, D., Blum, J., Denitto, M.,  
701 Farinelli, A., 2018a. Activity recognition for autonomous water drones  
702 based on unsupervised learning methods. In: *Proc. 4th Italian Workshop*  
703 *on Artificial Intelligence and Robotics (AI\*IA 2017)*. Vol. 2054. pp. 16–21.

- 704 Castellini, A., Beltrame, G., Bicego, M., Blum, J., Denitto, M., Farinelli, A.,  
705 2018b. Unsupervised activity recognition for autonomous water drones. In:  
706 Proc. Symposium on Applied Computing, SAC 2018. ACM, pp. 840–842.
- 707 Castellini, A., Bicego, M., Bloisi, D., Blum, J., Masillo, F., Peignier, S.,  
708 Farinelli, A., 2019a. Subspace clustering for situation assessment in aquatic  
709 drones: A sensitivity analysis for state-model improvement. *Cybernetics  
710 and Systems* 50 (8), 658–671.
- 711 Castellini, A., Chalkiadakis, G., Farinelli, A., 2019b. Influence of State-  
712 Variable Constraints on Partially Observable Monte Carlo Planning. In:  
713 Proc. 28th International Joint Conference on Artificial Intelligence (IJCAI  
714 2019). pp. 5540–5546.
- 715 Castellini, A., Masillo, F., Bicego, M., Bloisi, D., Blum, J., Farinelli, A.,  
716 Peigner, S., 2019c. Subspace clustering for situation assessment in aquatic  
717 drones. In: Proc. Symposium on Applied Computing, SAC 2019. ACM,  
718 pp. 930–937.
- 719 Castellini, A., Masillo, F., Sarteau, R., Farinelli, A., 2019d. eXplainable Mod-  
720 eling (XM): Data Analysis for Intelligent Agents. In: Proceedings of the  
721 18th International Conference on Autonomous Agents and Multiagent Sys-  
722 tems (AAMAS 2019). IFAAMAS, pp. 2342–2344.
- 723 Castellini, A., Paltrinieri, D., Manca, V., 2015. MP-GeneticSynth: inferring  
724 biological network regulations from time series. *Bioinformatics* 31, 785–87.
- 725 Chen, L., Hoey, J., Nugent, C. D., Cook, D. J., Yu, Z., 2012. Sensor-based ac-  
726 tivity recognition. *IEEE Transactions on Systems, Man, and Cybernetics,  
727 Part C (Applications and Reviews)* 42 (6), 790–808.
- 728 Chiu, B., Keogh, E., Lonardi, S., 2003. Probabilistic discovery of time series  
729 motifs. In: Proc. 9th ACM SIGKDD Int. Conf. on Knowledge Discovery  
730 and Data Mining. KDD '03. ACM, New York, USA, pp. 493–498.
- 731 Davies, D. L., Bouldin, D. W., Feb. 1979. A cluster separation measure. *IEEE  
732 Trans. Pattern Analysis Machine Intelligence* 1 (2), 224–227.
- 733 Dhiman, C., Vishwakarma, D. K., 2019. A review of state-of-the-art tech-  
734 niques for abnormal human activity recognition. *Engineering Applications  
735 of Artificial Intelligence* 77, 21 – 45.

- 736 Endsley, M. R., 1995. Toward a theory of situation awareness in dynamic  
737 systems. *Human Factors* 37 (1), 32–64.
- 738 Farinelli, A., Nardi, D., Pigliacampo, R., Rossi, M., Settembre, G. P., 2012.  
739 Cooperative situation assessment in a maritime scenario. *International*  
740 *Journal of Intelligent Systems* 27 (5), 477–501.
- 741 Fox, E. B., Sudderth, E. B., Jordan, M. I., Willsky, A. S., 2008. An HDP-  
742 HMM for systems with state persistence. In: *Proceedings of the 25th Inter-*  
743 *national Conference on Machine Learning. ICML '08. ACM*, pp. 312–319.
- 744 Fu, T.-c., 2011. A review on time series data mining. *Engineering Applica-*  
745 *tions of Artificial Intelligence* 24 (1), 164 – 181.
- 746 Hallac, D., Bhooshan, S., Chen, M., Abida, K., Sasic, R., Leskovec, J., 2018.  
747 Drive2vec: Multiscale state-space embedding of vehicular sensor data. In:  
748 *Int. Conf. Intelligent Transportation Systems. IEEE*, pp. 3233–3238.
- 749 Hallac, D., Nystrup, P., Boyd, S., 2016a. Greedy gaussian segmentation  
750 of multivariate time series. *Advances in Data Analysis and Classification*  
751 13 (3), 727–751.
- 752 Hallac, D., Sharang, A., Stahlmann, R., Lamprecht, A., Huber, M., Roehder,  
753 M., Sasic, R., Leskovec, J., 2016b. Driver identification using automobile  
754 sensor data from a single turn. In: *19th Int. Conf. Intelligent Transporta-*  
755 *tion Systems. IEEE*, pp. 953–958.
- 756 Hallac, D., Vare, S., Boyd, S., Leskovec, J., 2017. Toeplitz inverse covariance-  
757 based clustering of multivariate time series data. In: *Proc. 23rd ACM*  
758 *SIGKDD. KDD '17. ACM*, pp. 215–223.
- 759 Hastie, T., Tibshirani, R., Friedman, J., 2001. *The elements of statistical*  
760 *learning. Springer Series in Statistics. Springer, New York, USA.*
- 761 Hong, Y., Kwong, S., Chang, Y., Ren, Q., 2008. Consensus unsupervised  
762 feature ranking from multiple views. *Pattern Rec. Let.* 29 (5), 595 – 602.
- 763 Kaelbling, L. P., Lozano-Perez, T., 2013. Integrated task and motion plan-  
764 ning in belief space. *International Journal of Robotics Research* 32 (9-10).



- 765 Kelley, L. A., Gardner, S. P., Sutcliffe, M. J., 11 1996. An automated ap-  
766 proach for clustering an ensemble of NMR-derived protein structures into  
767 conformationally related subfamilies. *Protein Engineering, Design and Se-*  
768 *lection* 9 (11), 1063–1065.
- 769 Kim, E., Helal, S., Cook, D., 2010. Human activity recognition and pattern  
770 discovery. *IEEE Pervasive Computing* 9 (1), 48–53.
- 771 Kwon, Y., Kang, K., Bae, C., 2014. Unsupervised learning for human activity  
772 recognition using smartphone sensors. *Expert Systems with Applications*  
773 41 (14), 6067 – 6074.
- 774 Montanez, G., Amizadeh, S., Laptev, N., 2015. Inertial Hidden Markov Mod-  
775 els: Modeling change in multivariate time series. In: *Proc. AAAI Conf.*  
776 *Artificial Intelligence. AAAI '15.* pp. 911–916.
- 777 Moshtaghi, M., Bezdek, J. C., Erfani, S. M., Leckie, C., Bailey, J., 2019.  
778 Online cluster validity indices for performance monitoring of streaming  
779 data clustering. *Int. Journal of Intelligent Systems* 34 (4), 541–563.
- 780 Peignier, S., Rigotti, C., Rossi, A., Beslon, G., 2018. Weight-based search to  
781 find clusters around medians in subspaces. In: *Proceedings of the Sympo-*  
782 *sium on Applied Computing, SAC 2018.* ACM, pp. 471–480.
- 783 Rabiner, L. R., Feb 1989. A tutorial on hidden markov models and selected  
784 applications in speech recognition. *Proc. of the IEEE* 77 (2), 257–286.
- 785 Rousseeuw, P. J., 1987. Silhouettes: A graphical aid to the interpretation  
786 and validation of cluster analysis. *Journal of Computational and Applied*  
787 *Mathematics* 20, 53 – 65.
- 788 Russell, S., Norvig, P., 2009. *Artificial Intelligence: A Modern Approach*, 3rd  
789 Edition. Prentice Hall Press, Upper Saddle River, NJ, USA.
- 790 Trabelsi, D., Mohammed, S., Chamroukhi, F., Oukhellou, L., Amirat, Y.,  
791 2013. An unsupervised approach for automatic activity recognition based  
792 on hidden markov model regression. *IEEE Trans. Automation Science and*  
793 *Engineering* 10 (3), 829–835.
- 794 van der Maaten, L., Hinton, G., 2008. Visualizing data using t-SNE. *Journal*  
795 *of Machine Learning Research* 9, 2579–2605.