

Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

## Data &amp; Knowledge Engineering

journal homepage: [www.elsevier.com/locate/datak](http://www.elsevier.com/locate/datak)

## Circular Bayesian classifiers using wrapped Cauchy distributions

Ignacio Leguey<sup>a,b,\*</sup>, Concha Bielza<sup>b</sup>, Pedro Larrañaga<sup>b</sup><sup>a</sup> Departamento de Economía Financiera y Contabilidad e Idioma Moderno, Facultad de Ciencias Jurídicas y Sociales, Universidad Rey Juan Carlos de Madrid, Spain<sup>b</sup> Departamento de Inteligencia Artificial, Universidad Politécnica de Madrid, Campus de Montegancedo, 28660 Boadilla del Monte, Madrid, Spain

## ARTICLE INFO

## Keywords:

Data mining  
 Classification  
 Circular statistics  
 Wrapped Cauchy distribution  
 Bayesian networks  
 Cortical layer

## ABSTRACT

Capturing the dependences among circular variables within supervised classification models is a challenging task. In this paper, we propose four different supervised Bayesian classification algorithms where the predictor variables follow all circular wrapped Cauchy distributions. For this purpose, we introduce four wrapped Cauchy classifiers. The bivariate wrapped Cauchy distribution is the only bivariate circular distribution whose marginals and conditionals are also wrapped Cauchy distributions, a property that makes it possible to define these models easily. Furthermore, the wrapped Cauchy tree-augmented naive Bayes classifier requires the definition of a conditional circular mutual information measure between variables that follow wrapped Cauchy distributions. Synthetic data is used to illustrate, compare and evaluate the classification algorithms (including a comparison with the Gaussian TAN classifier, decision tree, random forest, multinomial logistic regression, support vector machine and simple neural network), leading to satisfactory predictive results. We also use a real neuromorphological dataset obtained from juvenile rat somatosensory cortex cells, where we measure the bifurcation angles of the dendritic basal arbors.

## 1. Introduction

Circular data is ubiquitous, present in many different areas such as biology, geology, medicine, oceanography, geophysics, meteorology, astronomy, ecology, neuroscience and geography. Some examples are directions of flight of homing pigeons [1], characterization of the phenology of species [2], formation of feldspar laths in basalt rocks [3], paleomagnetism in red slits and claystones [4,5], also in political sciences studying the gun crimes occurred in an specific period of time [6], directional word vectors in text mining [7], wildfire orientation in order to prevent fire propagation [8], wind and waves direction analysis [9,10], study and prediction of protein dihedral angles structure [11,12], and neuronal basal dendritic bifurcation angles analysis [13,14] among many others. The natural periodicity of circular data sometimes makes traditional statistics methods ineffective, since they ignore this characteristic. For instance, when dealing with circular data,  $0^\circ$  and  $360^\circ$  are considered as the same point, whereas if considered non-circular data, they are different points. Thus, circular data analysis is distinct from and more challenging than non-circular data. It should be noted that circular data has been studied extensively [3,15,16].

Probabilistic graphical models [17] are useful tools for data modeling that connect probability theory with graph theory. There are many advantages of using probabilistic graphical models, such as the fact that they are easily interpreted, they handle missing data effectively and they treat inference and learning tasks together. Bayesian networks [18] are one of the most commonly used probabilistic graphical models due to their factorization and domain representation properties. Bayesian networks have the

\* Corresponding author at: Departamento de Economía Financiera y Contabilidad e Idioma Moderno, Facultad de Ciencias Jurídicas y Sociales, Universidad Rey Juan Carlos de Madrid, Spain.

E-mail addresses: [ignacio.vitoriano@urjc.es](mailto:ignacio.vitoriano@urjc.es) (I. Leguey), [mcbielza@fi.upm.es](mailto:mcbielza@fi.upm.es) (C. Bielza), [pedro.larranaga@fi.upm.es](mailto:pedro.larranaga@fi.upm.es) (P. Larrañaga).

<https://doi.org/10.1016/j.datak.2019.05.005>

Received 23 October 2017; Received in revised form 24 November 2018; Accepted 25 May 2019

Available online 28 May 2019

0169-023X/© 2019 Elsevier B.V. All rights reserved.

characteristic that each variable is conditionally independent of those that are non-descendants in the graph given the value of their parents. Therefore the joint probability distribution is expressed as the product of the local distributions conditioned to their parents. For these reasons, Bayesian networks can deal efficiently with supervised classification i.e., the Bayesian network classifiers [13] and offer an explicit, graphical and interpretable representation of uncertain knowledge, which has made it possible to successfully apply them to real-world problems.

Supervised classification [19] deals with the problem of assigning a label to an instance, based on a set of variables that characterize it. Yet circular data has been commonly treated as linear data in supervised classification tasks. Only a few circular classifiers exist, and almost none of them are based on the principles of Bayesian networks, capable of capturing multivariate relationships among variables. Most of them focus on discriminant analysis and assume several circular distributions such as the von Mises distribution [20], later extended to the von Mises–Fisher distribution [21]. There are also circular discriminant analysis studies for the Watson, Selby and Arnold distributions on the sphere [22,23]. SenGupta and Roy [24] used a classification discriminant rule based on the mean chord-length to classify a new observation into one of two different circular populations that are von Mises, when training samples are available for each of them. Also a likelihood ratio test based on a bootstrapping approach for classifying into two populations was proposed for linear and circular data [25]. Kirby and Miranda [26] proposed a variation of a neural network, including a circular node, which was able to keep and send circular information. More recently, Fernandes and Cardoso [27] proposed a binary circular logistic regression as the discriminative counterpart to the naive Bayes model, which does not make assumptions on the input data distribution. López-Cruz et al. [28] is the only study in which Bayesian classifiers were used. For the von Mises and von Mises–Fisher distributions, López-Cruz et al. proposed an adaptation of the naive Bayes classifier and selective naive Bayes classifier, which are two of the simplest and best-known supervised classification models based on Bayesian network principles.

The lack of Bayesian supervised classifiers for circular data is due to the absence of circular Bayesian network models, which are very difficult to develop because of their circular multivariate distribution nature. A family of distributions is said to be closed under marginalization and conditioning when the marginals and conditionals of the multivariate distribution follow the same distribution. However, the marginals and conditionals of most circular distributions do not belong to the same family of distributions, making the modeling phase and posterior inference processes difficult.

The von Mises distribution [29], which is the analogue of the univariate Gaussian distribution, is the best-known circular model. A bivariate von Mises distribution also exists and was introduced by Mardia [30], who subsequently extended it to the multivariate case [31]. He showed that the conditional distributions are also von Mises distributions. Nevertheless, the marginal distributions are either unimodal or bimodal, and only the unimodal case could be approximated to a von Mises distribution when the concentration parameter is large. Therefore, as explained in [13] for discrete distributions, it would be much more complicated to achieve an efficient learning and inference. Therefore, we ruled out the use of von Mises distributions for our particular purpose. Another popular univariate symmetric circular distribution is the wrapped Cauchy, which was introduced by Lévy [32], and further studied by Wintner [33]. It was later obtained by mapping Cauchy distributions onto the circle [34]. Kato and Pewsey [35] developed a five-parameter bivariate wrapped Cauchy distribution for toroidal data, whose marginals and conditionals follow univariate wrapped Cauchy distributions. This family of bivariate wrapped Cauchy distributions is therefore closed under conditioning and marginalization. Leguey et al. [36] proposed a tree-structured Bayesian network model that deals with circular data which follows wrapped Cauchy distribution. However, this model only accounts for the discovery of conditional independence relationships of a set of random variables, without considering any as a class variable. This is a specificity of supervised classification problems that requires special learning algorithms.

Building on previous work regarding supervised classification using Bayesian networks for circular statistics, the novelty of this work lies in the proposal of four circular Bayesian classification models capable of dealing with supervised data following wrapped Cauchy distributions. The models to be presented are called wrapped Cauchy naive Bayes (wCNB), wrapped Cauchy selective naive Bayes (wCsNB), wrapped Cauchy semi-naive Bayes (wCsmNB) and wrapped Cauchy tree-augmented naive Bayes (wCTAN) classifiers. Even though the simplest of our proposals (i.e., wCNB) is a straightforward naive Bayes classifier extension to using wrapped Cauchy distributions, this has never been attempted before to the best of our knowledge.

The remainder of this paper is organized as follows. Section 2 reviews the bivariate wrapped Cauchy distribution of Kato and Pewsey [35]. Section 3 describes the four novel wrapped Cauchy classifiers proposed here. In Section 4, we assess the four models in synthetic domains, requiring the design of a simulation method for these wrapped Cauchy Bayesian network classifiers. Section 5 addresses a real-world neuromorphology data problem using the wrapped Cauchy classifiers. Finally, Section 6 provides concluding remarks and proposals for future work.

## 2. Wrapped Cauchy distribution

### 2.1. Definitions

A random variable  $\theta$  that follows a wrapped Cauchy distribution [32], denoted  $wC(\mu, \varepsilon)$ , has a density function

$$f(\theta) = \frac{1}{2\pi} \frac{1 - \varepsilon^2}{1 + \varepsilon^2 - 2\varepsilon \cos(\theta - \mu)}, \quad \theta, \mu \in (-\pi, \pi], \varepsilon \in [0, 1) \quad (1)$$

where  $\mu$  is the mean angle and  $\varepsilon$  the concentration parameter.  $f$  in Eq. (1) is unimodal and symmetric about  $\mu$  unless  $\varepsilon = 0$ , which yields the circular uniform distribution (i.e.,  $f(\theta) = 1/2\pi$ ).

A five-parameter bivariate wrapped Cauchy distribution was proposed by Kato and Pewsey [35]. A random vector  $(\theta_1, \theta_2)$  follows a bivariate wrapped Cauchy distribution, denoted  $buC(\mu_1, \mu_2, \varepsilon_1, \varepsilon_2, \rho)$ , if its density function is given by

$$f(\theta_1, \theta_2) = c[c_0 - c_1 \cos(\theta_1 - \mu_1) - c_2 \cos(\theta_2 - \mu_2) - c_3 \cos(\theta_1 - \mu_1) \cos(\theta_2 - \mu_2) - c_4 \sin(\theta_1 - \mu_1) \sin(\theta_2 - \mu_2)]^{-1}, \theta_1, \theta_2 \in (-\pi, \pi], \tag{2}$$

where  $c, c_0, c_1, c_2, c_3$  and  $c_4$  are

$$\begin{aligned} c &= (1 - \rho^2)(1 - \varepsilon_1^2)(1 - \varepsilon_2^2)/4\pi^2 \\ c_0 &= (1 + \rho^2)(1 + \varepsilon_1^2)(1 + \varepsilon_2^2) - 8|\rho|\varepsilon_1\varepsilon_2 \\ c_1 &= 2(1 + \rho^2)\varepsilon_1(1 + \varepsilon_2^2) - 4|\rho|(1 + \varepsilon_1^2)\varepsilon_2 \\ c_2 &= 2(1 + \rho^2)(1 + \varepsilon_1^2)\varepsilon_2 - 4|\rho|\varepsilon_1(1 + \varepsilon_2^2) \\ c_3 &= -4(1 + \rho^2)\varepsilon_1\varepsilon_2 + 2|\rho|(1 + \varepsilon_1^2)(1 + \varepsilon_2^2) \text{ and} \\ c_4 &= 2\rho(1 - \varepsilon_1^2)(1 - \varepsilon_2^2), \end{aligned}$$

with  $\mu_1, \mu_2 \in (-\pi, \pi]$ ,  $\varepsilon_1, \varepsilon_2 \in [0, 1)$  and  $\rho \in (-1, 1)$ .  $\varepsilon_1$  and  $\varepsilon_2$  regulate the concentration of the marginal distributions, and  $\rho$  is the parameter controlling the association between  $\theta_1$  and  $\theta_2$ , from total independence ( $\rho = 0$ ) to perfect correlation ( $\rho = \pm 1$ ). When  $\varepsilon_1 > 0$  and  $\varepsilon_2 > 0$ ,  $f$  in Eq. (2) is unimodal and pointwise symmetric about  $(\mu_1, \mu_2)$ .

From McCullagh [34], we know that representing wrapped Cauchy models in complex form simplifies the computation in many cases. Let  $Z = \exp(i\theta)$ , where  $\theta$  is distributed as in Eq. (1). Therefore the density function of  $Z$  is

$$f(z; \lambda) = \frac{1}{2\pi} \frac{|1 - |\lambda|^2|}{|z - \lambda|^2}, \quad z \in \Omega, \lambda \in \hat{\mathbb{C}} \setminus \Omega, \tag{3}$$

where  $\lambda = \varepsilon \exp(i\mu)$ ,  $\hat{\mathbb{C}} = \mathbb{C} \cup \{\infty\}$  and  $\Omega = \{z \in \mathbb{C} : |z| = 1\}$ , with  $|z|$  the module of the complex number. We use the notation  $Z \sim C^*(\lambda)$  to denote that  $Z$  is distributed as in Eq. (3).

Similarly, let  $(Z_1, Z_2) = (\exp(i\theta_1), \exp(i\theta_2))$ , where  $(\theta_1, \theta_2)$  is distributed as in Eq. (2). Therefore the density of  $(Z_1, Z_2)$  is

$$f(z_1, z_2) = \frac{(4\pi^2)^{-1}(1 - \rho^2)(1 - \varepsilon_1^2)(1 - \varepsilon_2^2)}{|a_{11}(\bar{z}_1\eta_1)^q z_2\bar{\eta}_2 + a_{12}(\bar{z}_1\eta_1)^q + a_{21}z_2\bar{\eta}_2 + a_{22}|^2}, z_1, z_2 \in \Omega, \tag{4}$$

where  $q$  is the sign of  $\rho$ ,  $\eta_k = \exp(i\mu_k)$  with  $k \in \{1, 2\}$ ,  $\bar{z}_k$  is the complex conjugate of  $z_k$ ,  $a_{11} = \varepsilon_1\varepsilon_2 - |\rho|$ ,  $a_{12} = |\rho|\varepsilon_2 - \varepsilon_1$ ,  $a_{21} = |\rho|\varepsilon_1 - \varepsilon_2$ ,  $a_{22} = 1 - |\rho|\varepsilon_1\varepsilon_2$ ,  $\varepsilon_1, \varepsilon_2 \in [0, 1)$ ,  $\rho \in (-1, 1)$  and  $\eta_1, \eta_2 \in \Omega$ .

Following the complex notation, we denote  $(Z_1, Z_2) \sim bC^*(\eta_1, \eta_2, \varepsilon_1, \varepsilon_2, \rho)$  if  $(Z_1, Z_2)$  is distributed as in Eq. (4). This five-parameter bivariate wrapped Cauchy complex form representation verifies the following result:

**Theorem 1** (Kato and Pewsey [35]). *A random vector  $(Z_1, Z_2)$  with density given by Eq. (4) has marginals  $Z_1 \sim C^*(\varepsilon_1\eta_1)$  and  $Z_2 \sim C^*(\varepsilon_2\eta_2)$ , and conditionals  $Z_1|Z_2 = z_2 \sim C^*(-\eta_1[A \circ (z_2\bar{\eta}_2)^q])$  and  $Z_2|Z_1 = z_1 \sim C^*(-\eta_2[A^T \circ (z_1\bar{\eta}_1)^q])$ , where  $A$  is the matrix with elements  $a_{11}, a_{12}, a_{21}$  and  $a_{22}$  defined in Eq. (4),  $A^T$  is the transpose of  $A$ , and*

$$A \circ z = \begin{pmatrix} a_{11}z + a_{12} \\ a_{21}z + a_{22} \end{pmatrix}.$$

As far as we know, there is no other bivariate circular distribution for which conditional and marginal distributions belong to the same family. Therefore, wrapped Cauchy is the best choice given no better alternative, as the requirements for the classifier structures that we will develop are of at most a tree-structure (i.e., only bivariate, marginal and conditional densities are required). Furthermore, we require the definition of a conditional circular mutual information measure between variables that follow wrapped Cauchy distributions.

### 2.2. Parameter estimation

Working with the density given by Eq. (2), numerical methods have to be used to find the parameter estimates, since there is no closed-form expression for the maximum likelihood estimates. Kato and Pewsey [35] demonstrated that the method of moments [37] is more efficient; it is computationally very fast, easy to implement and with closed form formulas for the parameter estimates.

Let  $\{(\theta_{1j}, \theta_{2j}), j = 1, \dots, N\}$  be a random sample from a  $buC(\mu_1, \mu_2, \varepsilon_1, \varepsilon_2, \rho)$  as stated in Eq. (2). Therefore the estimators obtained from the method of moments for  $\mu_1, \mu_2, \varepsilon_1, \varepsilon_2$  and  $\rho$  are [35]

$$\begin{aligned} \hat{\mu}_r &= \arg(\bar{R}_r), \quad \hat{\varepsilon}_r = |\bar{R}_r|, \quad \text{with} \quad \bar{R}_r = \frac{1}{N} \sum_{j=1}^N \exp(i\theta_{rj}), \\ \hat{\rho} &= \frac{1}{N} \left( \left| \sum_{j=1}^N \exp(i(\Phi_{1j} - \Phi_{2j})) \right| - \left| \sum_{j=1}^N \exp(i(\Phi_{1j} + \Phi_{2j})) \right| \right), \\ \text{with} \quad \Phi_{rj} &= 2 \arctan \left( \frac{1 + \hat{\varepsilon}_r}{1 - \hat{\varepsilon}_r} \tan \left( \frac{\theta_{rj} - \hat{\mu}_r}{2} \right) \right) \quad \text{and} \quad r = 1, 2. \end{aligned} \tag{5}$$

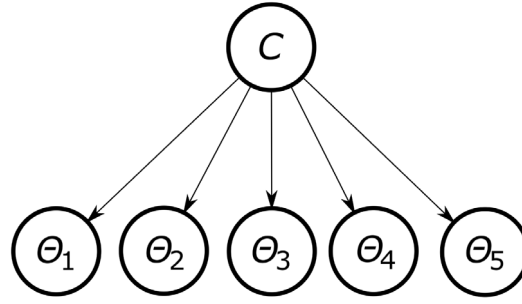


Fig. 1. wCNB structure with five circular predictor nodes, from which  $p(c|\theta) \propto p(c)f_{\Theta_1|c}(\theta_1|c)f_{\Theta_2|c}(\theta_2|c)f_{\Theta_3|c}(\theta_3|c)f_{\Theta_4|c}(\theta_4|c)f_{\Theta_5|c}(\theta_5|c)$ .

### 3. Wrapped Cauchy classifiers

Let  $\Theta = (\Theta_1, \dots, \Theta_n)$  be a vector of circular predictor random variables or features, and let  $C$  be a discrete class variable which takes values (labels) in the set  $\Lambda(C)$ . Given a sample of  $N$  labeled instances  $(\Theta^1, C^1), \dots, (\Theta^N, C^N)$ , the supervised classification problem consists of developing a model capable of assigning a class label to a new object based on the values of its features.

Bayesian network classifiers [13] have been used to solve classification problems with linear data, because of their easy representation of the problem domain and the efficient computation of the algorithms associated with Bayesian networks techniques. Our novel purpose is to develop the circular domain counterpart of four well-known Bayesian network classifiers (naive Bayes, selective naive Bayes, semi-naive Bayes and tree-augmented naive Bayes) when the underlying variables follow wrapped Cauchy distributions.

#### 3.1. Wrapped Cauchy naive Bayes

The wrapped Cauchy naive Bayes (wCNB) classifier is the simplest of the four Bayesian network classifier models that we present in this paper, where  $C$  is the parent of all circular features and these are assumed to be conditionally independent among them given  $C$  (Fig. 1)

$$p(C = c|\Theta = \theta) \propto p(C = c) \prod_{i=1}^n f_{\Theta_i|C=c}(\theta_i|c). \tag{6}$$

The wCNB determines the class value  $c^*$  for a new instance using a maximum a posteriori decision rule

$$c^* = \arg \max_{c \in \Lambda(C)} p(C = c|\Theta = \theta).$$

Since each predictor variable  $\theta_i$  given  $C = c$  follows a wrapped Cauchy distribution with location parameter  $\mu_{i,c}$  and concentration parameter  $\varepsilon_{i,c}$ , we can express Eq. (6) as

$$p(c|\theta) \propto \frac{p(C = c) \prod_{i=1}^n \alpha_{i,c}}{\prod_{i=1}^n (1 - \beta_{i,c})}, \tag{7}$$

where  $\alpha_{i,c} = \frac{1 - \varepsilon_{i,c}^2}{2\pi(1 + \varepsilon_{i,c}^2)}$  and  $\beta_{i,c} = \frac{2\varepsilon_{i,c} \cos(\theta_i - \mu_{i,c})}{(1 + \varepsilon_{i,c}^2)}$ .

#### 3.2. Wrapped Cauchy selective naive Bayes

Sometimes there are several predictor variables that do not contribute to classification (i.e., they are redundant), and naive Bayes classifier is affected by such variables [38]. Determining which of them are unnecessary via the use of feature subset selection (FSS) techniques [39] could increase the accuracy of the classification model significantly [40]. Wrapped Cauchy selective naive Bayes (wCsNB) is a classification model with a structure similar to that of wCNB, but not all the variables are necessarily used by the classifier. FSS techniques were previously employed in a circular classification model with von Mises and von Mises–Fisher distributions in [28], where a filter-wrapper algorithm is applied to rank the variables according to the mutual information between them and the class, and therefore, using the ranking provided by the filter step, the variables are selected to induce a new classifier until the best model is achieved.

We also use a filter-wrapper algorithm at this point. The filter step is based on the computation of the mutual information (MI) between each circular variable and the class. There is no equation to compute the MI between circular variables and discrete variables. Therefore, we approach the problem using Monte Carlo methods, as in [28]; we model the conditional density functions of  $\theta_i|C = c$  as wrapped Cauchy distributions. Hence

$$MI(\theta_i, C) \approx \frac{1}{M} \sum_{j=1}^M \log \frac{\hat{f}_{\theta_i|c^{*(j)}}(\theta_i^{*(j)}|c^{*(j)}) \hat{p}(C = c^{*(j)})}{\hat{f}_{\theta_i}(\theta_i^{*(j)}) \hat{p}(C = c^{*(j)})}, \tag{8}$$

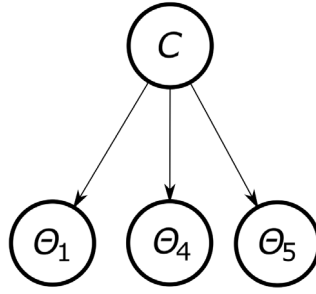


Fig. 2. wCsNB structure with three nodes selected from the original set of five predictive variables, from which  $p(c|\theta) \propto p(c)f_{\theta_1|c}(\theta_1|c)f_{\theta_4|c}(\theta_4|c)f_{\theta_5|c}(\theta_5|c)$ .

where  $M$  is the number of instances  $(\theta_i^{*(j)}, c^{*(j)})$  sampled from  $\hat{f}_{\theta_i|c}(\theta_i|c)\hat{p}(C=c)$ , with  $\hat{f}_{\theta_i|c}(\theta_i|c)$  the fitted wrapped Cauchy density function of the conditional density function of  $\theta_i$  given  $C=c$ , and  $\hat{p}(C=c)$  the relative frequency of instances that belong to class  $c$  in the training set.

The predictive variables are then ranked according to their MI values.

The wrapper step consists of creating a new classifier by deciding whether or not to include the ranked predictive variables from the filter step. Each iteration of the wrapper step induces a new classifier adding the next predictive variable from the list. If no accuracy improvement is achieved by including the next predictive variable from the ranked list, then the wrapper step finishes. This model is similar to the wCNB, but including only the selected wrapped Cauchy variables (set  $S$ ) (Fig. 2) and therefore

$$p(c|\theta) \propto p(c|\theta_S) = p(C=c) \prod_{i \in S} f_{\theta_i|C=c}(\theta_i|c). \tag{9}$$

As for the wCNB, the wCsNB determines the class value  $c^*$  for a new instance using a maximum a posteriori decision rule

$$c^* = \arg \max_{c \in \mathcal{A}(C)} p(C=c|\Theta_S = \theta_S).$$

Likewise for Eq. (7), we can express Eq. (9) as

$$p(c|\theta_S) = \frac{p(C=c) \prod_{i \in S} \alpha_{i,c}}{\prod_{i \in S} (1 - \beta_{i,c})},$$

where  $\alpha_{i,c} = \frac{1}{2\pi} \frac{(1 - \epsilon_{i,c}^2)}{(1 + \epsilon_{i,c}^2)}$  and  $\beta_{i,c} = \frac{2\epsilon_{i,c} \cos(\theta_{i,c} - \mu_{i,c})}{(1 + \epsilon_{i,c}^2)}$ .

### 3.3. Wrapped Cauchy semi-naive Bayes

Usually, the assumption of conditional independence between predictive variables given the class variable is dismissed. The semi-naive Bayes classification model [41] goes one step further and considers dependencies between predictive variables.

Our proposal for this model, called wrapped Cauchy semi-naive Bayes (wCsmNB) classifier, takes into account the possible dependence between predictive wrapped Cauchy variables by introducing new features obtained as the Cartesian product of two of the original circular predictor variables. Thus we work with a bivariate wrapped Cauchy distribution. These new features remains conditionally independent given the class variable.

Given  $L_k$  with  $k = 1, \dots, T$ , representing the  $k$ th feature (original or new features)

$$p(c|\theta) \propto p(C=c) \prod_{k=1}^T f_{\theta_{L_k}|C=c}(\theta_{L_k}|c).$$

To determine those original variables that are candidates to create new features from the Cartesian product between them, we develop an adaptation of the *forward sequential selection and joining* (FSSJ) algorithm [42] described in Algorithm 1. It is important to note that once the new features are created by joining two original features, these new features cannot be used to create others. However, these new features can be separated in order to use one of the two original features to create another new feature by joining with a different original feature that had not yet been added to the model. This algorithm may result in a selection of variables that provide the best achievable solution, before all of the original variables are included in the model (Fig. 3).

Again, as for the previous models presented in this section, the wCmNB determines the class value  $c^*$  for a new instance using a maximum a posteriori decision rule

$$c^* = \arg \max_{c \in \mathcal{A}(C)} p(C=c|\Theta = \theta_{L_k}).$$

**Algorithm 1** Adaptation of the FSSJ algorithm of [42]

- 1: Let  $T$  be the variable list, initialized as  $T = \emptyset$ .
- 2: Given  $\theta_1, \theta_2, \dots, \theta_n$  circular wrapped Cauchy predictor variables from a variable list  $A$ , move the first variable from  $A$  to  $T$ .
- 3: Move the next variable from  $A$  to  $T$ , considering:
  - Joining the variable to another variable currently in  $T$ . If the latter variable was previously joined to another variable from  $T$ , remove this from  $T$  and add it to  $A$ , and consider adding it later.
  - Add the variable as conditionally independent of the other variables given  $C$  to the current classifier.
- 4: Repeat Step 3 until the best model is achieved

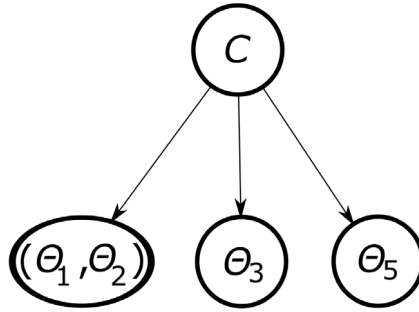


Fig. 3. wCsmNB structure with four nodes from the original set of five predictive variables, from which  $p(c|\theta) \propto p(c)f_{\theta_1, \theta_2|c}(\theta_1, \theta_2|c)f_{\theta_3|c}(\theta_3|c)f_{\theta_5|c}(\theta_5|c)$ .

3.4. Wrapped Cauchy tree-augmented naive Bayes

The tree-augmented naive Bayes (TAN) classifier [43] is a well-known Bayesian classifier with a tree-structure network for predictive features. Wrapped Cauchy tree-augmented naive Bayes (wCTAN) classifier is a variation of the TAN classifier with the novelty of the allowance of the use of wrapped Cauchy circular variables for predictive features. wCTAN assumes that the class variable has no parents, and the rest of the variables have at most one other variable as parent apart from  $C$  (Fig. 4).

The process for building a wCTAN is summarized in the following three steps:

- Step 1: The structure of the tree for predictive features is learned using Algorithm 2. We use the conditional circular mutual information, denoted as  $CMI(\theta_i, \theta_j|C)$ , which is defined as

$$CMI(\theta_i, \theta_j|C) = \sum_c CMI(\theta_i, \theta_j|C = c)p(C = c),$$

with

$$CMI(\theta_i, \theta_j|C = c) = \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} f(\theta_i, \theta_j|c) \log \left( \frac{f(\theta_i, \theta_j|c)}{f(\theta_i|c)f(\theta_j|c)} \right) d\theta_i d\theta_j.$$

where the marginal density functions given the class,  $f(\theta_i|c)$  and  $f(\theta_j|c)$ , and the joint density function given the class,  $f(\theta_i, \theta_j|c)$ , have been previously estimated from data. This structure learning algorithm (Algorithm 2) is based on score and search, where structure learning is posed as an optimization problem, using a maximum weighted spanning tree algorithm (where the weights are given by the CMI), a variant of the Chow Liu algorithm [44].

**Algorithm 2** Adaptation of the Chow Liu algorithm of [44]

- 1: Given  $\theta_1, \theta_2, \dots, \theta_n$  wrapped Cauchy variables, estimate the bivariate joint density function  $f(\theta_i, \theta_j|c)$  for all pairs of variables, and the marginals  $f(\theta_i|c)$ , for each  $c \in \Lambda(C)$ ,  $i, j = 1, \dots, n$
- 2: Using these, compute all conditional CMI( $\theta_i, \theta_j|C$ ) values, (i.e., the  $n(n - 1)/2$  edge weights) and order them
- 3: Assign the largest two edges to the undirected tree to be represented
- 4: Examine the next-largest edge, and add it to the tree unless it forms a loop, in which case discard it and examine the next largest edge
- 5: Repeat Step 4 until  $n - 1$  edges have been selected (and the spanning undirected tree is finished)

For Step 1 in Algorithm 2, the estimate of the bivariate and marginal densities are performed for each  $c$  using the methods explained in Section 2. Like the traditional mutual information measure for linear variables, the  $CMI(\theta_i, \theta_j|C)$  denotes the entropy reduction of  $\theta_i$  ( $\theta_j$ ) when the value of  $\theta_j$  ( $\theta_i$ ) is known given  $C$ , and represents the weight that links  $\theta_i$  and  $\theta_j$ .

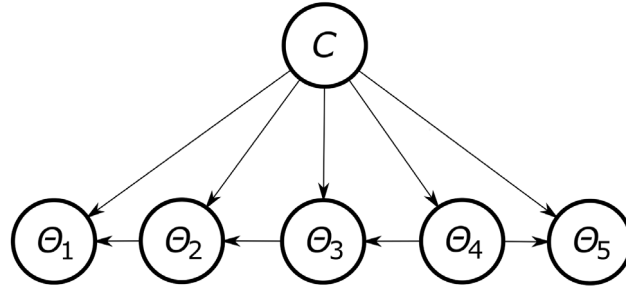


Fig. 4. wCTAN structure with five nodes, from which  $p(c|\theta) \propto p(c)f_{\Theta_1|c,\theta_2}(\theta_1|c,\theta_2)f_{\Theta_2|c,\theta_3}(\theta_2|c,\theta_3)f_{\Theta_3|c,\theta_4}(\theta_3|c,\theta_4)f_{\Theta_4|c}(\theta_4|c)f_{\Theta_5|c,\theta_4}(\theta_5|c,\theta_4)$ . The associated tree-structured Bayesian network has  $\Theta_4$  as its root node.

Once we have learned the undirected structure, a root node must be selected in order to determine the root of the tree by following the structure learned by Algorithm 2. Depending on the selected root node and given the undirected tree structure with  $n$  nodes, there are  $n$  possible resulting directed trees.

- Step 2: We add a class node  $C$  to the network structure. We connect this class node to every other node with an arc from  $C$  (Fig. 4).
- Step 3: Finally, we complete the classification model with the estimation of the parameters for each node given its parent node(s).

Therefore the conditional probability of  $C$  given the predictors is

$$p(C = c|\Theta = \theta) \propto p(C = c)f_{\Theta_{root}|C=c}(\theta_{root}|c) \prod_{i=1, i \neq root}^n f_{\Theta_i|C=c, Pa_{\Theta_i}}(\theta_i|c, Pa_{\Theta_i}),$$

where  $Pa_{\Theta_i}$  is the wrapped Cauchy parent of variable  $\Theta_i$  and  $\Theta_{root}$  is the root node of the tree.

Similar to the approach used in the rest of the models presented in this paper, the maximum a posteriori decision rule is used to determine the predicted class  $c^*$

$$c^* = \arg \max_{c \in \mathcal{A}(C)} p(C = c|\Theta = \theta).$$

#### 4. Experimental results

In this section, we report experiments carried out to show the behavior of each proposed classification model in Section 3. We include the comparison among the four circular classifiers and also with some of the best-known classification algorithms for linear data, such as decision tree (DTree), random forest (Rfor), multinomial logistic regression (MLG), support vector machine (SVM), simple neural network (Nnet) and the Gaussian tree-augmented naive Bayes classifier (GTAN) for continuous data, with the structure learned with the algorithm in [45] where predictor variables given the class value are assumed to follow Gaussian distributions.

The experiments were run using R software [46]. To generate the artificial datasets to test the models, we used the “Circular” R package for the simulation of circular data, and to implement the structure of the wCTAN classifier, we have adapted the “bnclassify” R package [47]. Simulating data that follows wrapped Cauchy distributions is easy and computationally very fast. Given the parameters, the “Circular” R package simulates wrapped Cauchy data by wrapping the simulation of a Cauchy distribution whose location parameter is the same as the wrapped Cauchy location parameter and the scale parameter is the negative logarithm of the wrapped Cauchy concentration parameter. If the wrapped Cauchy concentration parameter is equal to 1, then the value of the simulation will be the location parameter, whereas if the concentration parameter is equal to 0, the simulation is performed from a Uniform distribution in  $[0, 2\pi)$ .

In order to test the algorithms, we enforced dependence between nodes giving values of  $|\rho|$  in  $[0.5, 1)$ . The remaining parameters were assigned randomly to each node with  $-\pi < \mu < \pi$  and  $0 < \varepsilon < 1$ . For each classifier, we simulated 10 datasets each with 1000, 200 and 50 instances and 3, 5, 10, 20, 30, 45, 65, and 100 wrapped Cauchy predictor variables and a discrete class variable with 3, 6, 10, 15 and 20 different labels, so we simulated 1200 different datasets for each type of classifier. A 10-fold cross-validation was used to estimate the classification accuracy. Results for Bayesian network classifiers are shown in Table 1, while results for traditional linear classification algorithms are shown in Table 2.

We also applied the non-parametric Friedman test to detect statistically significant differences among our classification models as a whole set [48]. When the null hypothesis was rejected, we proceeded with post-hoc tests. We chose the Nemenyi test [49], as suggested by [50]. The significance level  $\alpha$  for all tests was 0.05.

Since multiple classifiers are compared, it is useful to represent the results of the post-hoc tests visually. The graph proposed by Demšar [50] is a simple diagram to easily represent these results. The top line is the axis on which we plot the average Friedman test ranks of the classifiers. The lowest (best) ranks are to the right, and we therefore consider the classifiers to the right as better. For the comparison results of all classifiers against each other, those that are not significantly different ( $p$ -value  $\geq 0.05$  in the Nemenyi post-hoc test) are connected.







**Table 3**

Mean  $\pm$  standard deviation accuracy of wCNB, wCsNB, wCsmNB, wCTAN, GTAN, DTree, Rfor, MLG, SVM and Nnet classifiers for different number of variables. Results are averaged from the classification performance from Tables 1 and 2 with 3, 6, 10, 15 and 20 different labels with 1000 instances. Bolded results are best performing classifiers.

		Classifiers				
		wCNB	wCsNB	wCsmNB	wCTAN	GTAN
No. of variables	3	0.735 $\pm$ 0.108	<b>0.755 <math>\pm</math> 0.097</b>	0.754 $\pm$ 0.109	0.743 $\pm$ 0.108	0.352 $\pm$ 0.171
	5	0.866 $\pm$ 0.069	0.877 $\pm$ 0.069	<b>0.879 <math>\pm</math> 0.066</b>	0.876 $\pm$ 0.074	0.409 $\pm$ 0.174
	10	0.976 $\pm$ 0.015	0.948 $\pm$ 0.021	<b>0.983 <math>\pm</math> 0.011</b>	0.974 $\pm$ 0.018	0.491 $\pm$ 0.159
	20	<b>0.998 <math>\pm</math> 0.001</b>	0.970 $\pm$ 0.017	<b>0.998 <math>\pm</math> 0.001</b>	<b>0.998 <math>\pm</math> 0.001</b>	0.610 $\pm$ 0.122
	30	0.998 $\pm$ 0.001	0.976 $\pm$ 0.015	<b>0.999 <math>\pm</math> 0.001</b>	<b>0.999 <math>\pm</math> 0.001</b>	0.674 $\pm$ 0.116
	45	<b>0.999 <math>\pm</math> 0.001</b>	0.980 $\pm$ 0.013	<b>0.999 <math>\pm</math> 0.001</b>	<b>0.999 <math>\pm</math> 0.001</b>	0.790 $\pm$ 0.091
	65	<b>0.999 <math>\pm</math> 0.001</b>	0.984 $\pm$ 0.013	<b>0.999 <math>\pm</math> 0.001</b>	<b>0.999 <math>\pm</math> 0.001</b>	0.824 $\pm$ 0.076
	100	<b>0.999 <math>\pm</math> 0.001</b>	0.989 $\pm$ 0.012	<b>0.999 <math>\pm</math> 0.001</b>	<b>0.999 <math>\pm</math> 0.001</b>	0.873 $\pm$ 0.062
		DTree	Rfor	MLG	SVM	Nnet
No. of variables	3	0.182 $\pm$ 0.021	0.678 $\pm$ 0.062	0.482 $\pm$ 0.070	0.627 $\pm$ 0.048	0.656 $\pm$ 0.055
	5	0.193 $\pm$ 0.017	0.803 $\pm$ 0.048	0.579 $\pm$ 0.053	0.718 $\pm$ 0.043	0.690 $\pm$ 0.051
	10	0.207 $\pm$ 0.012	0.921 $\pm$ 0.019	0.706 $\pm$ 0.046	0.818 $\pm$ 0.034	0.661 $\pm$ 0.044
	20	0.212 $\pm$ 0.013	0.984 $\pm$ 0.008	0.807 $\pm$ 0.032	0.905 $\pm$ 0.024	0.671 $\pm$ 0.065
	30	0.210 $\pm$ 0.009	0.996 $\pm$ 0.003	0.834 $\pm$ 0.030	0.947 $\pm$ 0.014	0.723 $\pm$ 0.055
	45	0.218 $\pm$ 0.010	<b>0.999 <math>\pm</math> 0.001</b>	0.880 $\pm$ 0.048	0.991 $\pm$ 0.005	0.749 $\pm$ 0.074
	65	0.217 $\pm$ 0.007	<b>0.999 <math>\pm</math> 0.001</b>	0.893 $\pm$ 0.049	0.996 $\pm$ 0.002	0.716 $\pm$ 0.129
	100	0.218 $\pm$ 0.006	<b>0.999 <math>\pm</math> 0.001</b>	0.941 $\pm$ 0.073	<b>0.999 <math>\pm</math> 0.001</b>	0.758 $\pm$ 0.135

#### 4.1. Comparison of classification models

In this section, we compare the performance of the wCNB, wCsNB, wCsmNB and wCTAN models, as well as the DTree, Rfor, MLG, SVM, Nnet and the GTAN algorithms, which ignores the circular nature of the data. We analyze the results of the simulation with 1000 instances. Additionally, we analyze the Bayesian network classifiers performance for 50 and 200 instances.

Table 3 shows the mean  $\pm$  standard deviation accuracy for each classifier for different number of variables. Each mean  $\pm$  standard deviation accuracy values was obtained from the results of 50 independent 10-fold cross-validation procedures varying the number of labels of the class variable (3, 6, 10, 15 and 20 different labels) with 1000 instances.

The statistical analysis after Friedman test rejection ( $p$ -value = 0.000000005) reveals (Fig. 5A) that, varying the number of variables, the best classifiers are wCsmNB, wCTAN, wCNB, Rfor and wCsNB with no statistically significant differences among them, whereas the DTree, GTAN, Nnet and MLG classifiers are the worst, presenting significant differences with respect to the rest of the classifiers but for the SVM (which does not present statistical differences with the Nnet and MLG) and demonstrating that treating circular data as linear-continuous is not effective. The SVM also presents statistical differences when compared with the wCsmNB, wCTAN and wCNB classifiers, which outperforms the SVM results. Nevertheless, there were no significant differences between the SVM and the remaining circular classifier (i.e., wCsNB) and the Rfor classifier.

Performing the same statistical analysis among Bayesian network classifiers with 50 and 200 instances yields similar results. The Friedman test null hypothesis that there is no significant difference is rejected for both ( $p$ -value = 0.00021 and  $p$ -value = 0.00004, respectively). The post-hoc analysis displays quite similar results to the 1000 instances one; in both cases, there are no statistically significant differences among the wCsmNB, wCTAN and wCNB classifiers, which are the best. Nevertheless, for 50 and 200 instances, there are no significant differences among GTAN and wCsNB classifiers. Furthermore, there are significant differences between the wCsNB and the wCsmNB classifier for the analysis with 50 instances, whereas for 200 instances, statistical differences were seen between the wCsNB classifier and both the wCsmNB and the wCTAN.

We also calculated the mean accuracy for each classifier for different number of labels in the class variable (see Table 4). Each mean accuracy value was obtained from the results of 60 independent 10-fold cross-validation procedures varying the number of variables to be used: 3, 5, 10, 20, 30 and 45 with 1000 instances. We do not include the results of the experiments with more than 45 variables due to the high mean accuracy values obtained in most of the classifiers from Tables 1 and 2, which would bias the results.

Since the Friedman test null hypothesis that there is no significant difference was rejected ( $p$ -value = 0.0000011), we performed the corresponding Nemenyi post-hoc analysis. Statistical test results (Fig. 5B) reveal that based on changing the number of labels, the best classifiers are the circular Bayesian network classifiers (i.e., wCsmNB, wCTAN, wCNB and wCsNB), with no statistically significant differences between them. Again, DTree, GTAN, Nnet and MLG are the worst, with no significant differences among them. These classifiers shows significant differences with wCsmNB, wCTAN, wCNB, Rfor and wCsNB, whereas SVM only presents significant differences with the wCsmNB, wCTAN, wCNB, DTree and GTAN classifiers.

The analysis for Bayesian network classifiers with 50 and 200 instances again yielded quite similar results to those obtained for 1000 instances. After Friedman test rejections ( $p$ -value = 0.0325 for 50 instances, and  $p$ -value = 0.00066 for 200 instances), post-hoc tests for 200 instances reveal the same statistically significant differences as for 1000 instances, where there is no statistical differences among the wCsmNB, wCTAN and wCNB classifiers, which are the best. For 50 instances, wCsmNB, wCTAN and wCNB are also the best classifiers together with the wCsNB, with no statistically significant differences among them. Likewise, for 1000 instances, GTAN is the worst for the analysis with 50 instances as well as the 200 instances, with no significant differences with the wCsNB classifier.

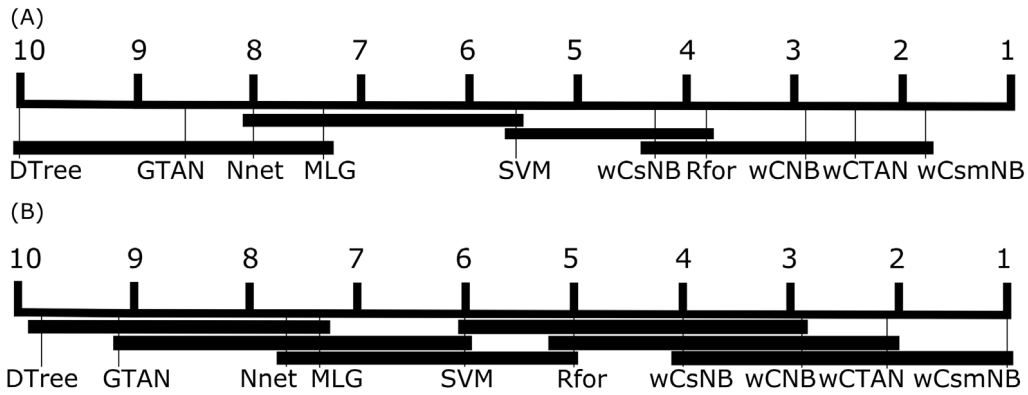


Fig. 5. Demšar diagrams presenting the statistical comparison among wCNB, wCsNB, wCsmNB, wCTAN, GTAN, DTree, Rfor, MLG, SVM and Nnet classification models for synthetic datasets with 1000 instances. Those classifiers that are not connected show differences that are statistically significant ( $p$ -value < 0.05). The lowest rank classifiers are to the right side of the graph (i.e., they can be considered the best). (A) Comparison varying the number of labels. (B) Comparison varying the number of variables.

Table 4

Mean  $\pm$  standard deviation accuracy of wCNB, wCsNB, wCsmNB, wCTAN, GTAN, DTree, Rfor, MLG, SVM and Nnet classifiers for different number of labels. Results are averaged from the classification performance with 3, 5, 10, 20, 30 and 45 different variables with 1000 instances. Bolded results are best performing classifiers.

		Classifiers				
		wCNB	wCsNB	wCsmNB	wCTAN	GTAN
No. of labels	3	0.973 $\pm$ 0.042	0.970 $\pm$ 0.034	<b>0.979 <math>\pm</math> 0.035</b>	0.976 $\pm$ 0.043	0.709 $\pm$ 0.061
	6	0.940 $\pm$ 0.083	0.939 $\pm$ 0.054	<b>0.952 <math>\pm</math> 0.069</b>	0.946 $\pm$ 0.074	0.526 $\pm$ 0.097
	10	0.910 $\pm$ 0.118	0.901 $\pm$ 0.089	<b>0.921 <math>\pm</math> 0.106</b>	0.918 $\pm$ 0.106	0.384 $\pm$ 0.124
	15	0.884 $\pm$ 0.146	0.862 $\pm$ 0.117	<b>0.906 <math>\pm</math> 0.128</b>	0.884 $\pm$ 0.148	0.311 $\pm$ 0.096
	20	0.858 $\pm$ 0.167	0.839 $\pm$ 0.134	<b>0.872 <math>\pm</math> 0.171</b>	0.860 $\pm$ 0.164	0.253 $\pm$ 0.098
		DTree	Rfor	MLG	SVM	Nnet
No. of labels	3	0.893 $\pm$ 0.051	0.968 $\pm$ 0.017	0.898 $\pm$ 0.041	0.954 $\pm$ 0.019	0.932 $\pm$ 0.026
	6	0.122 $\pm$ 0.014	0.938 $\pm$ 0.022	0.790 $\pm$ 0.049	0.887 $\pm$ 0.026	0.831 $\pm$ 0.053
	10	0.001 $\pm$ 0.001	0.898 $\pm$ 0.026	0.697 $\pm$ 0.045	0.831 $\pm$ 0.029	0.691 $\pm$ 0.071
	15	0.001 $\pm$ 0.001	0.860 $\pm$ 0.025	0.633 $\pm$ 0.045	0.777 $\pm$ 0.032	0.549 $\pm$ 0.067
	20	0.001 $\pm$ 0.001	0.821 $\pm$ 0.030	0.555 $\pm$ 0.052	0.723 $\pm$ 0.033	0.456 $\pm$ 0.070

### 5. Real data example

We applied our classifiers to a dataset of 3027 combinations of dendritic bifurcation angles coming from the basal arbors of 288 3D pyramidal neurons in layers II, III, IV, Va, Vb and VI (48 neurons per layer) of the 14-day-old (P14) rat hind limb somatosensory (S1HL) neocortex, recently published in [14] (Fig. 6).

We used the Bayesian network classification models presented in Section 3 and wrapped Cauchy distributions to model the bifurcation angles produced by the splitting of the dendritic segments of basal dendritic trees. The dendritic bifurcation angles are an important part of the geometry of pyramidal cell arbors. Since it is thought that these angles determine the space to be filled by the dendritic wiring, understanding and modeling them are crucial for advances in neuroscience to replicate brain functioning and structure in order to make further on how the brain processes information. This is important not only to understand it biologically (i.e., thoughts, emotions, feelings) but also technological, making essential contributions to new computing. Moreover, brain knowledge is basic for treating brain diseases such as Parkinson or Alzheimer.

Predicting which layer a neuron belongs to is an important task to help understand any neural circuit, and it represents part of the picture regarding the identification and characterization of all its components. To the best of our knowledge, there is no any supervised classification model that predicts the layer using circular predictive variables. Thus, we developed a classification model to predict which layer a given neuron belongs to, i.e.,  $A(C) = (II, III, IV, Va, Vb, VI)$ .

Following the notation used in [14],  $\theta_1$  will correspond to the first bifurcation angle (Order 1) generated for the first split of the dendritic segments starting from the soma. The second angle generated by the next consecutive splits will be represented as variable  $\theta_2$  (Order 2), etc. (Fig. 7). Angles of orders higher than six which were relatively scarce were not included in the model. For each set of angles of the same order, a wrapped Cauchy distribution was fitted (Table 5). We performed a goodness-of-fit test by transformation on the circle of the variables into circular uniform variables via  $2\pi F(\theta_1), \dots, 2\pi F(\theta_6)$ , where  $F$  is the cumulative distribution function, and applied Kuiper’s test [51] for circular uniformity with a significance level of  $\alpha = 0.05$ .

Note that in Table 5 the circular mean tends to decrease as the order increases. A neuroscientific explanation for this behavior relates to the fact that it is the first bifurcation orders that determine the volume of space to be filled by the dendritic trees [14]. This

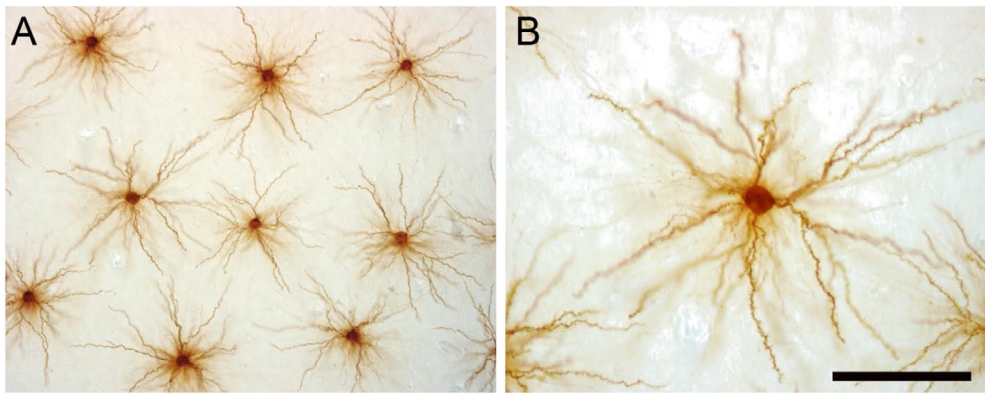


Fig. 6. (A) Low-power photomicrograph showing injected neurons in layers III from the S1HL region of P14 rats, as seen in the plane of section parallel to the cortical surface. (B) Higher magnification photomicrograph showing an example of a pyramidal cell basal dendritic arbor. Scale bar (in B) = 200  $\mu\text{m}$  in A; 90  $\mu\text{m}$  in B.

Source: Adapted from [14].

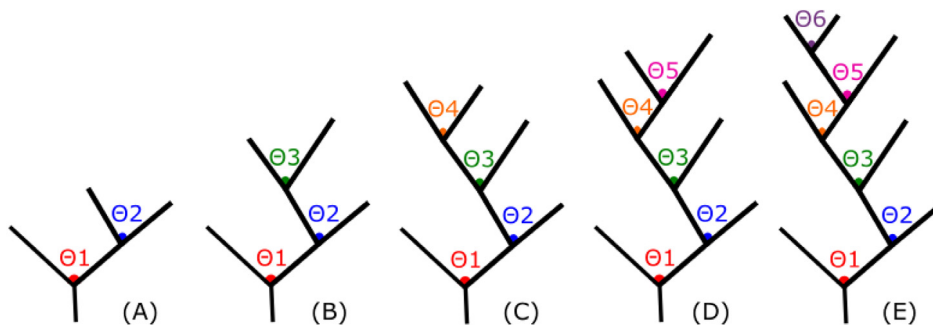


Fig. 7. Angles of different branch orders (from 1 to 6) measured between sibling segments in a dendritic arbor. The dendritic arbor has a maximum branching order of (A) 2 (B) 3 (C) 4 (D) 5 (E) 6.

**Table 5**  
Characteristics of the six different branching orders shown in Fig. 7.

Bifurcation order	Variable	No. of angles	$\hat{\mu}$ (in radians)	$\hat{\epsilon}$
1	$\theta_1$	1607	1.02	0.90
2	$\theta_2$	2072	0.90	0.91
3	$\theta_3$	1773	0.82	0.92
4	$\theta_4$	998	0.78	0.92
5	$\theta_5$	382	0.77	0.92
6	$\theta_6$	106	0.81	0.92

regulates the dendritic branching development rules that seem to determine the synaptic connectivity of pyramidal neurons. We also observe that the concentration values are high (around 0.91) and quite similar in every bifurcation order. This fact demonstrates that the dendritic structure (in terms of bifurcation angles) is determined by the location parameter.

Since not all dendritic arbors present angles of all orders, one classifier for the whole dataset is not suitable. Therefore, for each classification model proposed in this paper, we created a battery of five classifiers depending on the maximum bifurcation order of the arbor, when this is higher than 1 (Fig. 8). Before predicting class  $c^*$ , we have to check the maximum bifurcation order of the instance to be classified. For the wCTAN and GTAN structures (which require a root node in addition to the class node) we select as root node  $\theta_2$  for every classifier of the battery. We performed 10 fold cross-validation procedures in order to obtain the mean classification accuracy values for each classifier and maximum bifurcation order (Table 6).

We observe in Table 6 that the wCsNB classifier leads to the best results for arbors with a maximum branching order of  $\theta_2$  and  $\theta_3$ . Furthermore, for arbors with a maximum branching order of  $\theta_4$ ,  $\theta_5$  or  $\theta_6$ , the wCsmNB seems to perform best in terms of classification accuracy. The wCTAN and wCNB classifiers also report acceptable values in comparison with the highest ones for each maximum bifurcation arbor, although the wCsNB or wCsmNB classification models always perform better for this neuronal dataset. Comparing the accuracy results with the random label assignment (i.e.,  $1/6 = 0.16$ ), we observe that all of these results are over 0.16. In addition, for every case, the GTAN classifier exhibits the lowest accuracy values, below 0.16 except for  $\theta_4$ . This classifier was especially inaccurate for arbors that had maximum branching order of 6; the mean accuracy value was 0.047 for such cases.

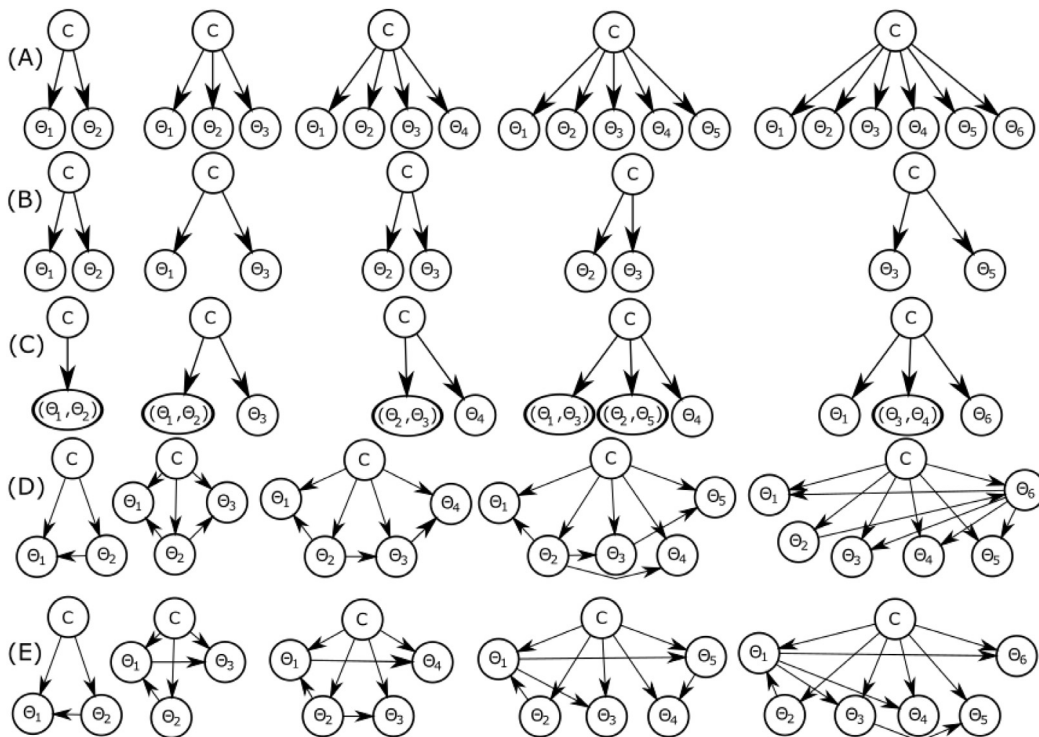


Fig. 8. Bayesian network classifier structures associated with the battery of classifiers depending on the maximum bifurcation order, for each type of classification algorithm: (A) wCNB, (B) wCsNB, (C) wCsmNB, (D) wCTAN, (E) GTAN.

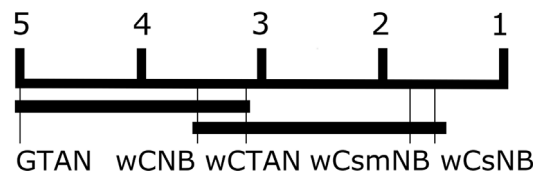


Fig. 9. Demšar diagram for the comparison of wCNB, wCsNB, wCsmNB, wCTAN and GTAN classification models using Friedman test and Nemenyi post-hoc test.

Table 6

Mean ± standard deviation of layers II, III, IV, Va, Vb and VI classification accuracy results of the battery of classifiers for each type of classifier applied over the dataset of dendritic bifurcation angles coming from the basal arbors of 288 3D pyramidal neurons of P14 rat SIHL neocortex. Bolded results are best performing classifiers.

	Classifiers					
	wCNB	wCsNB	wCsmNB	wCTAN	GTAN	
Max. bifurc. angle	2	0.182 ± 0.034	<b>0.184 ± 0.055</b>	0.182 ± 0.039	0.182 ± 0.034	0.158 ± 0.038
	3	0.191 ± 0.035	<b>0.219 ± 0.035</b>	0.203 ± 0.057	0.205 ± 0.057	0.113 ± 0.030
	4	0.222 ± 0.016	0.224 ± 0.034	<b>0.239 ± 0.057</b>	0.222 ± 0.046	0.212 ± 0.081
	5	0.196 ± 0.091	0.235 ± 0.127	<b>0.239 ± 0.063</b>	0.189 ± 0.055	0.128 ± 0.131
	6	0.220 ± 0.113	0.270 ± 0.094	<b>0.290 ± 0.137</b>	0.240 ± 0.126	0.047 ± 0.069

We applied the Friedman non-parametric test to detect statistically significant differences in the results provided by our algorithms. Since the null hypothesis that there is no significant difference was rejected ( $p$ -value = 0.004), we used Nemenyi post-hoc test to determine which pairwise of algorithms was the cause of the Friedman test rejection. In Fig. 9, the statistically significant differences between our classifiers are represented as a Demšar diagram. We noted that there are no statistically significant differences between our classification algorithms except for two cases; between the wCTAN and wCsNB and between GTAN and the wCsmNB.

Therefore, we can conclude that (i) apart from the difficulty identifying the layer a case belongs to, it seems reasonable to use any of our four proposed circular classifiers for this neuronal dataset, since there are no any statistically significant differences between them and (ii) GTAN is never recommended.

## 6. Conclusions and future work

Introducing the first set of supervised Bayesian classification models capable of dealing with circular wrapped Cauchy predictive variables was the main objective of this paper. We have presented four models and their algorithms, designed to perform classification. We demonstrated using synthetic data that these models could perform classification accurately given circular datasets. We also provided evidence of the improvement of the circular classifiers over linear classifiers for datasets of circular nature that follow wrapped Cauchy distributions.

We performed statistical comparisons among the classifiers using synthetic data with 50, 200 and 1000 instances. Based on the results, we realized that the wCsmNB, the wCTAN and the wCNB are the best classification models for circular data that follows wrapped Cauchy distributions, with no statistically significant differences among them. The linear classifier never outperformed any of the wrapped Cauchy classifiers.

For each of our new proposals, we evaluated a battery of classifiers using a real-world neuroscience dataset, in order to predict the layer that an instance belongs to. Results revealed that all of our four classification models are suitable. Performing Friedman test and its corresponding Nemenyi post-hoc test after rejection, we realized that there are no any statistically significant differences between wCNB, wCsNB, wCsmNB and wCTAN for this dataset. Wrapped Cauchy classifiers always outperformed their linear (Gaussian) counterparts.

The models shown in this paper are limited to no more than bivariate relationships. In future work, we intend to develop multivariate models in order to extend the Bayesian network classifiers for circular data to other more-sophisticated Bayesian network models (like k-dependence Bayesian network classifiers) capable of representing and taking into account multivariate relationships between circular variables — a difficult task due to the non-closed nature of the circular families that are known to date.

## Acknowledgments

This work has been partially supported by the Spanish Ministry of Economy and Competitiveness through the TIN2016-79684-P and Cajal Blue Brain (C080020-09) projects. This project has received funding from the European Union's Horizon 2020 Framework Programme for Research and Innovation under Specific Grant Agreement No. 785907 (HBP SGA2). I.L. is supported by the Spanish Ministry of Education, Culture and Sport Fellowship (FPU13/01941). The authors thankfully acknowledge the Cortical Circuits Laboratory (CSIC-UPM) for providing the neuron dataset.

## References

- [1] K. Schmidt-Koenig, On the role of the loft, the distance and the site of release in pigeon homing, *Biol. Bull.* 125 (1963) 154–164.
- [2] L. Morellato, L. Alberti, I. Hudson, Application of circular statistics in plant phenology: a case studies approach, in: *Phenological Research: Methods for Environmental and Climate Change Analysis*, Springer, 2010, pp. 339–359.
- [3] N.I. Fisher, *Statistical Analysis of Circular Data*, Cambridge University Press, 1995.
- [4] F. Fisher, Dispersion on a sphere, *Proc. Roy. Soc. Lond. A Math. Phys. Eng. Sci.* 217 (1130) (1953) 295–305.
- [5] J. Graham, The stability and significance of magnetism in sedimentary rocks, *J. Geophys. Res.* (1949) 131–167.
- [6] J. Gill, D. Hangartner, Circular data in political science and how to handle it, *Political Anal.* 18 (3) (2010) 316–336.
- [7] A. Banerjee, I. Dhillon, J. Ghosh, S. Sra, Clustering on the unit hypersphere using von Mises-Fisher distributions, *J. Mach. Learn. Res.* 6 (2005) 1345–1382.
- [8] A. Barros, J. Pereira, U. Lund, Identifying geographical patterns of wildfire orientation: A watershed-based analysis, *Forest Ecol. Manag.* 264 (2012) 98–107.
- [9] J. Bowers, I. Morton, G. Mould, Directional statistics of the wind and waves, *Appl. Ocean Res.* 22 (1) (2000) 13–30.
- [10] J.A. Carta, C. Bueno, P. Ramirez, Statistical modelling of directional wind speeds using mixtures of von Mises distributions: Case study, *Energy Convers. Manage.* 49 (5) (2008) 897–907.
- [11] E. Batschelet, *Circular Statistics in Biology*, Academic Press, 1981.
- [12] T. Harder, W. Boomsma, M. Paluszewski, J. Frellsen, K.E. Johansson, T. Hamelryck, Beyond rotamers: a generative, probabilistic model of side chains in proteins, *BMC Bioinformatics* 11 (1) (2010) 306.
- [13] C. Bielza, P. Larrañaga, Discrete Bayesian network classifiers: A survey, *ACM Comput. Surv.* 47 (1) (2014) 5.
- [14] I. Leguey, C. Bielza, P. Larrañaga, A. Kastanauskaitė, C. Rojo, R. Benavides-Piccione, J. DeFelipe, Dendritic branching angles of pyramidal cells across layers of the juvenile rat somatosensory cortex, *J. Comp. Neurol.* 524 (13) (2016) 2567–2576.
- [15] S.R. Jammalamadaka, A. SenGupta, *Topics in Circular Statistics*, World Scientific, 2001.
- [16] K.V. Mardia, P.E. Jupp, *Directional Statistics*, John Wiley & Sons, 2009.
- [17] D. Koller, N. Friedman, *Probabilistic Graphical Models: Principles and Techniques*, The MIT Press, 2009.
- [18] J. Pearl, Probabilistic reasoning in intelligent systems: Networks of plausible reasoning, *Morgan Kaufmann* 23 (1988) 33–34.
- [19] R.O. Duda, P.E. Hart, D.G. Stork, *Pattern Classification*. 2nd edition, John Wiley & Sons.
- [20] J.E. Morris, P. Laycock, Discriminant analysis of directional data, *Biometrika* 61 (2) (1974) 335–341.
- [21] M. Romanazzi, Discriminant analysis with high dimensional von Mises-Fisher distributions, in: *8th Annual International Conference on Statistics*, 2014, pp. 1–16.
- [22] S. El Khattabi, F. Streit, Identification analysis in directional statistics, *Comput. Statist. Data Anal.* 23 (1) (1996) 45–63.
- [23] A. Figueiredo, P. Gomes, Discriminant analysis based on the Watson distribution defined on the hypersphere, *Statistics* 40 (5) (2006) 435–445.
- [24] A. SenGupta, S. Roy, A simple classification rule for directional data, in: *Advances in Ranking and Selection, Multiple Comparisons, and Reliability*, Springer, 2005, pp. 81–90.
- [25] A. SenGupta, F.I. Ugwuowo, A classification method for directional data with application to the human skull, *Comm. Statist. Theory Methods* 40 (3) (2011) 457–466.
- [26] M.J. Kirby, R. Miranda, Circular nodes in neural networks, *Neural Comput.* 8 (2) (1996) 390–402.
- [27] K. Fernandes, J.S. Cardoso, Discriminative directional classifiers, *Neurocomputing* 207 (2016) 141–149.
- [28] P.L. López-Cruz, C. Bielza, P. Larrañaga, Directional naive Bayes classifiers, *Pattern Anal. Appl.* 18 (2) (2015) 225–246.
- [29] R. von Mises, Über die “Ganzzahligkeit” der Atomgewichte und verwandte Fragen, *Z. Phys.* 19 (1918) 490–500.

- [30] K.V. Mardia, Statistics of directional data, *J. R. Stat. Soc. Ser. B Stat. Methodol.* (1975) 349–393.
- [31] K.V. Mardia, G. Hughes, C.C. Taylor, H. Singh, A multivariate von Mises distribution with applications to bioinformatics, *Canad. J. Statist.* 36 (1) (2008) 99–109.
- [32] P. Lévy, L'addition des variables aléatoires définies sur une circonférence, *Bull. Soc. Math. France* 67 (1939) 1–41.
- [33] A. Wintner, On the shape of the angular case of cauchy's distribution curves, *The Annals of Mathematical Statistics* 18 (4) (1947) 589–593.
- [34] P. McCullagh, Möbius transformation and Cauchy parameter estimation, *Ann. Statist.* 24 (2) (1996) 787–808.
- [35] S. Kato, A. Pewsey, A Möbius transformation-induced distribution on the torus, *Biometrika* 102 (2) (2015) 359–370.
- [36] I. Leguey, C. Bielza, P. Larrañaga, Tree-structured Bayesian networks for wrapped Cauchy directional distributions, in: *Advances in Artificial Intelligence*, Vol. 9868, Springer, 2016, pp. 207–216.
- [37] K. Bowman, L. Shenton, Methods of moments, *Encyclopedia Stat. Sci.* 5 (1985) 467–473.
- [38] P. Langley, S. Sage, Induction of selective Bayesian classifiers, in: *Proceedings of the Tenth International Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, 1994, pp. 399–406.
- [39] Y. Saeys, I. Inza, P. Larrañaga, A review of feature selection techniques in bioinformatics, *Bioinformatics* 23 (19) (2007) 2507–2517.
- [40] R. Blanco, I. Inza, M. Merino, J. Quiroga, P. Larrañaga, Feature selection in Bayesian classifiers for the prognosis of survival of cirrhotic patients treated with TIPS, *J. Biomed. Inform.* 38 (5) (2005) 376–388.
- [41] I. Kononenko, Semi-naive Bayesian classifier, in: *European Working Session on Learning*, Springer, 1991, pp. 206–219.
- [42] M.J. Pazzani, Constructive induction of Cartesian product attributes, in: *Feature Extraction, Construction and Selection*, Springer, 1998, pp. 341–354.
- [43] N. Friedman, D. Geiger, M. Goldszmidt, Bayesian network classifiers, *Mach. Learn.* 29 (1997) 131–163.
- [44] C. Chow, C. Liu, Approximating discrete probability distributions with dependence trees, *IEEE Trans. Inform. Theory* 14 (3) (1968) 462–467.
- [45] D. Geiger, D. Heckerman, Learning Gaussian networks, in: *Proceedings of the Tenth International Conference on Uncertainty in Artificial Intelligence*, Morgan Kaufmann, 1994, pp. 235–243.
- [46] R Development Core Team, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing (2008).
- [47] B. Mihaljevic, C. Bielza, P. and Larrañaga, bnclassify: Learning discrete Bayesian network classifiers from data, R package version 0.3.2 (2015). URL <https://cran.r-project.org/src/contrib/Archive/bnclassify/>.
- [48] M. Friedman, The use of ranks to avoid the assumption of normality implicit in the analysis of variance, *J. Amer. Statist. Assoc.* 32 (200) (1937) 675–701.
- [49] P. Nemenyi, Distribution-free multiple comparisons, *Biometrics* 18 (2) (1962) 263.
- [50] J. Demšar, Statistical comparisons of classifiers over multiple data sets, *J. Mach. Learn. Res.* 7 (2006) 1–30.
- [51] N.H. Kuiper, Tests concerning random points on a circle, *Indag. Math. (N.S.)* 63 (1960) 38–47.

**Ignacio Leguey** is currently a Visiting Professor with the Department of Economía Financiera y Contabilidad e Idioma Moderno, Universidad Rey Juan Carlos de Madrid, Spain. His research interests include directional statistics, supervised classification, particularly Bayesian network classifiers, probabilistic graphical models, and real applications, like meteorology and neuroscience.

**Concha Bielza** is currently a Full Professor with the Department of Artificial Intelligence, Universidad Politécnica de Madrid, Spain. Her research interests include probabilistic graphical models, decision analysis, metaheuristics for optimization, classification models, and real applications, like bioinformatics, neuroscience and industry. She was awarded the 2014 UPM Research Prize.

**Pedro Larrañaga** is currently a Full Professor with the Department of Artificial Intelligence, Universidad Politécnica de Madrid, Spain. His research interests include probabilistic graphical models, data mining, classification models, and real applications, like biomedicine and neuroscience. He is ECCAI fellow since 2012 and he has been awarded the 2013 Spanish National Prize in Computer Science.