

1 **Using General linear model, Bayesian networks and Naive Bayes classifier for prediction**
2 **of *Karenia selliformis* occurrences and blooms**

3

4 **Wafa Feki-Sahnoun^{1*}, Hasna Njah², Asma Hamza¹, Nouha Barraï³, Mabrouka**
5 **Mahfoudi¹, Ahmed Rebai⁴, Malika Bel Hassen³**

6

7 ¹ Institut National des Sciences et Technologies de la Mer, Centre de Sfax, Rue Madagascar,
8 BP 1035, Sfax, CP 3018, Tunisie.

9 ² Faculté des Sciences Économiques et de Gestion de Sfax, Route de l'Aéroport Km 4 Sfax
10 3018 Tunisie. Laboratoire de Multimedia, InfoRmation Systems and Advanced Computing
11 Laboratory, Pôle technologique de Sfax, Route de Tunis Km 10 BP 242 CP 3021, Sfax,
12 Tunisie

13 ³ Institut National des Sciences et Technologies de la Mer (INSTM), 28 rue 2 mars 1934,
14 Salammbô 2025, Tunisie

15 ⁴ Centre de Biotechnologie de Sfax, Route Sidi Mansour Km 6, BP 1177, 3018 Sfax, Tunisie

16

17 *Corresponding Author: **Wafa Feki-Sahnoun**

18 wafafeki@yahoo.fr

19 Tel: +216 20247428; fax: +216 74497989.

20 njah.hasna@gmail.com (H. Njah)

21 asma.hamza@instm.rnrt.tn (A. Hamza)

22 barraï.nouha@instm.rnrt.tn (N. Barraï)

23 mabroukamahfoudi@yahoo.fr (M. Mahfoudi)

24 ahmed.rebai@cbs.rnrt.tn (A. Rebai)

25 belhassen.malika@instm.rnrt.tn (M. Bel Hassen)

26 **Abbreviations**

27 Generalized Linear Mixed-effect Model (GLMM)

28 *Karenia selliformis* (*K. selliformis*)

29 Analysis of Variance (ANOVA)

30 General Linear Model (GLM)

31 Bayesian Network (BN)

32 Naive Bayes classifiers (NB)

33 Tunisian National Meteorological Institute (INM)

34 Akaike Information Criterion (AIC)

35 Bayesian Information Criteria (BIC)

36 SamIam (Sensitivity Analysis, Modeling, Inference And More software)

37 Directed Acyclic Graph (DAG)

38 Conditional Probability Tables (CPTs).

39 Maximum likelihood (ML)

40 Maximum A Posteriori (MAP)

41 Evaporation (Evap)

42 Insolation (Insol)

43 Salinity (Sal)

44 Humidity (Humid)

45 Water Temperature (WatT)

46

47

48

49

50 **Research highlights**

51

- 52 ▪ Bayesian networks (BN) performed better than Naive Bayes and General linear
- 53 model.
- 54 ▪ All models incriminate high salinity in the blooms and site for the occurrences.
- 55 ▪ In BN, all relationships are more explanatory due to their inter-dependencies.
- 56 ▪ A non-linear behavior between salinity and species concentrations was observed.

57

58 **Abstract**

59 The prediction of the dinoflagellate red tide forming *Karenia selliformis* is a relevant task to
60 aid optimized management decisions in marine coastal water. The objective of the present study
61 is to compare different modeling approaches for prediction of *Karenia selliformis* occurrences
62 and blooms. A set of physical parameters (salinity, temperature and tide amplitude),
63 meteorological constraints (evaporation, air temperature, insolation, rainfall, atmospheric
64 pressure and humidity), sampling months and sampling sites are used. The model prediction
65 included General Linear Model (GLM), Bayesian Network (BN) and the simplest BN type
66 which is, Naive Bayes classifier (NB). The results showed that three models incriminated high
67 salinity in *Karenia selliformis* blooms and the sampling sites, mainly Boughrara lagoon, in the
68 occurrences. The BN performed better than linear models (NB and GLM) for both *Karenia*
69 *selliformis* occurrences and blooms prediction. This later is related to the facts that BN
70 considered the inter-independency between predictive variables and that the relationships
71 between the variables and the outcome are often non-linear such us; the transition to bloom
72 situations appeared to be triggered by a salinity threshold. This study is useful in the

73 management of this ecosystem so as to use the best disposal options in the early prediction of
74 the toxic blooms.

75 **Keywords:** *Karenia selliformis*, Naive Bayes classifier, General linear model, Bayesian
76 Network, hydro-meteorological parameters, the Gulf of Gabès.

77

78 **1. Introduction**

79 The harmful *Karenia* species are found throughout the world. They have been described
80 as a result of investigations into extensive animal mortalities or human health problems. They
81 have become the most studied species of harmful algae with extensive investigations on the
82 physiology and bloom formation. *Karenia selliformis* (Hansen et al. 2004) has been the most
83 abundant species causing severe harmful blooms in the Gulf of Gabès (Hamza and El Abed,
84 1994). Since the year 1990, *K. selliformis* blooms have occurred annually along the coast, and
85 represented over 64% of the reported blooms in this area (Feki et al. 2008, 2013). It has been
86 argued that their proliferation has been usually related to shellfish toxicity (Ben Naila et al.
87 2012; Marrouchi et al. 2009; Medhioub et al. 2010). Little is known about the effect of
88 environmental factors on *K. selliformis* occurrences and blooms. Information on optimal growth
89 of this species under controlled temperature and salinity was documented in culture experiments
90 (Medhioub et al. 2009). Tentatively models have been developed to apprehend the effects of
91 physical and meteorological variables on *Ks* occurrences and blooms in the Gulf of Gabès. A
92 generalized linear mixed-effect model (GLMM) incriminated mainly water temperature and
93 some nutrients in the spatiotemporal occurrences of *Karenia selliformis* (Feki et al. 2013).
94 Bayesian network approach showed that the bloom can be predicted based on salinity threshold

95 (Feki-Sahnoun et al. 2017). However, the best performance model having the highest goodness
96 of fit on the prediction of *Karenia* blooms and occurrences still needs to be determined.

97 To date, a large variety of statistical models are available to analyze a relationship
98 between environmental variables and species distributions (i.e. cross validation criteria,
99 ANOVA) and one of the most popular is the General linear model (GLM) (e.g. McCulloch et
100 al. 2008). Typically, the biological and physical processes, generating this data, are highly
101 complex, resulting in multiple correlations/dependencies between covariates and also between
102 outcome variables. Standard statistical approaches have a limited ability to describe such
103 interdependent multi-factorial relationships. In the last decades, Bayesian Network's modelling
104 (BN) has been widely used in solving environmental problems (Borsuk et al. 2004, 2006;
105 Bromley et al. 2005; Feki-Sahnoun et al. 2017; Pollino et al. 2007; Smith et al. 2007; Zaffalon,
106 2005) to analyze multi-dimensional data. Among the BN models, the naïve Bayes (NB) model
107 appears to be the most popular (Aguilera et al. 2013; Fytilis and Rizzo, 2013; Markus et al.
108 2010; Roperio et al. 2014, 2015). Despite its linearity and simplicity, the NB has been found to
109 perform surprisingly well (Friedman et al. 1997) particularly in many complex situations
110 (Domingos and Pazzani, 1997; Zhang, 2004; Boets et al. 2015).

111 The present study compares results from Bayesian networks [BN, Bishop, 2006; Pearl,
112 1985], naive Bayes classifier [NB, Friedman et al. 1997] and General linear model [GLM,
113 McCullagh and Nelder, 1989; Nelder and Wedderburn, 1972] in the modeling of *Ks*
114 occurrences and blooms; which is based on a set of physical and meteorological variables. In
115 contrast to Feki-Sahnoun et al. (2017) who used the same dataset and focused only on the
116 Bayesian network, this research elaborated upon the comparison between the three models
117 based on a new threshold concentration for bloom and non-bloom conditions established using
118 the Feki-Sahnoun et al. (2017) output results. The choice of these models is due to the fact that
119 GLM analysis is a well-known method that allows predicting a target variable (in this case,

120 *Karenia selliformis*) conditionally on a set of variables by fitting a regression model using least
121 squares (Hastie et al. 2009). Many different models could be fitted by including different
122 subsets of the available observed variables and interaction terms between them. The challenge
123 in implementing this technique is the selection of the best subset of variables, as the inclusion
124 of correlated predictors may result in increased standard errors of the regression coefficients
125 which cause prediction's sensitivity to model changes (Burnham and Anderson, 2002).
126 Whereas GLMs are typically data-driven models techniques (Hastie and Tibshirani, 1990;
127 Madsen and Thyregod, 2011; McCullagh and Nelder, 1989; Nelder and Wedderburn, 1972),
128 BNs can be based on expert knowledge and/or stakeholders (Aguilera et al. 2011; Cain et al.
129 2003; Chan et al. 2012; Haapasaari and Karjalainen, 2010; Haines-Young, 2011; McVittie et
130 al. 2015; Wang et al. 2009) in order to analyse how a system works, to seek information on the
131 appropriate probabilities, or to explore the usefulness of the decision process (Gonzalez-Redin
132 et al. 2016). BN is an alternative modeling approach applying different criteria for model
133 selection. This method consists of a graphical modeling tool that can be used to construct a
134 predictive model by factorizing the posterior density distribution function assuming a set of
135 stable conditional dependencies (Jensen, 1996; Lauritzen and Spiegelhater, 1988; Pearl, 1988).
136 Among BN models, naive Bayes classifier is the simplest; because of the fact that its variables
137 are conditionally independent given the class variable. It has the strength to not only provide a
138 prediction, but also the estimated probability associated with each possible outcome (Fernandez
139 et al. 2010).

140 The objective of the present study is to compare the efficiency of the linear (GLM, NB)
141 and non-linear (BN) models in order to predict *K. selliformis* occurrences and blooms in the
142 Gulf of Gabès using a set of meteorological (rainfall, atmospheric pressure, insolation, ...) and
143 physical (temperature, salinity, tide) variables. By comparing BN with GLM and NB using

144 identical data, the likely impact of such an analytical difference on the inferences was
145 highlighted.

146 **2. Material and methods**

147 **2.1. Study area**

148 The present study was carried out in the Gulf of Gabès (Tunisia, Eastern Mediterranean
149 Sea), in a wide continental shelf area extending from 9.5 to 12°E in longitude and from 33 to
150 35.5°N in latitude and sheltering various islands (Kerkennah and Djerba) and lagoons
151 (Boughrara and El Bibane). The Boughrara lagoon is surrounded by solar saltern areas where
152 net evaporation is very high and soil washing is negligible due to limited freshwater supplies
153 coupled with scanty rainfall.

154 The bathymetry is characterized by a low slope; the 50-m depth is reached at a distance
155 of 100 km (Fig. 1). The tide is semidiurnal, being among the highest in the Mediterranean and
156 having a maximum range of about 2 m generated by resonance phenomena (Sammari et al.
157 2006). The climate is dry and sunny with strong easterly winds.

158 Despite being oligotrophic, this Gulf has the peculiarity of being highly productive and
159 accounts for 65% of Tunisian fish production (DGPA, 2005–2009) and is a well-known habitat
160 for marine turtles (Jribi et al. 2008; Lotze and Worm, 2009; Maffucci et al. 2006). Although it
161 has huge importance for the local economy and wildlife conservation, the Gulf of Gabès has
162 been subject to the strong pressures of urbanization, industry, fisheries, tourism and
163 anthropogenic releases (Béjaoui et al. 2004; Ben Brahim et al. 2010; Rekik et al. 2012). It has
164 been the subject of increased anthropogenic interference threatening the entire ecosystem.

165 **2.2. Analysis process**

166 The flowchart in [Figure 1](#) was followed in order to conduct a consistent analysis. The
167 analysis ([Fig. 1](#)) includes three major steps: Data preparation, Discretization and Model
168 Prediction.

169 **2.2.1. Data preparation**

170 Data were collected in the framework of the National Phytoplankton Monitoring Program
171 in the shellfish harvest areas. This program has been operating since 1995 and covers fifteen
172 sites with weekly measures ([Fig. 2](#)). The monitoring was performed on a regular schedule all
173 year round. Each site can include more than one sampling station totalizing thirty-two sampling
174 stations ([Fig. 2](#)). The temporal windows of the dataset used for the analysis; that covered the
175 period ranging from 1997 to 2007; are considered as the most consistent in terms of data
176 shortage and because the identification of the species was carried out by the same team of
177 scientists.

178 Water for phytoplankton identification (1l) was sampled with a sampling water-bottle.
179 Temperature and salinity were measured for each water sample using a Handheld
180 Multiparameter Instrument: WTW Multi 340i/SET. Samples were fixed with lugol (4%)
181 solution and phytoplankton was enumerated by an inverted microscope using the Utermöhl's
182 method ([Utermöhl, 1958](#)).

183 Physical and meteorological parameters were simultaneously collected. Meteorological
184 data were provided by the Tunisian National Meteorological Institute (INM) and consisted of
185 air temperature, rainfall, evaporation, insolation, humidity and atmospheric pressure collected
186 from 15 meteorological stations. The tide amplitude, difference between high and ebb tides
187 considered as the variable indicating the tide effect, was obtained from the tide gauge stations
188 located in the main Tunisian ports and operated by the Tunisian Hydrographic and
189 Oceanographic Services.

190 The toxic dinoflagellate *Karenia selliformis* abundance was considered as the
191 investigated biological variable. The blooms generally define phytoplankton concentrations
192 having more than 10^6 cells L⁻¹ (Lassus, 1988). In this study, a concentration of 10^5 cells L⁻¹ was
193 considered as indicative of *K. selliformis* blooms since this concentration limit was responsible
194 for the red water discoloring in this area (Feki et al. 2008).

195 The dataset was filtered by removing the missing data i.e., the cases where some information
196 about the studied factors was unavailable. The final dataset contains over 1900 observations.

197 **2.2.2. Data discretization**

198 The models were applied to discrete variables. The choice of the class number and class
199 limits was already defined in a previous analysis (Feki-Sahnounet al. 2017). With regards to
200 *Karenia selliformis*, two discretization schemes were chosen. The first one used
201 presence/absence of the species to assess factors influencing the species occurrences whereas
202 the second one used bloom/non-bloom forming species to assess the factor affecting the species
203 abundance (Table 1).

204 **2.2.3. Model Prediction**

205 A set of linear and nonlinear models were used as prediction tools. For modeling the
206 effects of multiple factors (physical, meteorological and spatiotemporal parameters) on *K.*
207 *selliformis* occurrences and blooms, General linear model (GLM), Bayesian network (BN) and
208 naive Bayes classifier (NB) were developed.

209 Adjusted coefficient of determination (R^2), Akaike (AIC) (Akaike, 1974) and Bayesian
210 Information Criteria (BIC) (Schwartz, 1978) were used as model quality indicators to best
211 predict the phenomena. R^2 is based on the ratio between the variance explained by the model to
212 the total variance, while AIC and BIC are based on the likelihood of the data given the model.

213 GLM analysis was performed using the R package *MASS* (Venables and Ripley, 2002). BN and
214 NB models used algorithms implemented in the *bnlearn* package for R (R Development Core
215 Team, 2017; Scutari, 2010) and *SamIam* (Sensitivity Analysis, Modeling, Inference And More
216 software) (Chan and Darwiche, 2004) used to visualize the results.

217 **2.2.3.1. General linear model**

218 The univariate general linear models (GLM) allows categorical predictors to be included.
219 It was performed with stepwise procedures. The best-approximating GLM model with the
220 highest R^2 in combination with the lowest AIC obtained for each combination of studied factors
221 were selected as the final model (Burnham and Anderson, 2002). A final model included
222 marginally significant descriptors ($0.05 < P < 0.1$), if they explained deviance and/or AIC
223 notably improved with the descriptors in the model. GLM models were validated using
224 graphical tests implemented in R (plot function of a GLM object) to test residuals normality
225 and independence .

226 **2.2.3.2. Bayesian Networks**

227 Bayesian Networks (BNs) have two components: a qualitative one represented by a
228 Directed Acyclic Graph (DAG) depicting dependence relationships between variables and a
229 quantitative one represented by Conditional Probability Tables (CPTs).

230 **Figure 3A** shows a simple BN of five variables {U, V, W, X, Y} and four edges $U \rightarrow W$,
231 $V \rightarrow W$, $W \rightarrow Y$ and $X \rightarrow Y$. In a DAG, each node represents a variable of interest (i.e., U, V, W,
232 X and Y, **Fig. 3A**). Each variable is expressed according to its possible values continuous or
233 discrete (i.e., the values of X are x_1 and x_2). Nodes are connected to each other by edges
234 (directed arcs) to represent dependency relations between them (e.g., in **Fig. 3A**, X and W are
235 the parent of Y).

236 The CPT contain, for each possible value of the variable associated to a node, all the
237 conditional probabilities with respect to all the values' combinations of the variables associated
238 to the parent nodes. Generally, leaf nodes in a BN (nodes which have parents and that do not
239 have children) are nodes that are used to make a decision.

240 Concerning structure learning, the used score is the Bayesian Information Criterion
241 (BIC) (Schwartz, 1978) which is a penalized version of the log-likelihood. The search
242 algorithm; that is opted for; is the Tabu search (Glover, 1989), which was considered as a fast
243 and determinist algorithm.

244 For parameters' learning, a Maximum likelihood (ML) inference estimating the
245 probability of an event based on its occurrence frequency in the dataset (Neapolitan, 2003;
246 Pearl, 1988) was adopted.

247 2.2.3.3. Naive Bayes classifier

248 Naive Bayes classifier (NB) networks are the simplest model among BN. The term
249 'Naïve' refers to the strong independence assumption between variables (Friedman et al. 1997).
250 NB is graphically represented by a hierarchical structure where the class node is the parent of
251 all attribute nodes (Duda et al. 2001; Friedman et al. 1997). Therefore, the NB is
252 probabilistically defined by the conditional probabilities of each attribute given the class node
253 (Cooper and Herskovits, 1992). Each node; that is representing a given attribute; is associated
254 to Conditional Probability Table (CPT) containing, for each possible value of the corresponding
255 attribute, all the conditional probabilities with respect to the class nodes values.

256 Figure 3B shows an example of a class which is noted as "Y" and a set of four variables
257 {U, V, W, X}. The corresponding NB structure is composed of five nodes {U, V, W, X, Y}
258 and four edges {Y->U, Y->V, Y->W and Y->X}. Each variable is expressed according to its
259 possible values continuous or discrete (*i.e.*, the values of X are x1 and x2). The nodes are

260 connected to class "Y" by edges in order to represent the dependency relations to the class. Such
261 a structure is very useful for determining the impact of the attributes' values on the class' value.

262 The Maximum A Posteriori (MAP) estimation (Kramer and Sorenson, 1988) for
263 estimating the CPT of each node in the NB models was adopted.

264 The goal in this section is to determine the posterior probabilistic information. Thus, the
265 model can be used to predict the impact (in terms of probability) of introducing evidence for
266 certain variables (model a posteriori). For example, if the evidence (the observation that the
267 species is present), $P(\text{species}=\text{presence}) = 1$, is set in the class variable, the density functions
268 of the remaining environmental variables will be modified. In this way, an approximation to
269 the most probable configuration for the species presence can be obtained.

270 **3. Results**

271 **3.1. *Karenia selliformis* occurrences**

272 **3.1.1. General linear model**

273 The results of GLM analysis in Table 2 and from Equation (1) show that *Karenia*
274 *selliformis* occurrence depends mainly on Site, Month, Rain, Evaporation, Insolation and
275 Salinity.

276 High significant relationships were observed for sites belonging to the central and
277 southern part of the Gulf of Gabès namely G2, G3, M1, M2, M3, M6, S5 and S6 ($P < 0.1$). Late
278 summer, autumn and winter were generally the periods of high frequency in this area ($P < 0.1$).
279 High salinity and medium evaporation levels were responsible for Ks occurrences ($P < 0.1$)
280 (Table 2, Eq. 3).

281 Sampling sites (M6 and S5) and evaporation (medium level) show a negative relationship
282 with Ks occurrences. Whereas the relationships with G2, G3, M1, M2, M3, S6, from August to
283 December and salinity (high level) are positive (Eq. 2).

284 The stepwise GLM model selection allowed decreasing the AIC from 2158.81 to 2144.24.

285 The retained GLM model has the following structure:

$$286 \text{ Ks} = \text{Intercept} + \text{Site} + \text{Month} + \text{Rain} + \text{Evap} + \text{Insol} + \text{Sal} \quad (1)$$

$$287 = -1.84 + 1.16(\text{G2}) + 1.22(\text{G3}) + 0.76(\text{M1}) + 1.44(\text{M2}) + 0.87(\text{M3}) + (-1.22) (\text{M6}) + (-$$
$$288 0.56) (\text{S5}) + 0.5(\text{S6}) + 0.54(\text{Aug}) + 0.52(\text{Dec}) + 0.94(\text{Nov}) + 1.31(\text{Oct}) + 0.54(\text{Sep}) + (-0.36)$$
$$289 (\text{Evap b}) + 0.74(\text{Sal c}) \quad (2)$$

290 A numerical example of *Karenia selliformis* occurrences estimation in the M2 site and
291 during September is shown below:

$$292 \text{ Ks} = -1.84 + 1.16 (0) + 1.22 (0) + 0.76 (0) + 1.44 (1) + 0.87 (0) - 1.22 (0) - 0.56 (0) + 0.56 (0) +$$
$$293 0.54 (1) + 0.52 (0) + 0.94 (0) + 1.31 (0) + 0.54 (0) + -0.36 (\text{Evap b}) + 0.74 (\text{Sal c})$$
$$294 = -1.84 + 1.44 + 0.54 - 0.36 (\text{Evap b}) + 0.74 (\text{Sal c}) \quad (3)$$

295 From Eq. (3) and as shown in Table 2, the estimated values *Karenia selliformis*
296 occurrence is dependent upon the levels of evaporation and salinity.

297 3.1.2. Naive Bayes classifier

298 The resulting NB model is shown in Fig. 4. The introduction of the evidence “presence
299 of Ks” changes the probability distribution of the features because they are directly connected
300 with the Ks variable.

301 The marginal probability distributions of the NB model shows that it is likely to find the
302 Ks during October and July in the southern part of the Gabès’ Gulf both in M2 and M3 sites
303 (Fig. 4A and B).

304 The model also suggests that it is likely to find the species where there is a low rainfall
305 and slight high air temperature levels. The remaining variables (humidity, evaporation,
306 insolation, tide amplitude, salinity, water temperature and atmospheric pressure) vary in the
307 medium levels (b Class) (Fig. 4B).

308 **3.1.3. Bayesian networks**

309 The BN linking the physical and meteorological variables and the presence/absence of
310 the species shows (Fig. 5) that only the factor site influences the *Karenia selliformis* occurrence.
311 CPTs show that the posterior probability of the species presence was high in sites M2, M3, G3,
312 G2, G1, S6 and M1 (Figs. 5A and B), all of these sites belong to the central and southern part
313 of the Gulf of Gabès.

314 **3.1.4. Models comparison**

315 The main difference of BN with respect to the NB is the relationship between the
316 variables. In fact; this distinction increases the number of arcs in the structure and its
317 complexity, but also improves the accuracy and expressivity of the model.

318 BN had not detected the associations founded on GLM and NB. Hence, while GLM and
319 NB established a direct association between three and six variables, BN identified only site-
320 dependent factor linked to the outcome.

321 M2 and M3 show posterior probability similar to the NB model (27.29 % and 10.51 %
322 respectively) (Fig. 5B and 4B). In general, both NB and BN models show similar probabilities
323 trends (Fig. 5 and 4). However, BN varies in the definition of relationships between the
324 variables in the model, so that each variable is influenced, not only by the class variable Ks, but
325 also by the variables directly connected with it in the network.

326 The prediction performance of the GLM model is moderate having adjusted R^2 value of
327 0.17 (Table 4). The NB and BN models appeared to have the best fit and performed almost
328 equally well having an adjusted R^2 of 0.69 and 0.70 respectively, while BN was found to
329 perform slightly better (Table 4).

330 **3.2. *Karenia selliformis* blooms**

331 3.2.1. General linear model

332 GLM shows that *Karenia selliformis* blooms mainly depends on Site, Month,
333 Evaporation, Humidity, Water Temperature and Salinity (Table 3; Eq. (4)).

334 Humidity (Class b and c) shows a negative relationship with *Karenia selliformis* blooms,
335 whereas evaporation (Class c) and water temperature (Class c) show positive correlations
336 (Table 3, Eq. 5).

337 The GLM model selection, allows decreasing the AIC from 588.69 to 565.76.

338 The selected GLM model has the following structure:

$$339 K_s \sim \text{Intercept} + \text{Site} + \text{Month} + \text{Evap} + \text{Humid} + \text{WatT} + \text{Sal} \quad (4)$$

$$340 K_s = 22.54 + 1.59 (\text{Evap } c) - 2.12 (\text{Humid } b) - 1.83 (\text{Humid } c) + 1.57 (\text{WatT } c) \quad (5)$$

341 A numerical example of *Karenia selliformis* blooms estimation in the M2 site is shown
342 below:

$$343 K_s = 22.54 + 1.59 (\text{Evap } c) - 2.12 (\text{Humid } b) - 1.83 (\text{Humid } c) + 1.57 (\text{WatT } c) \quad (6)$$

344 From Eq. (6) and as shown in Table 3, the estimated value *Karenia selliformis* blooms
345 are dependent upon the levels of Evaporation, Water temperature and Humidity.

346 3.2.2. Naive Bayes classifier

347 The model suggests that it is likely to find blooms not only when there is a high salinity
348 and air temperature, but also when there are remaining variables that are situated in medium
349 level (b class) (Fig. 6B). It can also be seen regarding bloom a slightly high conditional
350 probability (38%) for low tide amplitude level (a Class). As K_s occurrences, NB model also
351 shows that it is likely to find the bloom in M2 (85.29%) between the two-temporal scenario of
352 summer and autumn (Fig. 6B).

353 3.2.3. Bayesian networks

354 The BN schema illustrating the physical, meteorological variables and the bloom/non-
355 bloom of the *K. selliformis* (Fig. 7) shows that salinity directly affects Ks blooms and represents
356 its single parent (Fig. 7). Salinity has site as parent, and site is dependent on tide, atmospheric
357 pressure, air temperature and month. CPTs show, that the probabilities of low salinity (a Class)
358 are very low and are not exceeding 6% (Fig. 7 A and B). It is worth noting that the absence of
359 blooms is high (70%) for medium salinity levels (b Class) particularly in M2 site (25%). (Fig.
360 7A). During Ks blooms, the probability to have high salinity (Class c) significantly increased
361 (80%). The highest probabilities are recorded in M2 (56%) (Fig. 7B). These results show that
362 the blooms were largely determined by salinity and were expected to occur in salinities higher
363 than 42.5, which are detected in site M2 (Fig. 7B).

364 **3.2.4. Models comparison**

365 BN do not detect the associations found on GLM and NB; except for salinity (Fig. 7).
366 Based on the prediction performance indicator, all models perform moderately with the better
367 performance for BN model ($R^2=0.58$) followed by the NB ($R^2=0.54$) and GLM models
368 ($R^2=0.48$) (Table 4).

369 **4. Discussion**

370 **4.1. Models comparison**

371 In computing the tree models, several strengths and constrains relative to one or other
372 models were raised. They are summarized in [Table 5](#). One of these constraints was the exclusive
373 use of discrete variables which was considered as an important weakness of BNs ([Landuyt et](#)
374 [al. 2013](#)). Most environmental variables are continuous whilst GLM can use continues or
375 discrete data, Bayesian networks usually build the model over discrete domains, so that
376 continuous variables need to be first discretized ([Uusitalo, 2007](#)). Discretization implies
377 capturing only rough characteristics of the original distribution ([Friedman and](#)

378 [Goldszmidt,1996](#)) and loss of statistical information ([Aguilera et al. 2011, 2010](#); [Alameddine](#)
379 [et al. 2011](#); [Chen and Pollino, 2012](#); [Hamilton et al. 2015](#); [Jensen, 2001](#); [Lucena-Moya et al.](#)
380 [2015](#); [Meineri et al. 2015](#); [Meyer et al. 2014](#); [Nielsen and Jensen, 2007](#); [Uusitalo, 2007](#)).
381 However, the discretization is beneficial in this study and the BN has succeeded in providing
382 reliable results and generalized BNs which do not provide biased results. Thus, the
383 discretization can mitigate the impact of the noisy data (incorrectly saved values / typos /
384 wrong values).

385 The models also confirmed the role of salinity as a key factor for Ks bloom ([Feki et al.](#)
386 [2008](#); [Feki-Sahnoun et al.2017](#)). The results of GLM analysis in [Table 2](#) and from Equation (1)
387 refutes the exclusion of salinity and the retention of water temperature from the generalized
388 linear mixed-effect model (GLMM) established for the *K. selliformis* occurrences in the Gulf
389 of Gabès ([Feki et al. 2013](#)). Therefore, BN, NB and GLM are likely to identify the same factors
390 when associations are strong and highly significant. The advantage of BN, comparing to the
391 two others models, is manifested in its ability to easily establish a direct association between
392 bloom and salinity and between occurrence and sampling site. Indeed, BN identifies a network
393 of inter-dependent factors linked to the outcome. It is more informative about potential causal
394 pathways in the Ks occurrences and blooms than GLM and NB ([Table 5](#)). This potential
395 advantage of the BN method would specifically be useful for observational studies with large
396 number of variables, where causal and time relationships are often unknown ([Pittavino et al.](#)
397 [2017](#)). Moreover, in the retained GLM model ([Table 2](#)) some variables are interdependent such
398 as insolation and evaporation both were related to salinity. Hence, the interdependencies
399 between variables were revealed in BN that might not be discovered in GLM and NB, as the
400 latter impose a linear relationship between covariates and the outcome. Thus, some hidden
401 dynamic was still not detected while using GLM and NB when interaction terms were
402 considered ([Table 5](#)). For instance; the relationships between tide, site and salinity ([Fig. 7](#)) are

403 explained by the fact that the species achieved its peak density in a high salinity and in a semi-
404 enclosed lagoon (M2), compared to the open coastline, incriminating the low water dilution
405 rate. In this case low tide amplitude, due to the weak water advection out of the lagoon, allows
406 to sustain the growth and the bloom's maintenance of the species (Feki-Sahnoun et al. 2017).
407 Hence, BN had an advantage over GLM and NB because it can disentangle the complex nature
408 of the data, stratifying further the presented internal mechanisms between the variables
409 (Pittavino et al. 2017).

410 The BN and NB performed almost similarly for both occurrences and blooms (Table 4).
411 The two models also yielded very close range of probabilities regarding all influencing
412 variables in the BN of Ks occurrences such as sampling site identified as the only parent (e.g.
413 given the evidence of presence: M2=27%) (Figs.4B and 5B) and salinity identified as the only
414 parent in the BN regarding Ks bloom (given the evidence of presence: c=79%) (Figs.6B and
415 7B). This could be in part explained by the reduced complexity of the studied dataset as shown
416 by Aguilera et al. (2010). Moreover, even performance indicators were rather close between the
417 two models, it was slightly higher for BN (Table 4). This could be explained by the
418 interdependency between variables considered in NB which might put more confidence on the
419 BN predicted probabilities (Boets et al. 2015). Indeed, BN models may better reflect the joint
420 probability distribution over the system's variables compared to a Naive Bayes classifier, as the
421 assumption of conditional independence among predictor variables, may not hold in reality
422 (Boets et al. 2015). In others words, while NB may predict occurrences or blooms better (Boets
423 et al. 2015), BN models may predict occurrence or blooms probabilities more accurately. This
424 may also explain the slightly weaker performance of the NB model compared to the BN model
425 to predict both presence/absence and bloom/non-bloom (Table 4). Although the number of arcs
426 in the structure and its complexity increases, no significant effect on model performance was
427 found and the BN improves usually the accuracy, flexibility and expressivity (Pearl, 1988).

428 Nevertheless, it is expected that; by increasing the number of states beyond the relevant range
429 (more than three) can result in significant performance loss. Increasing the number of states
430 would result in large CPTs wherein the probabilities that need to be estimated are conditional
431 on very specific combinations of environmental conditions. This increases the chance of having
432 no or only a limited number of data records available to learn these conditional probabilities.
433 This would result in a lower predictive performance on an independent test dataset (Boets et al.
434 2015).

435 More importantly; the NB model also shows that the introduction of the evidence
436 “presence of Ks” changes the probability distribution of several features mainly salinity, site,
437 month, air temperature, tide amplitude compared to the evidence of the Ks being absent (Fig.
438 4). The same figure was observed on the introduction of evidence “bloom of Ks” compared to
439 the evidence of the Ks being non -blooming (Fig. 6). However, the BN schemes demonstrated
440 that introduction of the evidence “presence of Ks” or “bloom of Ks” do not change the
441 probabilities except that for the parent of Ks (sampling site and salinity respectively) (Figs. 5
442 and 7). These outcomes arise from the facts that the NB may overamplify the weight of the
443 evidence of each attribute on the class (Friedman et al. 1997).

444 The probability of salinity exceeding 42.5 (c class) was very high (above 70%) in *K.*
445 *selliformis* blooms (Figs. 6 and 7). This result was explained by a non-linear behavior between
446 salinity and species abundances, and it suggests that the transition to high biomasses appeared
447 to be triggered by a salinity threshold (Feki-Sahnounet al. 2017). This circumstance, might
448 advantage the use of BN as a powerful computational technique for modeling complex
449 relationships in situations; where the proper form of the relation between the variables is
450 unknown or nonlinear (Chen and Billings, 1992; Felipe et al. 2014; Schmitt and Brugere, 2013).
451 The conventional techniques such as GLM and NB, although they are widespread within
452 ecology, but they present some short comings related to the facts that the relationships between

453 variables in environmental sciences are often non-linear and that data rarely have normal errors
454 (Chen and Billings, 1992).

455 The overall conclusion is that BN are the recommended models since it has less
456 constraints than the others (Table 5) and provide the highest statistical performance (Table 4).

457 **4.2. Ecological implications and management recommendations**

458 The three models incriminated the variable site in Ks occurrence confirming the north-
459 south gradient of Ks occurrence already pointed out in the Gulf (Feki et al. 2013). This is due
460 to the influence of the site M2 (the Boughrara lagoon) which is totalizing to itself more than
461 27% of the species occurrences (Figs. 4B and 5B) and more than 85% of the species blooming
462 (Fig. 6B). One of the direct implication of this result is the consideration of the site M2 as a
463 "hot spot " for the species proliferation which might imply paying more attention to the
464 variability of the physical and meteorological variables identified in the bloom forming in this
465 specific location.

466 As it happens salinity appears as the parameter that effect directly Ks blooms.
467 Nevertheless, GLM retained water temperature as an explicative factor in Ks Blooms (Table 3,
468 Eq (4)) and NB showed an increase in the associated probabilities of the blooms regarding the
469 highest temperature level (Fig. 6).Whereas no link was pointed out by BN between Ks bloom
470 and water temperature (Fig. 7). The correlations between high water temperature level and Ks
471 blooms were often highlighted in several ecosystems (Carreto et al. 2001; Clement et al. 2001;
472 Uribe and Ruiz, 2001) and in the Gulf of Gabès (Dammak-Zouari et al. 2009) stressing the role
473 of this parameter in the bloom enhancement. In the BN's variables interconnections, water
474 temperature itself is affected by air temperature the one of Ks's parents (Fig. 7).This suggests
475 that temperature on itself might not be the factor affecting Ks but it might rather impact on
476 salinity which in turn influenced Ks. The mechanism by which salinity could affect the bloom

477 enhancement is not yet well identified: is this a purely physiological effect or could involve
478 also physical conditions is still to be investigated in order to statute whether the apparent direct
479 effect between Ks blooms and salinity was rather an association than a causality.

480 The identification of a salinity threshold, revealed by BN, is a relevant information for
481 the prediction of the blooms in the studied ecosystem, and it can be used to design and set up
482 an early warning system for Ks blooms based on a real time observation of salinity. One of the
483 direct application of such a system is the suspension of shellfish exploitation during the
484 suspected period of high salinity; pending the species abundance determination and the
485 associated toxicity tests performed.

486 **Conclusions**

487 In terms of predictive ability, BN performed better than linear models (NB and GLM)
488 regarding Ks occurrences and blooms prediction, probably due to the existence of nonlinear
489 relationships with the salinity key variable.

490 BN and NB performed quasi equally in terms of performance indicator. BN has an
491 advantage over NB because of its capacity to capture and illustrate graphically the data's natural
492 complexity more effectively. In BN, all relationships between variables are modeled, which
493 appears to be more explanatory in the view of the inter-dependencies between variables studies.
494 BN can work together with NB for pre-selection of variables inputs.

495 The three investigate models converge on the identification of salinity as a key variable
496 for the prediction of Ks occurrences and blooms. The salinity is in turn site dependent with
497 more that 55% of the bloom concentrated in Boughrara lagoon which suggests using this
498 parameter for the control of the Ks bloom in this specific location identified as a hot spot area.

499 **Acknowledgments**

500 This work was supported by the PASRI (L'Agence Nationale de Promotion de la Recherche
501 scientifique)/MOBIDOC (Mobilisation de docteur pour la réalisation de Travaux de Recherche
502 dans l'Entreprise) funded Project (post-doctoral grant for the first author). The authors wish to
503 thank Mr. Hamdi DKHIL, English Teacher in Franklin Center Sfax (Tunisia) for having edited
504 this Paper.

505 **References**

- 506 Aguilera, P.A., Fernandez, A., Fernandez, R., Rumi, R., Salmeron, A., 2011. Bayesian networks
507 in environmental modelling. *Environ. Model. Softw.* 26, 1376–1388.
- 508 Aguilera, P.A., Fernandez, A., Reche, F., Rumi, R., 2010. Hybrid Bayesian network classifiers:
509 application to species distribution models. *Environ. Model. Softw.* 25, 1630–1639.
- 510 Aguilera, P.A., Fernandez, A., Roperro, R.F., Molina, L., 2013. Groundwater quality assessment
511 using data clustering based on hybrid Bayesian networks. *Stoch. Environ. Res. Risk Assess.* 27,
512 435–447.
- 513 Akaike, H., 1974. A new look at the statistical identification model. *IEEE Trans. Auto. Control.*
514 19, 716-723.
- 515 Alameddine, I., Cha, Y.K., Reckhow, K.H., 2011. An evaluation of automated structure
516 learning with Bayesian networks: An application to estuarine chlorophyll dynamics. *Environ.*
517 *Model. Softw.* 26, 163–172.
- 518 Béjaoui, B., Rais, S., Koutitonsky, V., 2004. Modélisation de la dispersion du phosphogypse
519 dans le Golfe de Gabès. *Bull. Inst. Natl. Sci. Technol. Mer de Salammbô* 31, 103–109.
- 520 Ben Brahim, M., Hamza, A., Hannachi, I., Rebai, A., Jarboui, O., Bouain, A., Aleya, L., 2010.
521 Variability in the structure of epiphytic assemblages of *Posidonia oceanica* in relation to human
522 interferences in the Gulf of Gabès, Tunisia. *Mar. Environ. Res.* 70, 411–421.

523 Ben Naila, I., Hamza, A., Gdoura, R., Diogene, J., Iglesia, P., 2012. Prevalence and persistence
524 of gymnodimines in clams from the Gulf of Gabès (Tunisia) studied by mouse bioassay and
525 LC-MS/MS. *Harmful Algae*. 18, 56–64.

526 Bishop, C.M., 2006. *Pattern Recognition and Machine learning*, Springer, New York.

527 Boets, P., Landuyt, D., Everaert, G., Broekx, S., Goethals, P.L.M., 2015. Evaluation and
528 comparison of data-driven and knowledge-supported Bayesian Belief Networks to assess the
529 habitat suitability for alien Macroinvertebrates. *Environ. Model. Softw.* 74, 92–103.

530 Borsuk, M., Reichert, P., Peter, A., Schager, E., Burkhardt-Holm, P., 2006. Assessing the
531 decline of brown trout (*Salmo trutta*) in swiss rivers using Bayesian probability network. *Ecol.*
532 *Model.* 192, 224–244.

533 Borsuk, M., Stow, C., Reckhow, K., 2004. A Bayesian network of eutrophication models for
534 synthesis, prediction, and uncertainty analysis. *Ecol. Model.* 173, 219–239.

535 Bromley, J., Jackson, N., Clymer, O., Giacomello, A., Jensen, F., 2005. The use of Hugin to
536 develop Bayesian networks as aid to integrated water resource planning. *Environ. Model.*
537 *Softw.* 20, 231–242.

538 Burnham, K.P., Anderson D.R., 2002. *Model Selection and Multimodel Inference: a Practical*
539 *Information-Theoretic Approach*, 2nd ed. Springer Verlag, New York.

540 Cain, J.D., Jinapala, K., Makin, I.W., Somaratna, P.G., Ariyaratna, B.R., Perera, L.R., 2003.
541 Participatory decision support for agricultural management. A case study from Sri Lanka.
542 *Agric. Syst.* 76, 457–482.

543 Carreto, J.L., Seguel, M., Montoya, N.G., Clément, A., Carignan, M.O., 2001. Pigment profile
544 of the ichthyotoxic dinoflagellate *Gymnodinium* sp. from a massive bloom in southern Chile. *J.*
545 *Plankton Res.* 23, 1171–1175.

546 Chan, H., Darwiche, A., 2004. Sensitivity analysis in Bayesian networks: From single to
547 multiple parameters, in: Chickering, M., Halpern, J. (Eds.), Proceedings of the 20th conference
548 on Uncertainty in Artificial Intelligence. AUAI Press, Arlington, VA, pp. 67–75.

549 Chan, T.U., Hart, B.T., Kennard, M.J., Pusey, B.J., Shenton, W., Douglas, M.M., Valentine, E.,
550 Patel, S., 2012. Bayesian network models for environmental flow decision making in the Daly
551 river, northern Territory, Australia. *River Res. Appl.* 28, 283–301.

552 Chen, S., Billings, S.A., 1992. Neural networks for nonlinear dynamic system modelling and
553 identification. *Internat. J. Control.* 56, 319–346.

554 Chen, S.H., Pollino, C.A., 2012. Good practice in Bayesian network modelling. *Environ.*
555 *Model. Softw.* 37, 134–145.

556 Clement, A., Miriam, S., Arzul, G., Guzman, L., Alarcon, C., 2001. Widespread outbreak of a
557 haemolytic, ichthyotoxic *Gymnodinium* sp. in southern Chile, in: Hallegraeff, G.M., Blackburn,
558 S.I., Bolch, C.J., Lewis, R.J. (Eds.), Harmful Algal Blooms 2000. IOC of UNESCO, Paris, pp.
559 66–69.

560 Cooper, G.F., Herskovits, E., 1992. A Bayesian method for the induction of probabilistic
561 networks from data. *Mach. Learn.* 9, 309–347.

562 Dammak-Zouari, H., Hamza, A., Bouain, A., 2009. Gymnodiniales in the Gulf of Gabès
563 (Tunisia). *Cah. Biol. Mar.* 50, 153–170.

564 DGPA., 2005–2009. Direction Générale de la pêche et de l’aquaculture. Ministère de
565 l’agriculture, Tunisie, annuaire statistique.

566 Domingos, P., Pazzani, M., 1997. On the optimality of the simple bayesian classifier under
567 zero-one loss. *Mach. Learn.* 29, 103–130.

568 Duda, R.O., Hart, P.E., Stork, D.G., 2001. Pattern Classification, second ed. Wiley Interscience,
569 New York.

570 Feki-Sahnoun, W., Hamza, A., Njah, H., Barraï, N., Mahfoudi, M., Rebai, A., Bel Hassen, M.,
571 A., 2017. Bayesian network approach to determine environmental factors controlling *Karenia*
572 *selliformis* occurrences and blooms in the Gulf of Gabès, Tunisia. *Harmful Algae*. 63, 119–132.

573 Feki, W., Hamza, A., Bel Hassen, M., Rebai, A., 2008. Les efflorescences phytoplanctoniques
574 dans le golfe de Gabès (Tunisie) au cours de dix ans de surveillance (1995–2005). *Bull. Inst.*
575 *Natl. Sci. Tech. Oceanogr. Peche Salammbô*. 35, 105–116.

576 Feki, W., Hamza, A., Frossard, V., Abdennadher, M., Hannachi, I., Jacquot, M., Bel Hassen,
577 M., Aleya, L., 2013. What are the potential drivers of blooms of the toxic dinoflagellate *Karenia*
578 *selliformis*? A 10-year study in the Gulf of Gabès, Tunisia, southwestern Mediterranean Sea.
579 *Harmful Algae*. 23, 8–18.

580 Felipe, V.P., Okut, H., Gianola, D., Silva, M.A., Rosa, G.J., 2014. Effect of genotype
581 imputation on genome-enabled prediction of complex traits: an empirical study with mice data.
582 *BMC Genet*. 15, 149.

583 Fernandes, J.A., Irigoien, X., Goikoetxea, N., Lozano, J.A., Inza, I., Pérez, A., Bode, A., 2010.
584 Fish recruitment prediction, using robust supervised classification methods. *Ecol. Model*. 221,
585 338–352.

586

587 Friedman, N., Geiger, D., Goldszmidt, M., 1997. Bayesian network classifiers. *Mach. Learn.*
588 29, 131–163.

589 Friedman, N., Goldszmidt, M., 1996. Discretization of continuous attributes while learning
590 Bayesian networks, in: Saitta, L. (Ed.), *Proceedings of the Thirteenth International Conference*
591 *on Machine Learning*. CA: Morgan Kaufmann, San Francisco, pp. 157–165.

592 Fytilis, N., Rizzo, D.M., 2013. Coupling self-organizing maps with a Naïve Bayesian classifier:
593 stream classification studies using multiple assessment data. *Water Resour. Res.* 49, 7747–
594 7762.

595 Glover, F., 1989. Tabu Search—Part I. *ORSA J. Comput.* 1, 190–206.

596 Gonzalez-Redin, J., Luque, S., Poggio, L., Smith, R., Gimona, A., 2016. Spatial Bayesian belief
597 networks as a planning decision tool for mapping ecosystem services trade-offs on forested
598 landscapes. *Environ. Res.* 144,15–26.

599 Haapasaari, P., Karjalainen, T.P., 2010. Formalizing expert knowledge to compare alternative
600 management plans: sociological perspective to the future management of Baltic salmon stocks.
601 *Mar. Policy* 34, 477–486.

602 Haines-Young, R., 2011. Exploring ecosystem service issues across diverse knowledge
603 domains using Bayesian Belief Networks. *Prog. Phys. Geogr.* 35, 681–699.

604 Hamilton, S.H., Pollino, C.A., Jakeman, A.J., 2015. Habitat suitability modelling of rare species
605 using Bayesian networks: model evaluation under limited data. *Ecol. Model.* 299, 64–78.

606 Hamza, A., El Abed, A., 1994. Les eaux colorées dans le golfe de Gabès: Bilan de six ans de
607 surveillance (1989–1994). *Bull. Inst. Natl. Sci. Tech. Oceanogr. Peche Salamambo.* 21, 66–72.

608 Hansen, G., Erard-Le Denn, E., Daugbjerg, N., Rodriguez, F., 2004. *Karenia selliformis*
609 responsible for the fish-kills in the gulf of Gabès, Tunisia 1994. Harmful algal Blooms. Program
610 and Abstracts of the 11th International Conference: Cape Town. Communication Ifremer 2004.

611 Hastie, T., Tibshirani, R., Friedman, J., 2009. *The elements of statistical learning*, Springer,
612 New York.

613 Hastie, T.J., Tibshirani, R.J., 1990. *Generalized Additive Models*. Chapman & Hall/CRC.

614 Jensen, F.V., 1996. *An introduction to Bayesian Networks*, Springer Verlag, New York.

615 Jensen, F.V., 2001. *Bayesian Networks and Decision Graphs*, Springer-Verlag, New York.

616 Jribi, I., Echwikhi, K., Bradai, M.N., Bouain, A., 2008. Incidental capture of sea turtles by
617 longlines in the Gulf of Gabès (South Tunisia): a comparative study between bottom and surface
618 longlines. *Sci. Mar.* 72, 337–342.

619 Kramer, S., Sorenson, H.W., 1988. Recursive Bayesian estimation using piece-wise constant
620 approximations. *Automatica*. 24, 789–801.

621 Landuyt, D., Broekx, S., D'Hondt, R., Engelen, G., Aertsens, J., Goethals, P.L.M., 2013. A
622 review of Bayesian belief networks in ecosystem service modelling. *Environ. Model. Softw.*
623 46, 1–11.

624 Lassus P., 1988. Plancton toxique et plancton d'eaux rouges sur les côtes européennes, Ed.
625 IFREMER, Paris.

626 Lauritzen, S.L., Spiegelhalter, D.J., 1988. Local computations with probabilities on graphical
627 structures and their applications to expert systems. *J. R. Stat. Soc.* 50, 157–224.

628 Lotze, H.K., Worm, B., 2009. Historical baselines for large marine animals. *Trends Ecol. Evol.*
629 24, 254–262.

630 Lucena-Moya, P., Brawata, R., Kath, J., Harrison, E., ElSawah, S., Dyer, F., 2015.
631 Discretization of continuous predictor variables in Bayesian networks: an networks is NP-hard.
632 *Artif. Intell.* 60, 141–153.

633 Madsen, H., Thyregod, P., 2011. Introduction to General and Generalized Linear Models.
634 Chapman & Hall/CRC.

635 Maffucci, F., Kooistra, W.H.C.F., Bentivegna, F., 2006. Natal origin of loggerhead turtles,
636 *Caretta caretta*, in the neritic habitat off the Italian coasts, Central Mediterranean. *Biol.*
637 *Conserv.* 127, 183–189.

638 Markus, M., Hejazi, M.I., Bajcsy, P., Giustolisi, O., Savic, D.A., 2010. Prediction of weekly
639 nitrate-N fluctuations in a small agricultural watershed in Illinois. *J. Hydroinformatics.* 12, 251–
640 261.

641 Marrouchi, R., Benoit, E., Kharrat, R., Molgó, J., 2009. Gymnodimines: a family of
642 phycotoxins contaminating shellfish, in: Benoit, E., Goudey-Perrière, F., Marchot, P., Servent,

643 D. (Eds.), *Toxins and Signalling*. Châtenay-Malabry, SFET Editions, Collection Rencontres en
644 Toxinologie, France, pp. 79–83.

645 McCulloch, C.E., Searle, S.R., Neuhaus, J.M., 2008. *Generalized, Linear, and Mixed Models*,
646 2nd edition, Hoboken, New Jersey: John Wiley & Sons.

647 McCullagh, P., Nelder, J., 1989. *Generalized Linear Models*, Second Edition, Boca Raton:
648 Chapman and Hall/CRC.

649 McVittie, A., Norton, L., Martin-Ortega, J., Siameti, I., Glenk, K., Aalders, I., 2015.
650 Operationalizing an ecosystem services-based approach using Bayesian Belief Networks: an
651 application to riparian buffer strips. *Ecol. Econ.* 110, 15–27.

652 Medhioub, A., Medhioub, W., Amzil, Z., Sibat, M., Bardouil, M., Ben Neila, I., Mezghani, S.,
653 Hamza, A., Lassus, P., 2009. Influence on environmental parameters on *Karenia selliformis*
654 toxin content in culture. *Cah. Biol. Mar.* 50, 333–342.

655 Medhioub, W., Guéguen, M., Lassus, P., Bardouil, M., Truquet, P., Sibat, M., Medhioub, N.,
656 Soudant, P., Kraiem, M., Amzil, Z., 2010. Detoxification enhancement in the gymnodimine
657 contaminated grooved carpet shell, *Ruditapes decussatus* (Linné). *Harmful Algae.* 9, 200–207.

658 Meineri, E., Dahlberg, C.J., Hylander, K., 2015. Using Gaussian Bayesian Networks to
659 disentangle direct and indirect associations between landscape physiography, environmental
660 variables and species distribution. *Ecol. Model.* 313, 127–136.

661 Meyer, S.R., Johnson, M.L., Lillholm, R.J., Cronan, C.S., 2014. Development of a
662 stakeholder-driven spatial modeling framework for strategic landscape planning using Bayesian
663 networks across two urban–rural gradients in Maine, USA. *Ecol. Model.* 291, 42–57.

664 Neapolitan, R., 2003. *Learning Bayesian networks*, Prentice Hall, New York.

665 Nelder, J.A., Wedderburn, R.W.M., 1972. *Generalized Linear Models*. *J. R. Stat. Soc. A*, 135,
666 370–384.

667 Nielsen, T.D., Jensen, F.V., 2007. Bayesian Networks and Decision Graphs, second ed.
668 Springer, New York.

669 Pearl, J., 1985. A model of self-activated memory for evidential reasoning In Proceedings of
670 the 7th Conference of the Cognitive Science Society, University of California, Irvine, CA:
671 University of California, pp. 329–334.

672 Pearl, J., 1988. Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference,
673 first ed. Morgan Kaufmann, San Mateo, CA.

674 Pittavino, M., Dreyfus, A., Heuer, C., Benschop, J., Wilson, P., Collins-Emerson, J., Torgerson,
675 P.R., Furrer, R., 2017. Comparison between Generalized Linear Modelling and Additive
676 Bayesian Network; Identification of Factors associated with the Incidence of Antibodies against
677 *Leptospira interrogans* sv *Pomona* in Meat Workers in New Zealand. Acta Trop. 173, 191–
678 199.

679 Pollino, C., White, A., Hart, B., 2007. Examination of conflicts and improved strategies for the
680 management of an endangered Eucalypt species using Bayesian networks. Ecol. Model. 201,
681 37–59.

682 R Development Core Team, 2017. R: A Language and Environment for Statistical Computing.
683 Vienna, Austria: the R Foundation for Statistical Computing. ISBN 3-900051-07-0.
684 <http://www.R-project.org/>(accessed 25.07.17).

685 Rekik, A., Drira, Z., Guerhazi, W., Elloumi, J., Maalej, S., Aleya, L., Ayadi, H., 2012. Impacts
686 of an uncontrolled phosphogypsum dumpsite on summer distribution of phytoplankton,
687 copepods and ciliates in relation to abiotic variables along the near-shore of the southwestern
688 Mediterranean coast. Mar. Pollut. Bull. 64, 336–346.

689 Ropero, R.F., Aguilera, P.A., Fernandez, A., Rumí, R., 2014. Regression using hybrid Bayesian
690 networks: modelling landscape-socioeconomy relationships. Environ. Model. Softw. 54, 127–
691 137.

692 Ropero, R.F., Aguilera, P.A., Rumí, R., 2015. Analysis of the socioecological structure and
693 dynamics of the territory using a hybrid Bayesian network classifier. *Ecol. Model.* 311, 73–87.

694 Sammari, C., Koutitonsky, V.G., Moussa, M., 2006. Sea level variability and tidal resonance in
695 the Gulf of Gabès, Tunisia. *Cont. Shelf. Res.* 26, 338–350.

696 Schmitt, LHM., Brugere, C., 2013. Capturing Ecosystem Services, Stakeholders' Preferences
697 and Trade-Offs in Coastal Aquaculture Decisions: A Bayesian Belief Network Application.
698 *PLoS ONE.* 8, e75956.

699 Schwarz, G., 1978. Estimating the dimension of a model, *Ann. Stat.* 6, 461–464.

700 Scutari, M., 2010. Learning Bayesian Networks with the bnlearn R Package. *J. Stat. Softw.*
701 35,1–22.

702 Smith, C., Howes, A., Price, B., McAlpine, C., 2007. Using Bayesian belief network to predict
703 suitable habitat of an endangered mammal-the Julia Creek dunnart (*Sminthopsis douglasi*).
704 *Biol. Conserv.* 139, 333–347.

705 Uribe, J.C., Ruiz, M., 2001. *Gymnodinium* brown tide in the Magellanic fjords, southern Chile.
706 *Rev. Biol. Mar. Oceanogr.* 36, 155–164.

707 Utermöhl, H., 1958. Zur vervollkommnung der quantitativen phytoplankton-methodik. *Mitt.*
708 *Int. Ver. Theor. Angew. Limnol.* 9, 1–38.

709 Uusitalo, L., 2007. Advantages and challenges of Bayesian networks in environmental
710 modelling. *Ecol. Model.* 203, 312–318.

711 Venables, W.N., Ripley, B.D., 2002. Random and Mixed Effects, in: *Modern Applied Statistics*
712 *with S. Statistics and Computing.* Springer, New York, NY.

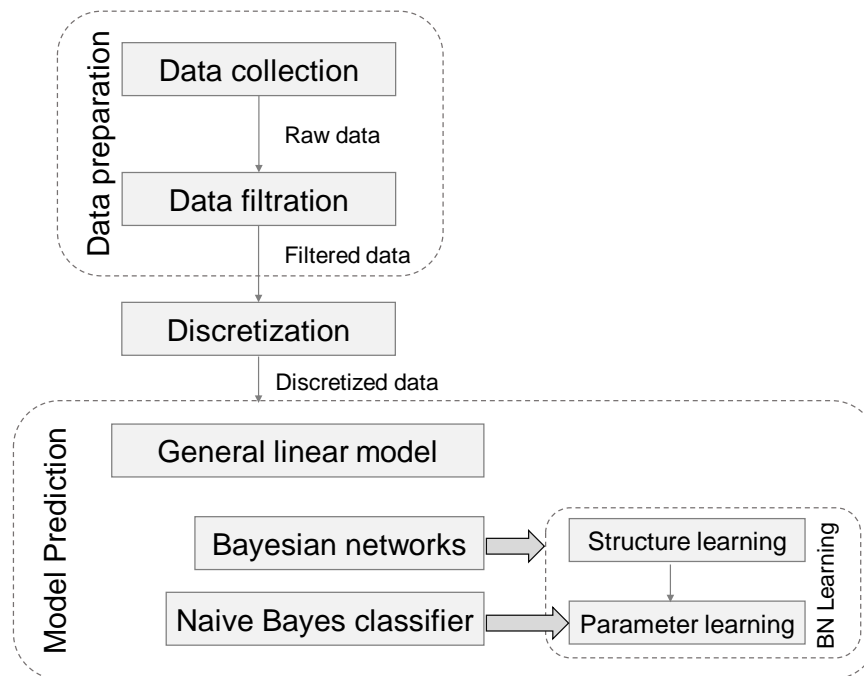
713 Wang, Q.J., Robertson, D.E., Haines, C.L., 2009. A Bayesian network approach to knowledge
714 integration and representation of farm irrigation: 1. Model development: knowledge integration
715 of farm irrigation, 1. *Water Resour. Res.* 45, W02409.

716 Zaffalon, M., 2005. Credible classification for environmental problems. *Environ. Model.*
717 *Softw.* 20, 1003–1012.

718 Zhang, H., 2004. The optimality of naive Bayes, in: *Proceedings of the 17th International Florida*
719 *Artificial Intelligence Research Society Conference*. AAAI Press, Florida, USA.

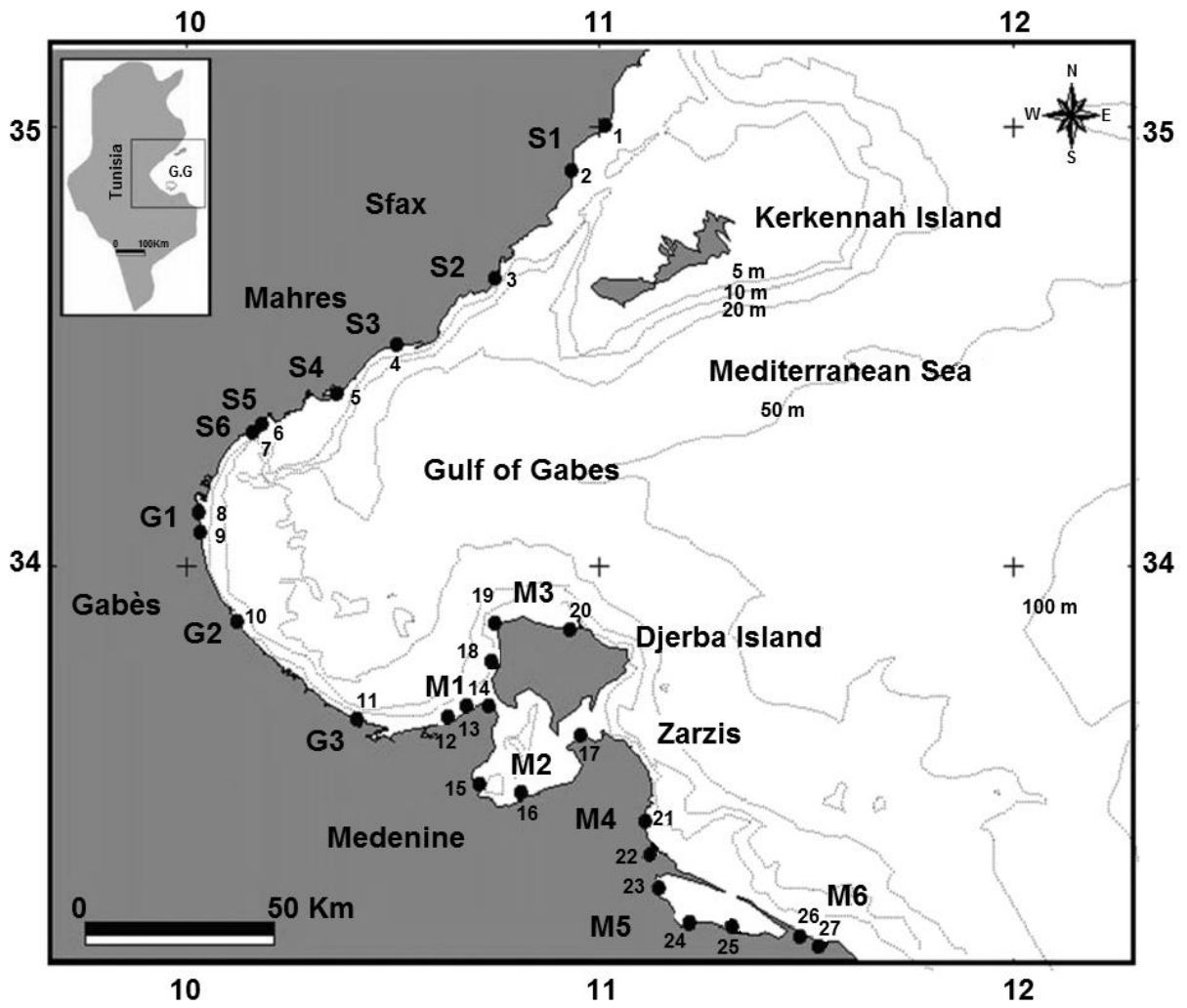
720 **Figures**

721 **Fig.1.** Framework of the proposed approach describing the three major steps: Data preparation,
722 Discretization and Model Prediction.



723

724 **Fig. 2.** Geographical map focusing on the monitoring network of phytoplanktonic sampling
725 stations in the Gulf of Gabès, Tunisia: 15 sampling sites and 28 sampling locations from 1997
726 to 2007. S1 (1-2), S2 (3), S3 (4), S4 (5), S5 (6), S6 (7), G1 (8-9), G2 (10), G3 (11), M1 (12-
727 13), M2 (14-15-16-17), M3 (18-19-20), M4 (21-22), M5 (23-24-25), M6 (26-27).



728

729

730

731

732

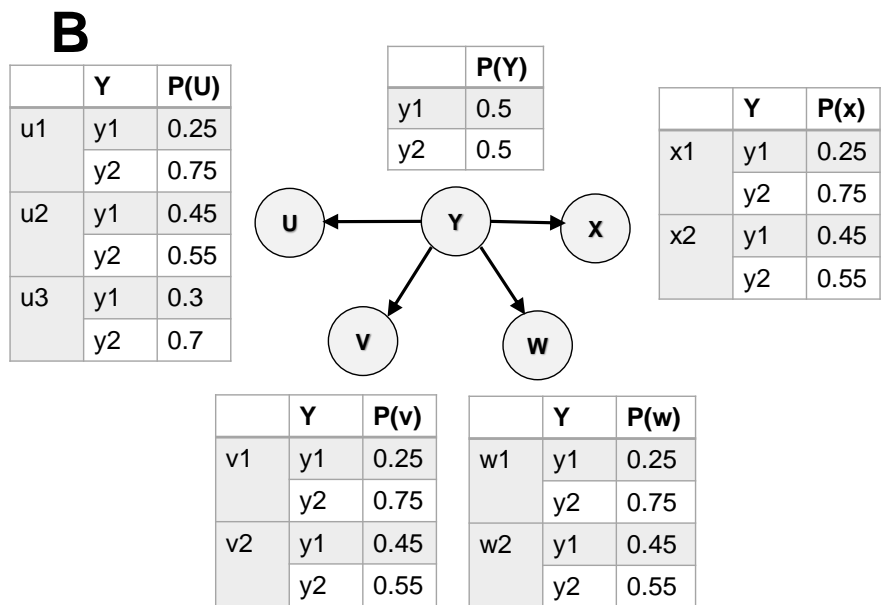
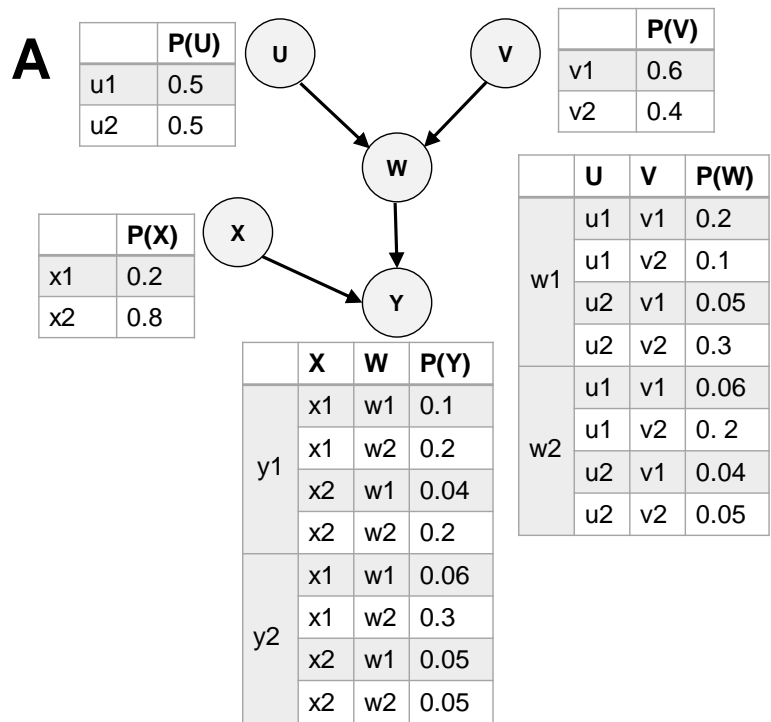
733

734

735

736

737 **Fig.3.** An example of A) a simple Bayesian Network and B) a Naïve Bayesian network

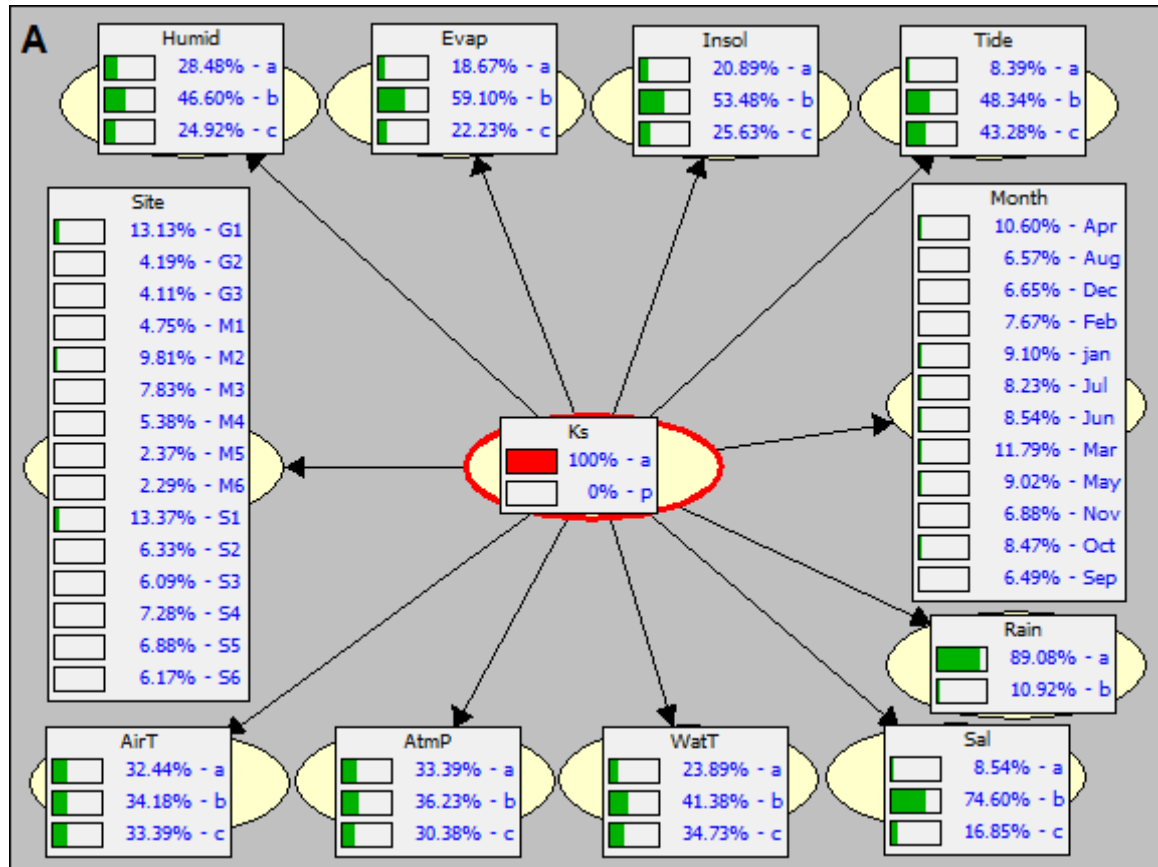


738

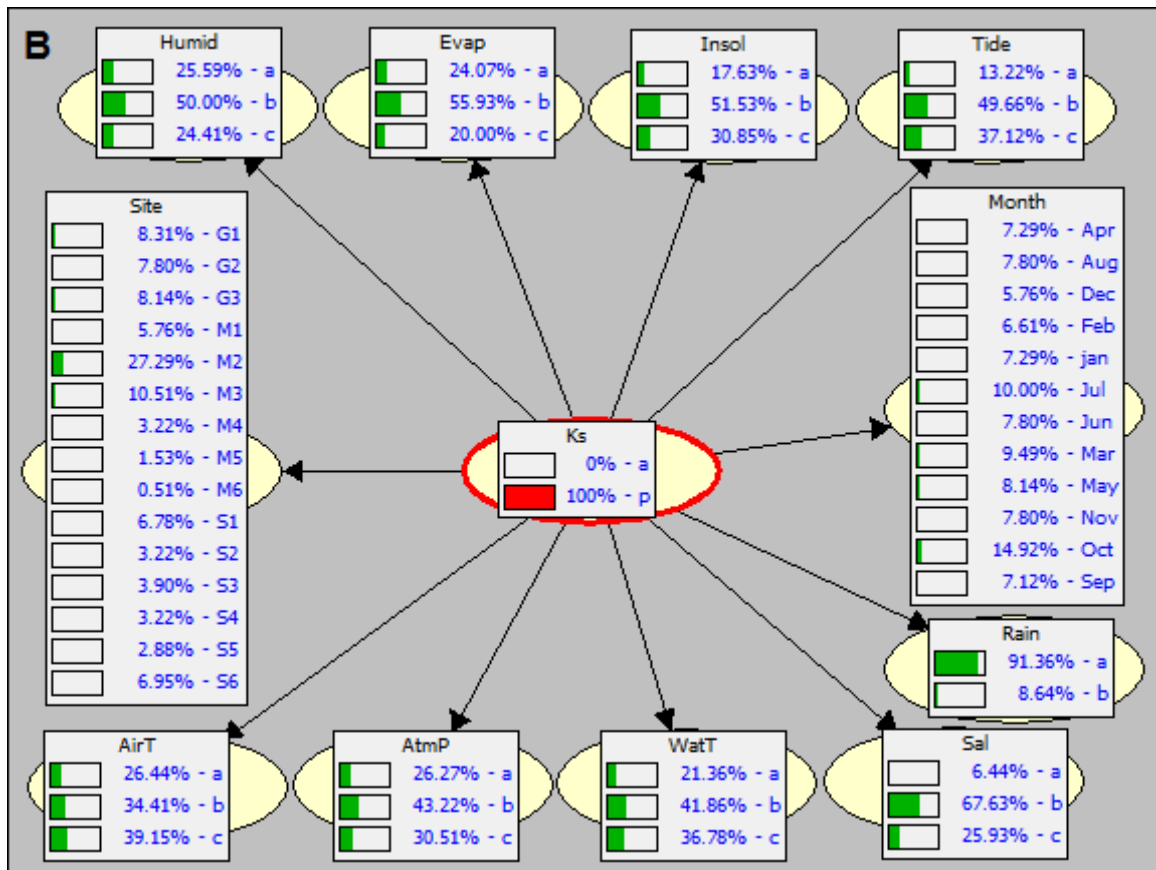
739

740

741 **Fig.4.** Naïve Bayesian network plot depicting the relationship between the hydro-
 742 meteorological parameters, spatio-temporal factors and the presence (A) and absence (B) of
 743 *Karenia selliformis* in the Gulf of Gabès.



744



745

746

747

748

749

750

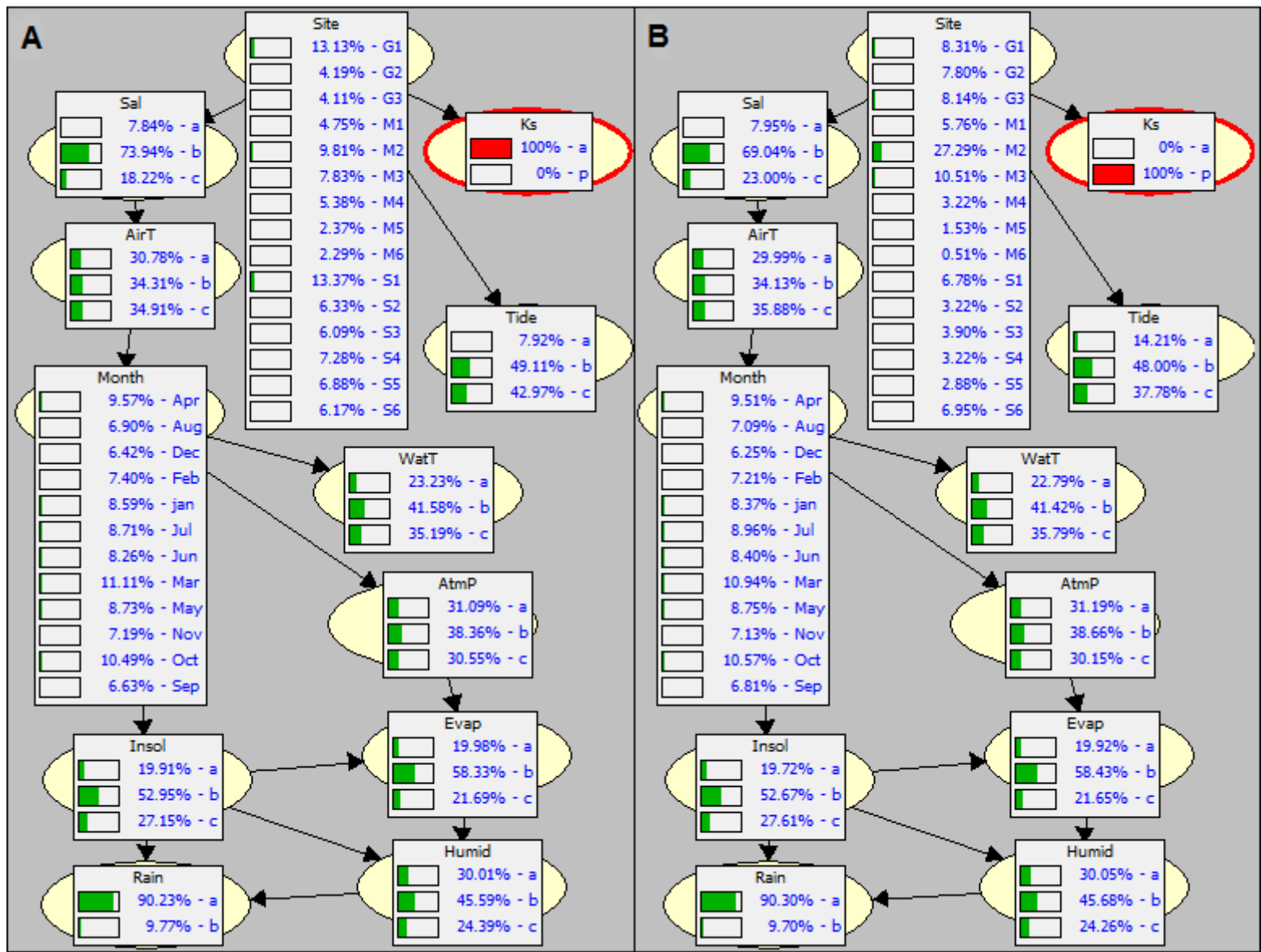
751

752

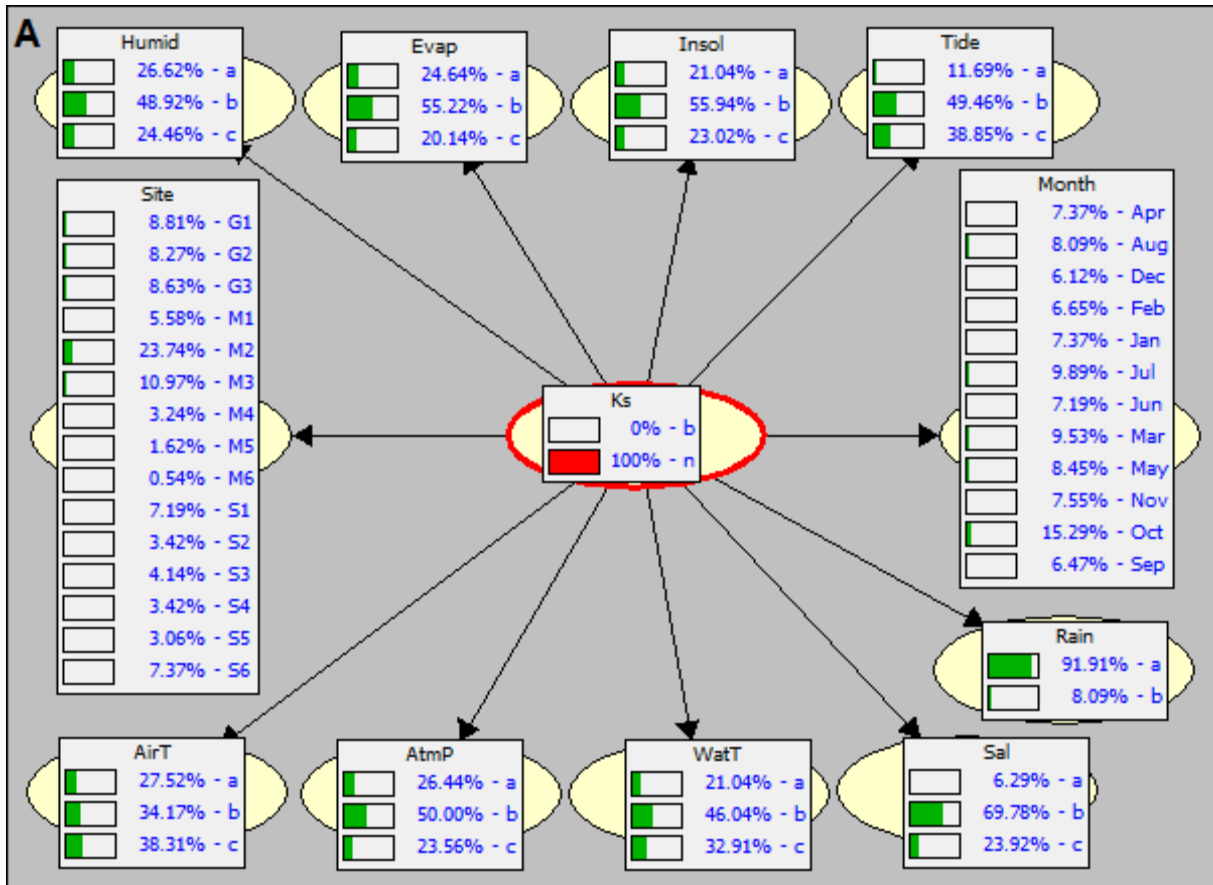
753 **Fig.5.** Bayesian network plot depicting the relationship between the hydro-meteorological

754 parameters, spatio-temporal factors and the presence (A) and absence (B) of *Karenia selliformis*

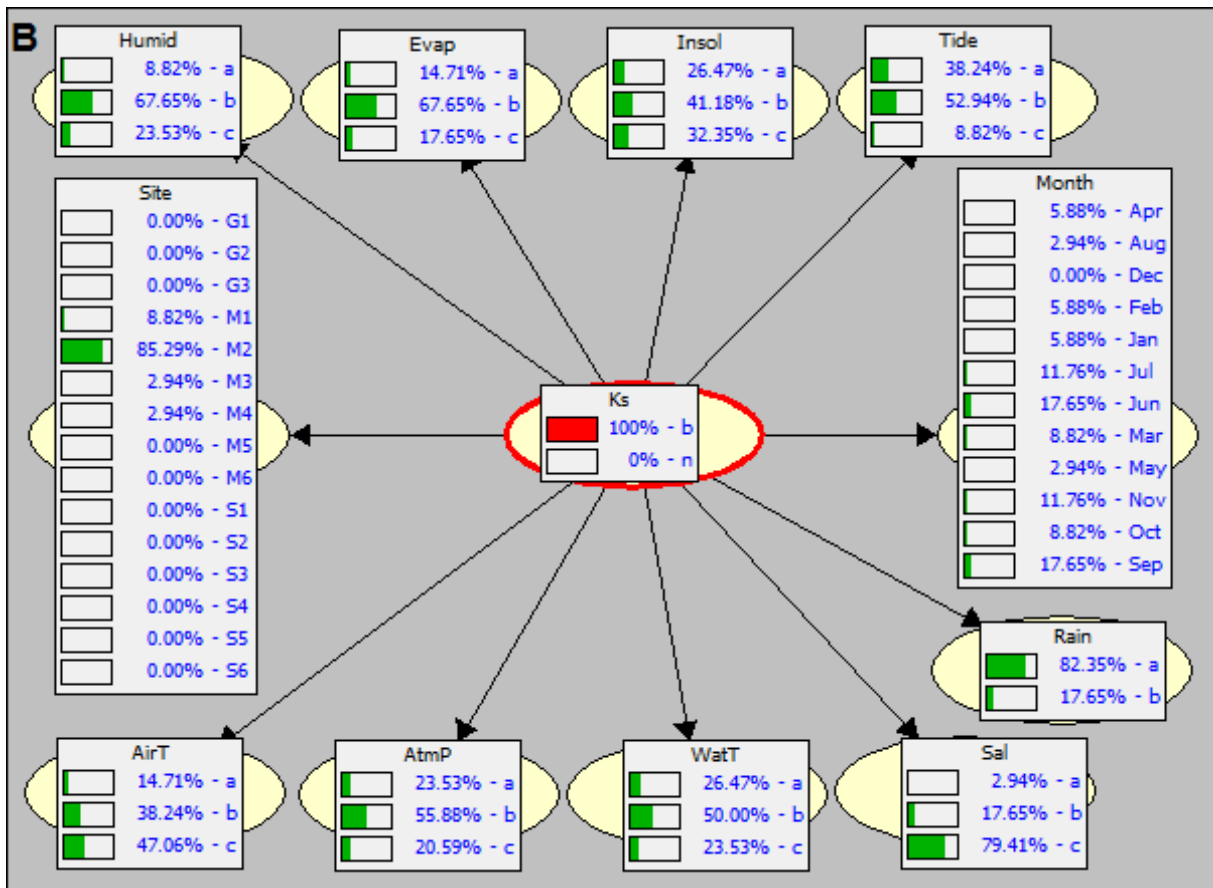
755 in the Gulf of Gabès.



757 **Fig.6.** Naïve Bayesian network plot depicting the relationship between the hydro-
 758 meteorological parameters, spatio-temporal factors and the non-bloom (A) and bloom (B) of
 759 *Karenia selliformis* in the Gulf of Gabès.



760



761

762

763

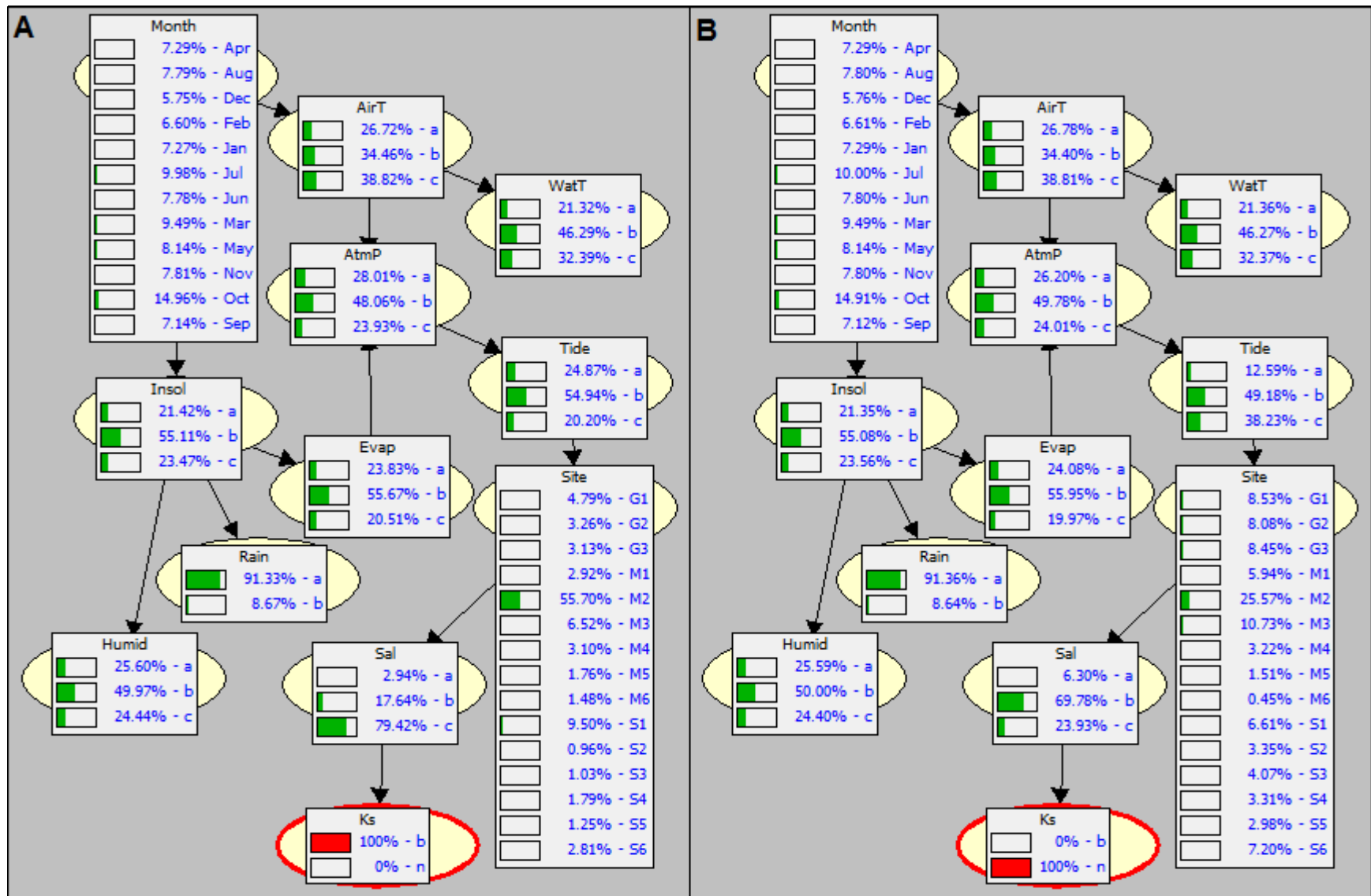
764

765

766

767

768 **Fig.7.** Bayesian network plot depicting the relationship between the hydro-meteorological
 769 parameters, spatio-temporal factors and the bloom (A) and non-bloom (B) of *Karenia*
 770 *selliformis* in the Gulf of Gabès.



771

772

773 **Table 1**

774 The biological and hydro-meteorological parameters discretized into intervals.

775 Rainfall (Rain), Evaporation (Evap), Air temperature (AirT), Insolation (Insol), Humidity (Humid), Atmospheric pressure (AtmP), Tide amplitude

776 (Tide), Water temperature (WatT), Salinity (Sal) and *Karenia selliformis* (Ks).

777

Rain	Evap	AirT	Insol	Humid	AtmP	Tide	WatT	Sal	Ks	Ks
(mm)	(mm)	(°C)	(h)	(%)	(Pa)	(m)	(°C)			
a=0	a<20	a<10	a<8	a<52	a<1011	a=0	a<14	a<37	a=absent	b=bloom
b>0	20≤b≤50	10≤b≤30	8≤b≤10	52≤b≤80	1011≤b≤1019	0≤b≤1.5	14≤b<29	37≤b<42.5	p=present	n=non-bloom
	c>50	c>30	c>10	c>80	c>1019	c>1.5	c≥29	c≥42.5		

778

Table 2

Model accounting for the observed variation of *Karenia selliformis* occurrences in the Gulf of Gabès according to the results of the general linear regression model analyses (GLM) with stepwise both selections of variables.

	AIC
Ks ~ Site + Month + Rain + Evap + AirT + Insol + Humid + AtmP + Tide + WatT + Sal	2158.81
Ks ~ Site + Month + Rain + Evap + AirT + Insol + AtmP + Tide + WatT + Sal	2154.98
Ks ~ Site + Month + Rain + Evap + Insol + AtmP + Tide + WatT + Sal	2151.16
Ks ~ Site + Month + Rain + Evap + Insol + Tide + WatT + Sal	2147.46
Ks ~ Site + Month + Rain + Evap + Insol + WatT + Sal	2144.66
Ks ~ Site + Month + Rain + Evap + Insol + Sal	2144.24

	Variable	Estimate	Std. error	z value	P value
(Intercept)	-1.84	0.36	-5.09	0	***
SiteG2	1.16	0.27	4.35	0	***
SiteG3	1.22	0.26	4.6	0	***
SiteM1	0.76	0.28	2.71	0.01	**
SiteM2	1.44	0.23	6.36	0	***
SiteM3	0.87	0.24	3.66	0	***
SiteM4	-0.1	0.32	-0.31	0.75	
SiteM5	-0.1	0.43	-0.23	0.82	
SiteM6	-1.22	0.64	-1.9	0.06	.
SiteS1	-0.32	0.25	-1.26	0.21	
SiteS2	-0.11	0.31	-0.36	0.72	
SiteS3	0.16	0.3	0.53	0.6	

SiteS4	-0.29	0.31	-0.95	0.34	
SiteS5	-0.56	0.32	-1.76	0.08	.
SiteS6	0.5	0.26	1.94	0.05	.
MonthAug	0.5	0.28	1.77	0.08	.
MonthDec	0.52	0.3	1.77	0.08	.
MonthFeb	0.45	0.28	1.61	0.11	
Monthjan	0.41	0.28	1.47	0.14	
MonthJul	0.23	0.28	0.84	0.4	
MonthJun	0	0.28	0.01	1	
MonthMar	0.3	0.26	1.18	0.24	
MonthMay	0.2	0.27	0.77	0.44	
MonthNov	0.94	0.28	3.37	0	***
MonthOct	1.31	0.26	5.06	0	***
MonthSep	0.54	0.28	1.92	0.05	.
Rainb	-0.31	0.21	-1.46	0.15	
Evapb	-0.36	0.15	-2.47	0.01	*
Evapc	-0.28	0.19	-1.48	0.14	
Insolb	-0.08	0.16	-0.5	0.62	
Insolc	0.37	0.23	1.63	0.1	
Salb	0.32	0.22	1.48	0.14	
Salc	0.74	0.26	2.85	0	**

Table 3

Model accounting for the observed variation of *Karenia selliformis* blooms in the Gulf of Gabès according to the results of the general linear regression model analyses (GLM) with stepwise both selections of variables.

	AIC
Ks ~ Site + Month + Rain + Evap + AirT + Insol + Humid + AtmP + Tide + WatT + Sal	232.86
Ks ~ Site + Month + Rain + Evap + AirT + Humid + AtmP + Tide + WatT + Sal	229.04
Ks ~ Site + Month + Rain + Evap + Humid + AtmP + Tide + WatT + Sal	226.91
Ks ~ Site + Month + Rain + Evap + Humid + Tide + WatT + Sal	224.22
Ks ~ Site + Month + Rain + Evap + Humid + WatT + Sal	222.14
Ks ~ Site + Month + Evap + Humid + WatT + Sal	220.15

	Variable	Estimate	Std. error	z value	P value
(Intercept)	22.54	3495.00	0.01	0.99	
SiteG2	-0.66	5103.00	0.00	1.00	
SiteG3	-0.54	5085.00	0.00	1.00	
SiteM1	-19.80	3495.00	-0.01	1.00	
SiteM2	-20.45	3495.00	-0.01	1.00	
SiteM3	-18.18	3495.00	-0.01	1.00	
SiteM4	-18.56	3495.00	-0.01	1.00	
SiteM5	-0.14	10050.00	0.00	1.00	
SiteM6	-0.26	16060.00	0.00	1.00	
SiteS1	0.12	5264.00	0.00	1.00	
SiteS2	-0.04	6711.00	0.00	1.00	
SiteS3	-1.74	6627.00	0.00	1.00	

SiteS4	-0.82	6673.00	0.00	1.00	
SiteS5	-0.28	7086.00	0.00	1.00	
SiteS6	-0.58	5319.00	0.00	1.00	
MonthAug	0.43	1.52	0.28	0.78	
MonthDec	19.68	3962.00	0.01	1.00	
MonthFeb	0.65	1.36	0.48	0.63	
MonthJan	0.98	1.26	0.78	0.44	
MonthJul	0.90	1.20	0.76	0.45	
MonthJun	-0.37	1.08	-0.34	0.73	
MonthMar	1.17	1.18	0.99	0.32	
MonthMay	1.88	1.40	1.34	0.18	
MonthNov	-1.74	1.10	-1.59	0.11	
MonthOct	1.03	1.23	0.84	0.40	
MonthSep	-0.86	1.17	-0.73	0.46	
Evapb	0.54	0.72	0.75	0.45	
Evapc	1.59	0.87	1.83	0.07	.
Humidb	-2.12	0.89	-2.40	0.02	*
Humidc	-1.83	0.97	-1.88	0.06	.
WatTb	0.40	0.74	0.54	0.59	
WatTc	1.57	0.90	1.75	0.08	.
Salb	0.41	1.38	0.30	0.77	
Salc	-1.02	1.39	-0.74	0.46	

Table 4

Adjusted R^2 values for the three models: General linear model (GLM). Bayesian networks (BN) and naive Bayes classifier (NB)

Adjusted R^2		
	Occurrence	Bloom
BN	0.70	0.58
NB	0.69	0.54
GLM	0.17	0.48

Table 5

Comparative analysis of the main fundamentals of the three predictive models: BN, NB and GLM. The lines separate models that do not have the same characteristics.

GLM	NB	BN
Can handle continues or discrete data	Exclusive use of discrete variables which often results in a significant information loss	
Not informative about potential causal pathways between predictive variables		More informative about potential causal pathways and able to easy establish direct associations between variables
The interdependencies between variables were not discovered		The interdependencies between variables were revealed
The relationships between variables and outcome are linear		Can handle nonlinear relationships between variables and outcome
	Less confidence on the predicted probabilities	More confidence on the predicted probabilities and may better reflect the joint probability distribution over the system's variables
	The number of arcs reduced the model performance	The number of arcs has no significant effect on model performance and the BN

		improves usually the accuracy, flexibility and expressivity
	The introduction of the evidence changes the probability distribution of several variables and amplify the weight of the evidence of each attribute on the class	The introduction of the evidence do not change the probabilities of variables except that for the parent of the class