

# Frobenius Norm Regularization for the Multivariate Von Mises Distribution

Luis Rodriguez-Lujan,<sup>\*</sup> Pedro Larrañaga,<sup>†</sup> Concha Bielza<sup>‡</sup>  
*Computational Intelligence Group, Departamento de Inteligencia Artificial,  
Universidad Politécnica de Madrid, Madrid, Spain*

Penalizing the model complexity is necessary to avoid overfitting when the number of data samples is low with respect to the number of model parameters. In this paper, we introduce a penalization term that places an independent prior distribution for each parameter of the multivariate von Mises distribution. We also propose a circular distance that can be used to estimate the Kullback–Leibler divergence between any two circular distributions as goodness-of-fit measure. We compare the resulting regularized von Mises models on synthetic data and real neuroanatomical data to show that the distribution fitted using the penalized estimator generally achieves better results than nonpenalized multivariate von Mises estimator. © 2016 Wiley Periodicals, Inc.

## 1. INTRODUCTION

Directional data appears in many science domains in the form of directions or angles, but also as any kind of periodical data like the hours of the day. Indeed, this periodicity is the main characteristic which differentiates directional data from linear data. Recent examples of directional data in the literature cover a plethora of topics, including wind direction,<sup>1,2</sup> handwriting recognition,<sup>3</sup> people orientation in computer vision,<sup>4</sup> animal orientation,<sup>5,6</sup> marine currents,<sup>6</sup> protein backbone angles,<sup>7–9</sup> and text mining similarity measures<sup>10</sup> among others.

Directional statistics<sup>11,12</sup> provides specific tools for modeling directional data. In the past years, new circular distributions and techniques have emerged in the literature for univariate circular data<sup>13,14</sup> but also for multivariate circular data.<sup>6,9,15,16</sup> The von Mises distribution, the circular analogue of the normal distribution, is still the distribution of choice in the directional statistics field.

The extension of the von Mises distribution to a multivariate distribution<sup>7</sup> presents the problem that no closed formulation is known for the normalization term when the number of variables is greater than two,<sup>17</sup> and therefore it cannot be

<sup>\*</sup>Author to whom all correspondence should be addressed; e-mail: luis.rodriguezl@upm.es.

<sup>†</sup>e-mail: pedro.larranaga@fi.upm.es.

<sup>‡</sup>e-mail: mcbielza@fi.upm.es.

easily fitted nor compared to other distributions. In this article, we reformulate the log pseudo-likelihood expression for the multivariate von Mises distribution in such a way that it becomes easier to compute.

In some application areas such as neuroanatomy, quality data are scarce and the process to obtain new data, a three-dimensional (3D) reconstruction of a neuron, could take up to several days. In this situation when we deal with a small-sample learning problem, penalizing the model complexity is needed to prevent overfitting or, as shown in this paper, to compensate the estimator bias. Although the usual approach is to use a uniform  $L_1$  penalization over the parameters of the model,<sup>8</sup> in the present paper we propose a more general penalization term that could, in some way, be aware of previous knowledge about the structure of dependencies between the variables in the data, for example, taking into account spatial relationships between variables. In the linear case, for the multivariate normal distribution, this structure-aware penalization paradigm has been applied to learn graphical models with hubs<sup>18</sup> or to penalize according to some defined distance between the variables.<sup>19</sup>

The results in this paper extend those in Rodriguez-Lujan et al.<sup>20</sup> The added contributions of the present paper are a redefinition of a penalization term for learning the parameters of the multivariate conditional distributions, a brief proof of consistence of the penalized estimator, and the study of its bias and variance properties through numerical experiments. The application to real-world data in neuroanatomy is extended to include new data sets from different species rather than focus exclusively on human neurons.

This paper is organized as follows: Section 2 reviews the univariate and multivariate von Mises distributions and defines the maximum pseudo-likelihood estimator. In Section 3, we propose a penalization term for the log pseudo-likelihood based on the Frobenius norm, prove its asymptotic convergence, and compare it against the nonpenalized estimator. Then, in Section 4, we compare the von Mises distribution and the Gaussian distribution over real data from human, rat, and mouse neurons using an approximation of the Kullback-Leibler (KL) divergence, introduced for the first time for circular data in this paper, as the assessment metric. We conclude the paper in Section 5 with a final discussion and some proposals for future work.

## 2. THE MULTIVARIATE VON MISES DISTRIBUTION

### 2.1. Definition of Density Functions

The von Mises distribution is one of the most relevant probability distributions in the field of directional statistics, and it is often considered as the normal distribution in the circumference.<sup>12</sup> Contrary to other directional distributions, such as wrapped<sup>11,14</sup> or projected distributions,<sup>15,21</sup> the von Mises distribution is a native directional distribution. Owing to this purely directional nature, the von Mises distribution has better mathematical properties than other nonnative circular distributions like the wrapped-normal distribution. Perhaps, one of the most important properties of the von Mises distribution is that it belongs to the canonical exponential family,

which is obvious when we examine its density function:

$$f_{VM}(\theta; \mu, \kappa) = \frac{1}{2\pi I_0(\kappa)} \exp\{\kappa \cos(\theta - \mu)\}; \quad \theta, \mu \in [0, 2\pi), \kappa > 0 \quad (1)$$

where  $\mu$  stands for the mean angle,  $\kappa$  measures the concentration around the mean, and  $I_0$  is the modified Bessel function of the first kind of order 0. The modified Bessel function of the first kind of order  $n$ , when  $n$  is a integer, can be expressed as the following integral formula:

$$I_n(z) = \frac{1}{\pi} \int_0^\pi e^{z \cos \theta} \cos(n\theta) d\theta \quad (2)$$

The density in Equation (1) is the angular interpretation of the von Mises distribution. It is easy to define the equivalent geometrical formulation using similar parameters over the unit circumference  $\mathcal{S}_1 = \{(x, y) \in \mathbb{R}^2 : x^2 + y^2 = 1\}$ :

$$f_{VM}(\vec{x}; \vec{\mu}, \kappa) = \frac{1}{2\pi I_0(\kappa)} \exp\{\kappa \vec{x} \cdot \vec{\mu}\}; \quad \vec{x}, \vec{\mu} \in \mathcal{S}_1, \kappa > 0 \quad (3)$$

Here  $\cdot$  is the canonical dot product in  $\mathbb{R}^2$  and  $\vec{\mu}$  is the mean direction. From this definition, we can observe that the probability density for each direction (unitary vector) is proportional to its projection onto the real line determined by the mean direction. For example, the von Mises–Fisher distribution,<sup>22</sup> the analogue of the von Mises distribution in higher dimensions, is usually defined from this geometric point of view. In this paper, however, we will follow the angular interpretation of the von Mises distribution in both univariate and multivariate cases.

Continuing the analogy between the normal distribution and the von Mises distribution, we can define a multivariate von Mises distribution (MVM)<sup>7</sup> as the directional equivalent to the multivariate normal distribution. The density function of the  $p$ -variate von Mises distribution is

$$\begin{aligned} f_{MVM}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\Lambda}) &= \frac{1}{C(\boldsymbol{\kappa}, \boldsymbol{\Lambda})} \exp\{\boldsymbol{\kappa} \cdot \cos(\boldsymbol{\theta} - \boldsymbol{\mu}) \\ &\quad + \frac{1}{2} \sin(\boldsymbol{\theta} - \boldsymbol{\mu}) \boldsymbol{\Lambda} \sin(\boldsymbol{\theta} - \boldsymbol{\mu})^T\}; \\ \boldsymbol{\theta}, \boldsymbol{\mu} &\in [0, 2\pi)^p, \boldsymbol{\kappa} = (\kappa_1, \dots, \kappa_p), \kappa_j > 0 \quad \forall j, \boldsymbol{\Lambda} \in S^p(\mathbb{R}) \end{aligned} \quad (4)$$

where  $S^p(\mathbb{R})$  is the set of real symmetric matrices of size  $p$ ,  $\boldsymbol{\mu}$  is a  $p$ -dimensional vector, the multivariate equivalent of the mean angle in Equation (1),  $\boldsymbol{\kappa}$  is the concentration vector,  $\boldsymbol{\Lambda} = (\lambda_{ij})$  is a square symmetric matrix of size  $p$  whose diagonal elements  $\lambda_{ii}$  are zero, and  $C(\boldsymbol{\kappa}, \boldsymbol{\Lambda})$  is the normalization term. Here  $\cos$  and  $\sin$  are applied entrywise to the  $p$ -dimensional vector  $\boldsymbol{\theta} - \boldsymbol{\mu}$ . The matrix  $\boldsymbol{\Lambda}$  can be seen as a dependency matrix, where the element  $\lambda_{ij}$  measures the conditional probabilistic dependency between the  $i$ th and the  $j$ th variables. Actually, if the

distribution is highly concentrated, that is, the fluctuations in each component are sufficiently small, we can express the correlation coefficient between the  $i$ th and the  $j$ th component as a function of  $\lambda_{ij}$ .<sup>7</sup>

$$\rho_{ij} = \frac{\lambda_{ij}}{\sqrt{\kappa_i \kappa_j}}$$

The normalization term  $C(\boldsymbol{\kappa}, \boldsymbol{\Lambda})$  in Equation (4) is the source of many problems that arise with the multivariate von Mises distribution. It does not have<sup>17</sup> a known closed-form formula for  $p > 2$ ; therefore, its exact value has to be computed through numerical approximation, which can be troublesome in high-dimensional settings. However, as we will see in the following subsections, we can overcome some of the difficulties derived from the lack of a closed formula for the normalization term using the univariate conditional distributions of a multivariate von Mises distribution, which are indeed univariate von Mises distributions with known parameters:

$$f_{MVM}(\theta_j | \theta_1, \dots, \theta_{j-1}, \theta_{j+1}, \dots, \theta_p; \boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\Lambda}) = f_{VM}(\theta_j; \mu_{\setminus j}, \kappa_{\setminus j})$$

$$\begin{aligned} \mu_{\setminus j}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\Lambda}) &= \mu_j + \arctan \left( \frac{1}{\kappa_j} \sum_{l \neq j} \lambda_{jl} \sin(\theta_l - \mu_l) \right) \\ \kappa_{\setminus j}(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\Lambda}) &= \sqrt{\kappa_j^2 + \left( \sum_{l \neq j} \lambda_{jl} \sin(\theta_l - \mu_l) \right)^2} \end{aligned} \quad (5)$$

It is also important to note that the multivariate von Mises distribution is symmetrical and rotationally equivariant with respect to the mean  $\boldsymbol{\mu}$ . This last property allows us to assume in the following sections that  $\mu_j = 0, \forall j = 1, \dots, p$  without loss of generality. We will also restrict the values of  $\boldsymbol{\kappa}$  and  $\boldsymbol{\Lambda}$  such that matrix

$$\mathbf{P} = \text{diag}(\kappa_1, \dots, \kappa_p) - \boldsymbol{\Lambda} \quad (6)$$

is positive definite, where  $\text{diag}(\kappa_1, \dots, \kappa_p)$  denotes the square diagonal matrix whose diagonal entries are  $\kappa_1, \dots, \kappa_p$ . The positive-definiteness of  $\mathbf{P}$  is a sufficient condition to ensure that the unique maximum of the multivariate von Mises distribution is attained at  $\boldsymbol{\mu}$ .<sup>23</sup>

## 2.2. Sampling

In this subsection, we present the two methods provided in the literature to generate samples from a multivariate von Mises distribution. Both methods rely on efficient sample generation from a von Mises distribution.<sup>24</sup> The first method

relies on the Gibbs sampling (GS) technique.<sup>7,8</sup> Our implementation of the GS is described in Algorithm 1, where  $k_{thinning}$  is introduced to break the dependence between consecutive draws. In practice, a certain number of the first draws ( $N_{burnin}$ ) is discarded to reduce the dependency of the generated samples on the starting point, so the selected samples are closer to the stationary distribution.

ALGORITHM 1.

Input: Parameters  $\boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\Lambda}$

Output:  $\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(N)}$ . That is  $N$  random samples from  $MVM(\boldsymbol{\mu}, \boldsymbol{\kappa}, \boldsymbol{\Lambda})$

Steps:

1. Initialize  $\boldsymbol{\theta}^{(0)} = (\theta_1^{(0)}, \dots, \theta_p^{(0)})$  to some arbitrary value and set  $i = 1$
2. Repeat  $N$  times:
  - (a) Update  $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}^{(i-1)}$
  - (b) Select one component  $\theta_j^{(i)}$  at random and update it as:
$$\theta_j^{(i)} \sim f_{MVM}(\theta_j^{(i)} | \theta_1^{(i)}, \dots, \theta_{j-1}^{(i)}, \theta_{j+1}^{(i)}, \dots, \theta_p^{(i)})^{24}$$
  - (c) Repeat the previous step a minimum of  $k_{thinning}$  times until all components have been updated at least once
  - (d) Update  $i = i + 1$
3. Return  $\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(N)}$

The second method is based on the rejection sampling algorithm using an independent product of one-dimensional von Mises distributions.<sup>23</sup> This method is only applicable when the matrix  $\mathbf{P}$  is positive-definite. Additionally, the acceptance probability of the sampler decreases exponentially on the number of variables  $p$ , but also depends on the eigenvalues of  $\mathbf{P}$ . In our experiments, we will use the rejection sampling whenever the number of variables is low ( $p \leq 10$ ), preferring the Gibbs sampler for medium and high dimensionality settings ( $p > 10$ ).

### 2.3. Parameter Learning: Maximum Pseudo-Likelihood Estimator

Using the maximum likelihood estimator for the multivariate von Mises distribution is an often complicated and costly process due to the lack of a known closed formula for the normalization term  $C(\boldsymbol{\kappa}, \boldsymbol{\Lambda})$  in Equation (4). Although with high concentration values and moderate correlation between variables  $C(\boldsymbol{\kappa}, \boldsymbol{\Lambda})$  can be approximated by a Taylor series expansion,<sup>9</sup> this procedure cannot be extended to the general case. As a consequence, for each possible value of  $\boldsymbol{\kappa}$  and  $\boldsymbol{\Lambda}$  the value of the  $C(\boldsymbol{\kappa}, \boldsymbol{\Lambda})$  has to be approximated numerically. Since we select the parametric configuration comparing the likelihood (or pseudo-likelihood) value, even relatively small errors in the estimation of the normalization term could lead to a completely wrong parameter estimation. Although there are methods that can provide an estimate avoiding the curse of dimensionality, e.g., Monte Carlo integration with uniform sampling, if we require an extremely precise estimation of the normalization term, the computational cost does not justify the precision gain compared to the use of the pseudo-likelihood. This cost just grows bigger as we need to recalculate the normalization term in each step of the minimization algorithm. For this reason, some authors have proposed the use of other estimators different from the

maximum likelihood estimator, like the pseudo-likelihood, to learn the parameters of a multivariate von Mises distribution given a set of data samples.<sup>6-8</sup>

The pseudo-likelihood of a multivariate density function<sup>25</sup> is defined as the product of all its univariate conditional densities  $f(\theta_j^{(i)} | \theta_1^{(i)}, \dots, \theta_{j-1}^{(i)}, \theta_{j+1}^{(i)}, \dots, \theta_p^{(i)})$ . Fortunately, as we have mentioned before, the univariate conditional densities of a multivariate von Mises distribution are von Mises distributions with known parameters. This makes the maximum pseudo-likelihood estimator computationally tractable, at the expense of providing less efficient estimations than the maximum likelihood.<sup>26</sup> The full pseudo-likelihood for  $N$  independent  $p$ -dimensional samples  $\Theta = \{\theta^{(1)}, \dots, \theta^{(N)}\}$  is expressed as

$$\mathcal{PL}(\Theta | \mu, \kappa, \Lambda) = (2\pi)^{-Np} \prod_{i=1}^N \prod_{j=1}^p \frac{1}{I_0(\kappa_{\setminus j}^{(i)})} \exp \{ \kappa_{\setminus j}^{(i)} \cos(\theta_j^{(i)} - \mu_{\setminus j}^{(i)}) \} \quad (7)$$

where  $\mu_{\setminus j}^{(i)}$  and  $\kappa_{\setminus j}^{(i)}$  are the univariate conditional parameters in Equation (5) given the  $i$ th data sample  $\theta^{(i)}$ . The usual approach is to maximize the natural logarithm of the pseudo-likelihood expression, which is more tractable and does not change the location of the maximum. In our case, the expression of the log pseudo-likelihood is

$$\begin{aligned} \log \mathcal{PL}(\Theta | \mu, \kappa, \Lambda) &= -Np \log(2\pi) \\ &+ \sum_{i=1}^N \sum_{j=1}^p (-\log I_0(\kappa_{\setminus j}^{(i)}) + \kappa_{\setminus j}^{(i)} \cos(\theta_j^{(i)} - \mu_{\setminus j}^{(i)})) \end{aligned} \quad (8)$$

Similarly to the multivariate normal distribution, the maximum likelihood estimator for the mean parameter  $\hat{\mu} = (\hat{\mu}_1, \dots, \hat{\mu}_p)$ , given a data sample, is the principal argument of the sample (circular) mean:

$$\hat{\mu}_j = \arg \sum_{k=1}^N e^{i\theta_j^{(k)}} \quad (9)$$

We can compute  $\hat{\mu}$  and then rotate (center) the data sample  $\Theta$ , so that we can assume  $\mu = 0$  for the rest of this section.

To find the maximum of function (8), we use the low memory extension of the widespread Broyden-Fletcher-Goldfarb-Shanno (BFGS) quasi-Newton method with simple box constraints (L-BFGS-B).<sup>27</sup> The concentration values are restricted to be positive, that is,  $\kappa_j > 0, \forall j = 1, \dots, p$ . An advantage of the quasi-Newton methods is that it suffices to compute first-order derivatives of the function to maximize (8), i.e., we need to compute the partial derivatives of the log pseudo-likelihood

function with respect to  $\kappa_j$  and  $\lambda_{jk}$ . The resulting equations are

$$\begin{aligned}
\frac{\partial \log \mathcal{P}\mathcal{L}}{\partial \kappa_j} &= \sum_{i=1}^N \left[ \frac{\partial \kappa_{\setminus j}^{(i)}}{\partial \kappa_j} (\cos(\theta_j^{(i)} - \mu_{\setminus j}^{(i)}) - A_1(\kappa_{\setminus j}^{(i)})) + \frac{\partial \mu_{\setminus j}^{(i)}}{\partial \kappa_j} \kappa_{\setminus j}^{(i)} \sin(\theta_j^{(i)} - \mu_{\setminus j}^{(i)}) \right] \\
\frac{\partial \log \mathcal{P}\mathcal{L}}{\partial \lambda_{jk}} &= \sum_{i=1}^N \left[ \frac{\partial \kappa_{\setminus j}^{(i)}}{\partial \lambda_{jk}} (\cos(\theta_j^{(i)} - \mu_{\setminus j}^{(i)}) - A_1(\kappa_{\setminus j}^{(i)})) + \frac{\partial \mu_{\setminus j}^{(i)}}{\partial \lambda_{jk}} \kappa_{\setminus j}^{(i)} \sin(\theta_j^{(i)} - \mu_{\setminus j}^{(i)}) \right] \\
\frac{\partial \kappa_{\setminus j}^{(i)}}{\partial \kappa_R} &= \delta(R, j) \frac{\kappa_j}{\kappa_{\setminus j}^{(i)}} \\
\frac{\partial \kappa_{\setminus j}^{(i)}}{\partial \lambda_{R,S}} &= \delta(R, j) \frac{\sin(\theta_S) \sum_{l \neq j} \lambda_{R,l} \sin \theta_l^{(i)}}{\kappa_{\setminus j}^{(i)}} \\
\frac{\partial \mu_{\setminus j}^{(i)}}{\partial \kappa_R} &= \delta(R, j) \frac{-\sum_{l \neq j} \lambda_{R,l} \sin \theta_l^{(i)}}{(\kappa_{\setminus j}^{(i)})^2} \\
\frac{\partial \mu_{\setminus j}^{(i)}}{\partial \lambda_{R,S}} &= \delta(R, j) \frac{\sum_{l \neq j} \lambda_{R,l} \sin \theta_l^{(i)}}{(\kappa_{\setminus j}^{(i)})^2} \tag{10}
\end{aligned}$$

where  $A_1 = \frac{I_1}{I_0}$  is the ratio of the modified Bessel functions of order one and zero, and  $\delta(x, y) = 1$  if  $x = y$  and zero otherwise. In these expressions,  $R$  and  $S$  can be substituted for any valid index from 1 to  $p$ .

### 2.3.1. Computational Complexity Reduction

Our goal in this section is to simplify function (8) and specially (10) from a computational point of view. Our first step is to express certain sums as a matrix product to take advantage of highly efficient implementations of linear algebra computations such as BLAS or LAPACK.<sup>28</sup> We define the  $N \times p$  auxiliary matrix  $\Psi$ :

$$\Psi = \sin(\Theta)\Lambda, \text{ that is, } \psi_{ij} = \sum_{l \neq j} \lambda_{j,l} \sin(\theta_l^{(i)}) \tag{11}$$

Then we focus on the right-hand term of (8). By applying some basic trigonometric identities, we obtain

$$\cos(\theta_j^{(i)} - \mu_{\setminus j}^{(i)}) = \cos \left( \theta_j^{(i)} - \arctan \left( \frac{1}{\kappa_j} \sum_{l \neq j} \lambda_{j,l} \sin(\theta_l^{(i)}) \right) \right)$$

$$\begin{aligned}
&= \cos(\theta_j^{(i)}) \cos\left(\arctan\left(\frac{1}{\kappa_j} \sum_{l \neq j} \lambda_{j,l} \sin(\theta_l^{(i)})\right)\right) \\
&\quad + \sin(\theta_j^{(i)}) \sin\left(\arctan\left(\frac{1}{\kappa_j} \sum_{l \neq j} \lambda_{j,l} \sin(\theta_l^{(i)})\right)\right) \\
&= \cos(\theta_j^{(i)}) \frac{\kappa_j}{\kappa_{\setminus j}^{(i)}} - \sin(\theta_j^{(i)}) \frac{\sum_{l \neq j} \lambda_{j,l} \sin(\theta_l^{(i)})}{\kappa_{\setminus j}^{(i)}} \quad (12)
\end{aligned}$$

If we substitute this expression in (8) and express the sums as elements of the  $\Psi$  matrix defined in Equation (11), we obtain a more compact version of the log pseudo-likelihood function that does not require to compute the inverse of the tangent:

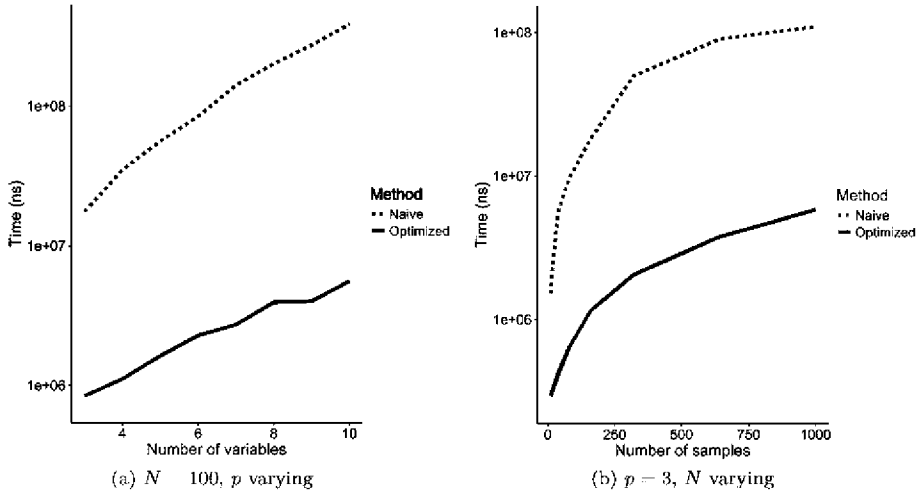
$$\begin{aligned}
\log \mathcal{P}\mathcal{L}_{op}(\Theta | \mu, \kappa, \Lambda) &= -Np \log(2\pi) + \sum_{i=1}^N \sum_{j=1}^p \left(-\log I_0(\kappa_{\setminus j}^{(i)})\right) \\
&\quad + \cos(\theta_j^{(i)}) \kappa_j - \sin(\theta_j^{(i)}) \psi_{ij} \quad (13)
\end{aligned}$$

Additionally, the partial derivatives obtained from Equation (13) are simpler partial derivatives than those in (10):

$$\begin{aligned}
\frac{\partial \log \mathcal{P}\mathcal{L}_{op}}{\partial \kappa_j} &= \sum_{i=1}^N \left[ \cos(\theta_j^{(i)}) - A_1(\kappa_{\setminus j}^{(i)}) \frac{\kappa_j}{\kappa_{\setminus j}^{(i)}} \right] \\
\frac{\partial \log \mathcal{P}\mathcal{L}_{op}}{\partial \lambda_{jk}} &= \sum_{i=1}^N \left[ \sin(\theta_k^{(i)}) \left( \sin(\theta_j^{(i)}) - A_1(\kappa_{\setminus j}^{(i)}) \frac{\psi_{ij}}{\kappa_{\setminus j}^{(i)}} \right) \right] \quad (14)
\end{aligned}$$

To evaluate the performance gain between expressions (8) and (13), we compare the average execution time to compute the maximum pseudo-likelihood estimation using each one. To maximize the likelihood function, we use the L-BFGS-B algorithm, a quasi-Newton method with reduced memory footprint. In both cases, the routines that perform most of the calculations are implemented in *ANSI C*. We propose two experiments: First, the number of variables  $p$  varies for a fixed number of samples  $N$  (Experiment A); and second, the other way around (Experiment B). We repeat each experiment 100 times for every one of the 10 randomly generated  $\mathbf{P}$  positive definite matrices to compute the average time as well as the standard deviation. The experimental configuration is summarized in Table I. The experiments were executed in a common desktop computer with an i7-4790k processor at 4 GHz and 8 GB of RAM at 1600 MHz.





**Figure 1.** Average fitting time for the naive implementation (8) and the optimized version (13) of the log pseudo-likelihood function. Each execution has been repeated 100 times with two additional warm-up iterations. We have plotted an envelope of width  $\sigma$  around the average time line, but, due to its small size, it can be hardly seen.

**Table I.** Experimental configuration summary for time comparison.

Experiment	$p$	$N$	Number of repetitions
A	$\{3, 4, \dots, 10\}$	100	10 random $\mathbf{P}$ matrices $\times$ 100 times
B	3	$\{10, 20, 40, \dots, 1000\}$	10 random $\mathbf{P}$ matrices $\times$ 100 times

Figure 1 shows that the optimized version is significantly faster in every case, and it scales better when the number of samples ( $N$ ) or variables ( $p$ ) increases, although both versions have the same computational complexity  $\mathcal{O}(p^2 N)$ . Based on these results, we will use this optimized version of the log pseudo-likelihood function in the next experiments.

### 2.3.2. Estimator Properties

Asymptotic consistency of an estimator is perhaps the most important property of any admissible estimator as it guarantees that the estimator is arbitrarily close to the estimated value provided sufficient number of samples. In other words, the estimation improves as the number of samples increases. The consistency of the maximum pseudo-likelihood estimator has been already proved by other authors.<sup>8</sup>

Although there are some related works on the analytical study of mean and variance properties on implicitly defined biased estimators,<sup>29</sup> we follow an empirical approach to evaluate the bias and the variance of the maximum pseudo-likelihood estimator when estimating  $\kappa$  and  $\mathbf{\Lambda}$ . Specifically, we aim to determine the

**Table II.** Experiments to estimate bias and variance properties.

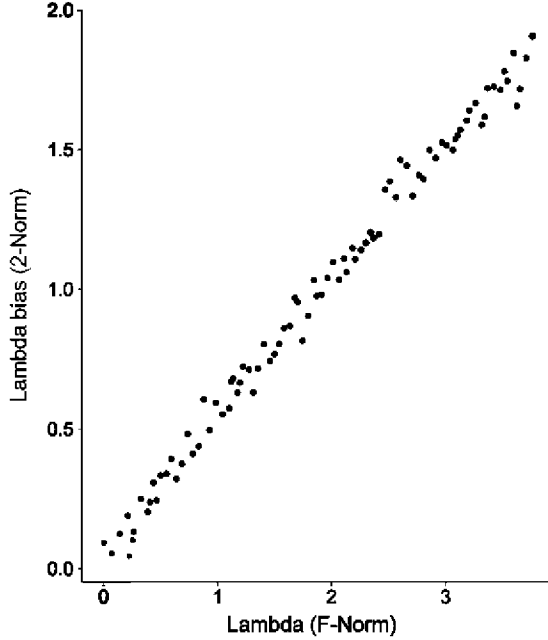
Experiment	$p$	$N$	$\kappa$	$\mathbf{\Lambda}$	Number of repetitions
A	3	10	(3, 3, 3)	$\mathbf{\Lambda} = \mathbf{0}$	5000
B	3	10	(3, 3, 3)	$\mathbf{\Lambda} \neq \mathbf{0}$	5000
C	3	10	(0.01, 0.01, 0.01)	$\mathbf{\Lambda} = \mathbf{0}$	5000
D	3	10	(0.01, 0.01, 0.01)	$\mathbf{\Lambda} \neq \mathbf{0}$	5000
E	10	50	(3, ..., 3)	$\mathbf{\Lambda} = \mathbf{0}$	5000
F	10	50	(3, ..., 3)	$\mathbf{\Lambda} \neq \mathbf{0}$	5000
G	50	100	(3, ..., 3)	$\mathbf{\Lambda} = \mathbf{0}$	10000
H	50	100	(3, ..., 3)	$\mathbf{\Lambda} \neq \mathbf{0}$	10000

properties of the estimator for low and medium concentration values ( $\kappa$ ), as well as for configurations when the variables are independent ( $\mathbf{\Lambda} = \mathbf{0}$ ) and dependent ( $\mathbf{\Lambda} \neq \mathbf{0}$ ). Based on the previous goals, we define the experimental configurations listed in Table II.

For each parametric configuration, we repeat the same process a high number of times: First,  $N$  independent samples are generated from a MVM( $\boldsymbol{\mu} = \mathbf{0}$ ,  $\kappa$ ,  $\mathbf{\Lambda}$ ). Then, the maximum pseudo-likelihood estimator is applied to get the estimated values  $\hat{\boldsymbol{\xi}}_{\mathcal{P}\mathcal{L}} = (\hat{\boldsymbol{\kappa}}, \hat{\mathbf{\Lambda}})$ . Finally, after all repetitions have been completed, we compute the bias and the variance of the estimator. Note that in the cases where  $\mathbf{\Lambda} = \mathbf{0}$  the pseudo-likelihood estimator is in fact the maximum likelihood estimator.

For each experiment, we apply the Hotelling's T-square test,<sup>30</sup> a multivariate hypothesis test, over  $\hat{\boldsymbol{\kappa}}$  and  $\hat{\mathbf{\Lambda}}$  independently to verify whether the estimator is unbiased for that configuration. For those cases where the unbiasedness is rejected, we perform a one-sided  $t$ -test for each individual parameter  $\xi_j$  to check if its estimator is overestimating the real value ( $\mathbb{E}[\frac{\hat{\xi}_j}{\xi_j}] > 1$ ) or either underestimating it. The  $\hat{\boldsymbol{\kappa}}$  unbiasedness hypothesis is rejected with a significance level  $\alpha = 0.05$  in every case, whereas it is not rejected for  $\hat{\mathbf{\Lambda}}$  in experiments A, C, and E and rejected in the rest with the same  $\alpha$ . Defying our belief that the  $\hat{\mathbf{\Lambda}}$  estimator is unbiased when  $\mathbf{\Lambda} = \mathbf{0}$ , the hypothesis is rejected in experiment G, even when the average bias value is low ( $\approx 10^{-2}$ ) for each  $\hat{\lambda}_{ij}$ . Therefore, due to this last result, we cannot draw any conclusion about the unbiasedness hypothesis. Finally, the overestimation hypothesis  $H_0 : \mathbb{E}[\frac{\hat{\xi}_j}{\xi_j}] > 1$  is not rejected for every single parameter where the unbiasedness hypothesis has been rejected, which means that the maximum pseudo-likelihood statistic overestimates the real value of the parameters.

Finally, we perform a last experiment. We want to check for a relation between the distance of  $\mathbf{\Lambda}$  to the zero matrix, measured as the Frobenius norm of  $\mathbf{\Lambda}$  ( $\|\mathbf{\Lambda}\|_F$ ), and  $Bias(\hat{\mathbf{\Lambda}})$ . To do so, we generate  $N = 10$  samples from a three-variate MVM with fixed  $\kappa = (3, 3, 3)$  and compute the Euclidean norm of the estimator bias while varying  $\|\mathbf{\Lambda}\|_F$  from 0 to 4, verifying in each case that the matrix  $\mathbf{P}$  is positive definite. For every parameter configuration, the process is repeated 10,000 times. As seen in Figure 2, there is a clear linear relationship between  $\|\mathbf{\Lambda}\|_F$  and  $\|Bias(\hat{\mathbf{\Lambda}})\|_2$ .



**Figure 2.** Relation between the distance to the zero matrix  $\|\mathbf{\Lambda}\|_F$  and  $\|Bias(\hat{\mathbf{\Lambda}})\|_2$ .

### 3. PENALIZED PARAMETER LEARNING

Learning the parameters of a multivariate von Mises distribution using a penalized maximum pseudo-likelihood estimator has been already proposed in the literature by other authors.<sup>8</sup> Their approach is to add a term in the pseudo-likelihood expression that penalizes the absolute value of all  $\lambda_{ij}$  elements equally. Since the pseudo-likelihood overestimates the parameters, specially when learning from small samples, it is broadly justified to apply a penalization on the parameter absolute value. The actual problem is that, if we apply the exact same penalty on each  $\lambda_{ij}$  we are implicitly assuming that all  $\lambda_{ij}$  are similar in size.

From the Bayesian data analysis perspective,<sup>31</sup> the parameters  $\kappa$  and  $\mathbf{\Lambda}$  are real-valued random variables with a specific distribution given the data  $\Theta$ , the posterior probability distribution  $p(\kappa, \mathbf{\Lambda} | \Theta)$ . This function is obtained from the prior probability distribution  $p(\kappa, \mathbf{\Lambda})$ , a distribution that reflects the prior knowledge about the value of the parameters, and the aforementioned likelihood function through Bayes' theorem:

$$p(\kappa, \mathbf{\Lambda} | \Theta) \propto p(\kappa, \mathbf{\Lambda})p(\Theta | \kappa, \mathbf{\Lambda}) \quad (15)$$

Hence, the maximum likelihood estimator corresponds to a scenario with an uninformative prior where maximizing the posterior probability actually corresponds to

maximizing the likelihood. Under this Bayesian framework, a uniform  $L_1$  regularization corresponds to setting the exact same prior distribution, centered at zero, for each parameter and to assuming that all  $\lambda_{ij}$  are similar in size. This issue has been already studied in the literature, but focusing on multivariate normal distributions.<sup>18,19</sup> Here we propose a penalized estimator based on the maximum pseudo-likelihood estimator where each parameter is penalized individually; in other words, the prior distribution for each parameter is independent from the rest.

### 3.1. F-Norm Penalized Pseudo-Likelihood

Let  $\mathbf{P}$  be the positive definite matrix defined in Equation (6),  $\Phi$  a  $p \times p$  real symmetric matrix, and  $\mathbf{H}$  a triangular real  $p \times p$  matrix, whose elements  $h_{ij}$  are nonnegative. Then, we add the following penalization term to the log pseudo-likelihood introduced in Equation (13):

$$pen(\boldsymbol{\kappa}, \boldsymbol{\Lambda}) = -\|(\mathbf{P} - \Phi) \circ \mathbf{H}\|_F \quad (16)$$

where  $\|A\|_F = \sqrt{\text{tr}(AA^T)}$  is the Frobenius matrix norm (F-norm) and  $\circ$  is the entrywise Hadamard product  $A \circ B = (a_{ij}b_{ij})$ . In the penalization term given by Equation (16), the matrix  $\Phi$  is our prior guess about the elements  $p_{ij}$  of  $\mathbf{P}$ ; whereas  $h_{ij}$ , the elements of the confidence matrix  $\mathbf{H}$ ,<sup>32</sup> measures our degree of confidence in the correctness of the value  $\phi_{ij}$ . Note that the triangular structure imposed to  $\mathbf{H}$  is not arbitrary since it prevents the same  $\lambda_{ij} = \lambda_{ji}$  parameter from being penalized twice.

The last requirement to use the penalized pseudo-likelihood with the term given by Equation (16) is to update the partial derivatives in Equation (14) by adding the following terms:

$$\begin{aligned} \frac{\partial pen(\boldsymbol{\kappa}, \boldsymbol{\Lambda})}{\partial \kappa_j} &= -\frac{h_{jj}^2(\kappa_j - \phi_{jj})}{\|(\mathbf{P} - \Phi) \circ \mathbf{H}\|_F} \\ \frac{\partial pen(\boldsymbol{\kappa}, \boldsymbol{\Lambda})}{\partial \lambda_{jk}} &= -\frac{h_{jk}^2(\lambda_{jk} - \phi_{jk})}{\|(\mathbf{P} - \Phi) \circ \mathbf{H}\|_F} \end{aligned} \quad (17)$$

Now, we face the problem of selecting values for  $\Phi$  and  $\mathbf{H}$  such that the estimation given by the penalized version is at least as good as the unpenalized one. Based on the results in the preceding section that show how the maximum pseudo-likelihood estimator overestimates the scale of  $\boldsymbol{\kappa}$  and  $\boldsymbol{\Lambda}$ , we can mimic the  $L_1$  penalization and set  $\Phi = \mathbf{0}$  if no prior information is available, but we still need to define the values of the confidence matrix for each parameter. A conservative approach, given the absence of reliable prior information, is to set the elements  $h_{ij}$  to some medium-low value close to 1; otherwise, there is a risk of underestimating the parameters. Obviously, the exact best value for each  $h_{ij}$  depends on the number of variables  $p$  and the real value of  $p_{ij}$ . In Section 4, we show a real example where the structure of this confidence matrix is clearly induced by the problem itself. Another interesting

approach for highly concentrated data, could be to use the inverse of the sample circular covariance matrix<sup>7</sup> as our  $\Phi$  matrix.

Going back to the Bayesian analysis framework, the penalization term in Equation (16) in the log pseudo-likelihood could be intuitively seen as placing an independent normal prior  $\mathcal{N}(\phi_{ij}, \sigma(h_{ij}))$  over each  $p_{ij}$ .

### 3.2. Properties

As we did with the maximum pseudo-likelihood estimator in Section (2.3.2), in this section we analyze the properties of the penalized estimator and compare the experimental results between the unpenalized and the penalized maximum pseudo-likelihood estimators. The following proposition addresses the consistency of the penalized estimator.

PROPOSITION 1. *The penalized maximum pseudo-likelihood estimator*

$$\hat{\xi}_{\mathcal{P}\mathcal{L}} := \operatorname{argmax}_{\xi=(\kappa, \Lambda)} \{ \log \mathcal{P}\mathcal{L}(\Theta | \mu, \kappa, \Lambda)_{op} + \operatorname{pen}(\kappa, \Lambda) \}$$

*is consistent.*

*Proof.* Let  $\hat{\xi}_{\mathcal{P}\mathcal{L}}^n$  and  $\hat{\xi}_{\mathcal{P}\mathcal{P}\mathcal{L}}^n$  be the unpenalized and penalized maximum pseudo-likelihood estimators, respectively, and  $\xi_0$  the real parameter that exists is unique and well-defined. Then, we can define that the maximum pseudo-likelihood estimator  $\hat{\xi}_{\mathcal{P}\mathcal{L}}^n$  as an M-estimator that maximizes the function:

$$\mathbb{M}_n(\xi) = \frac{1}{n} \sum_{i=1}^n (\log \mathcal{P}\mathcal{L}(\theta^i | \mu, \kappa, \Lambda))$$

If we restrict ourselves to a compact neighborhood of  $\xi_0$ ,  $\mathcal{B}(\xi_0)$ , and having in mind that the log-pseudo-likelihood function is continuous and square-integrable, then the conditions to apply the uniform strong law of large numbers are met on that neighborhood, as a result:

$$\sup_{\xi \in \mathcal{B}(\xi_0)} \|\mathbb{M}_n(\xi) - \mathbb{E}[\log \mathcal{P}\mathcal{L}(\theta | \mu, \kappa, \Lambda)]\| \xrightarrow{a.s.} 0$$

For convenience, in the rest of the proof we will refer to the function  $\mathbb{E}[\log \mathcal{P}\mathcal{L}(\theta | \mu, \kappa, \Lambda)]$  simply as  $M(\xi)$ . Please recall that  $\xi = (\kappa, \Lambda)$ .

On the other hand, let  $h_n(\xi) = \frac{\|(\mathbf{P}-\Phi) \circ \mathbf{H}\|_F}{N}$  be our complexity penalization function. It is clear that  $h_n$  is continuous, positive and the sequence decreases monotonically (i.e.,  $h_n < h_m, \forall m > n$ ). Since the numerator is obviously finite in any compact set of the parameter space,  $h_n \rightarrow 0$  pointwise. Furthermore, by applying Dini's theorem, we can state that  $h_n$  converges to 0 uniformly in any compact subspace of the parameter space.

As we did before, the penalized maximum pseudo-likelihood estimator  $\hat{\xi}_{\mathcal{P}\mathcal{P}\mathcal{L}}^n$  can be also defined as an M-estimator:

$$\mathbb{M}_n^{pen}(\xi) = \frac{1}{n} \sum_{i=1}^n (\log \mathcal{P}\mathcal{L}(\theta^i | \mu, \kappa, \Lambda)) + h_n(\xi)$$

Now, due to the fact that  $\xi_0$  exists, is unique and well-defined,  $\exists R$  such that  $\forall r < R \ S_r : \{\xi \mid \|\xi_0 - \xi\| \leq \frac{r}{2}\}$  verifies that

$$\forall \xi \notin S_r \ M(\xi) \leq M(\xi_0) - \epsilon_r$$

Therefore,  $\forall \delta > 0 \ \exists N_1 \mid \forall n_1 > N_1$  such that

$$P \left\{ \sup_{\xi \in \mathcal{B}(\xi_0)} \|\mathbb{M}_n(\xi) - M\| > \frac{\epsilon_r}{2} \right\} < \delta$$

In parallel,  $\sup_{\xi \in \mathcal{B}(\xi_0)} \|\mathbb{M}_n - \mathbb{M}_n^{pen}\| = \sup_{\xi \in \mathcal{B}(\xi_0)} h_n$ , as  $h_n$  converges uniformly to 0 in any compact, then  $\exists N_2 \mid \forall n_2 > N_2 \ \sup_{\xi \in \mathcal{B}(\xi_0)} h_n < \frac{\epsilon_r}{2}$ . It follows that

$$\begin{aligned} \sup_{\xi \in \mathcal{B}(\xi_0)} \|M - \mathbb{M}_n^{pen}\| &\leq \sup_{\xi \in \mathcal{B}(\xi_0)} \|M - \mathbb{M}_n\| + \sup_{\xi \in \mathcal{B}(\xi_0)} \|\mathbb{M}_n - \mathbb{M}_n^{pen}\| \\ &\Rightarrow P \left\{ \sup_{\xi \in \mathcal{B}(\xi_0)} \|M - \mathbb{M}_n^{pen}\| > \epsilon_r \right\} < \delta \end{aligned}$$

which implies that  $P\{\|\xi_0 - \hat{\xi}_{\mathcal{P}\mathcal{P}\mathcal{L}}^n\| > \frac{r}{2}\} < \delta \ \forall n > \max(N_1, N_2)$ . Then when can take  $r \rightarrow 0$  and  $\delta \rightarrow 0$  so the argument holds when  $n \rightarrow \infty$ .  $\square$

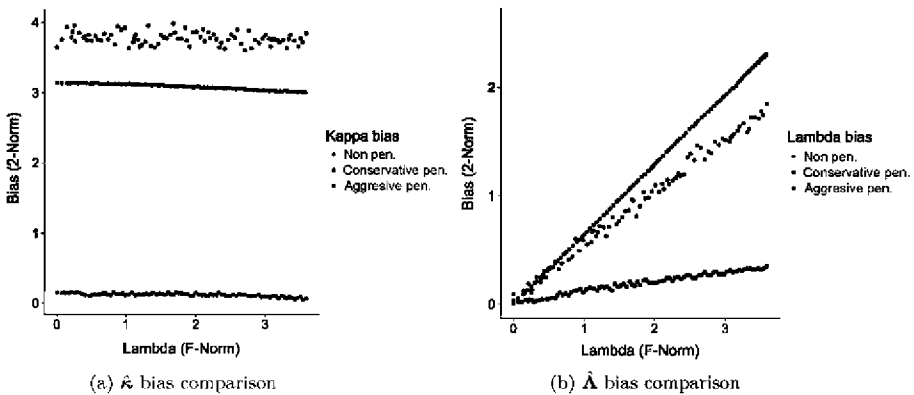
To experimentally study the properties of the penalized estimator and compare it with the unpenalized version, we repeat the experiments of Table II. Following our own recommendations, we will test the penalized estimator with prior matrix  $\Phi = 0$  and compare two penalization strategies: a conservative approach where the elements  $h_{i \geq j} = 0.5$ , and a more aggressive penalization with mid-high confidence values  $h_{i \geq j} = 3$ . Our expectation is that while the conservative approach will reduce both bias and variance with respect to the nonpenalized estimator, the aggressive approach will increase the bias due to underestimation in low-dimensional settings where the real parameters values are far from zero, but with a much lower variance. On the other hand, in the cases where the number of parameters is high we expect a sharp increase in the parameter overestimation produced by the pseudo-likelihood estimator. It is reasonable to believe that in those cases a higher penalization would perform better. Results in Table III agree with our expectations.

Finally, we also repeat the low-dimensional experiment shown in Figure 2 including the two penalized approaches. In Figure 3, we can see how the conservative approach clearly reduces the bias in both  $\hat{\Lambda}$  and  $\hat{\kappa}$  estimations and outperforms the

**Table III.** Euclidean norm of the bias, trace of the variance matrix, and mean square error (MSE) comparison between unpenalized, conservative penalization, and aggressive penalization for the experiments defined in Table (II).

Experiment	Nonpenalized			Penalized $h_{ij} = 0.5$			Penalized $h_{ij} = 3$		
	Bias	Var	MSE	Bias	Var	MSE	Bias	Var	MSE
A	3.75	61.69	75.80	<b>0.14</b>	<b>5.55</b>	<b>5.57</b>	3.13	0.13	9.90
B	3.94	80.57	96.70	<b>1.38</b>	<b>5.31</b>	<b>7.22</b>	4.57	0.18	21.07
C	1.17	3.99	5.37	1.08	1.90	3.07	<b>1.00</b>	<b>0.64</b>	<b>1.65</b>
D	1.31	4.04	5.77	1.24	2.19	3.74	<b>1.15</b>	<b>0.70</b>	<b>2.04</b>
E	1.88	31.15	34.70	1.27	24.08	25.69	<b>1.10</b>	<b>7.16</b>	<b>8.38</b>
F	2.24	61.15	66.18	<b>1.71</b>	43.79	49.72	2.04	<b>6.47</b>	<b>10.65</b>
G	14.34	1717.14	1922.76	14.10	1712.23	1910.97	<b>10.76</b>	<b>852.14</b>	<b>967.99</b>
H	14.23	1676.50	1879.01	13.64	1469.97	1655.95	<b>10.16</b>	<b>798.11</b>	<b>901.24</b>

Best results are in bold characters.



**Figure 3.** Bias comparison between unpenalized, conservative penalization, and aggressive penalization for a three-variate von Mises distribution with  $\kappa = (3, 3, 3)$  from  $N = 10$  samples and  $\|\Lambda\|_F$  varying from 0 to 4.

aggressive approach. As we anticipated before, Figure 3b shows how the most aggressive penalization actually increases the bias of  $\hat{\Lambda}$  due to underestimation, which it is not the case in Figure 3a.

#### 4. A REAL-WORLD APPLICATION IN NEUROANATOMY

As an application of the penalized estimator on real-world data, we use the von Mises multivariate distribution to model the angles between basal dendrites of pyramidal cells. A neuron can be divided into three main parts: the cell body or soma, dendrites, and axon. Among all the different types of neurons in the nervous system, pyramidal neurons stand out as one of the most important types. They play a key role in the connectivity of cortical columns which are the functional blocks in

the neocortex. Since the morphology of a neuron determines its connectivity pattern to a large extent,<sup>33</sup> we are interested in studying the pyramidal neurons from an anatomical point of view. A unique characteristic of the pyramidal neurons is that they have basal dendrites, so called because they grow from the base of the soma and spread out horizontally. Our goal is to characterize the angles between these basal dendrites and to find differences or similarities between species, brain regions, etc. Ultimately, this knowledge will lead to better single-cell simulations and will improve our understanding on how the brain works.

#### 4.1. Evaluation: Approximated KL Divergence

First of all, we need to determine how to assess the goodness-of-fit of our models. Here we face the problem of dealing with multivariate directional data, for which we have very few available tools compared to the well-known field of multivariate linear data. Other authors use a variety of techniques to evaluate the goodness of fit, such as using a multivariate normal approximation<sup>7</sup> that cannot be applied in the general case or they assess the log-likelihood value,<sup>8</sup> which can be hard to compute accurately for the multivariate von Mises distribution, and it is not appropriate to compare distributions from different families.

Here, we propose to use an approximation of the KL divergence for multivariate distributions that does not make any assumption on the underlying true distribution of the samples<sup>34</sup> as evaluation metric. Suppose that we have two sets of samples  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  and  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^M$  from two  $p$ -dimensional unknown probability distributions,  $P$  and  $Q$ , respectively. Then the estimated KL divergence is

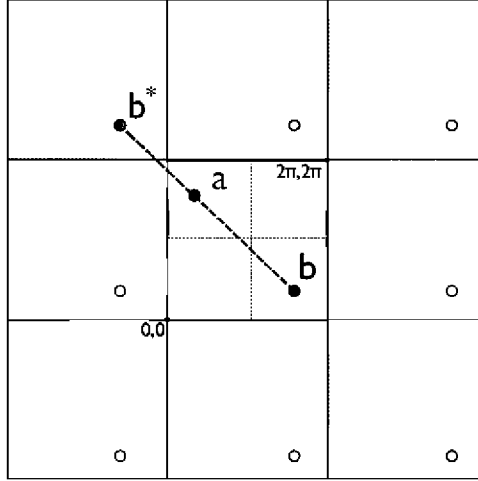
$$\hat{D}_k(P||Q) = \frac{p}{N} \sum_{i=1}^N \left[ \log \left( \frac{s_k(\mathbf{x}_i)}{r_k(\mathbf{x}_i)} \right) \right] + \log \frac{M}{N-1} \quad (18)$$

where  $r_k(\mathbf{x}_i)$  and  $s_k(\mathbf{x}_i)$  are the distance to the  $k$ th nearest neighbor of  $\mathbf{x}_i$  in  $\mathbf{X} \setminus \{\mathbf{x}_i\}$  and  $\mathbf{Y}$ , respectively.

With this approximation of the KL divergence between the unknown underlying distribution of the data and the fitted model, we can measure how far is the fitted distribution from the data just by taking one set of samples from the data and the other generated by sampling the fitted distribution using any method defined in Section 2.2 for the multivariate von Mises distribution. However, there is still a missing piece to compute the estimation: We need to define which distance are we going to use. Any linear norm such as the Euclidean norm is not adequate since they do not take into account the periodicity of the data. Another option would be to use the geodesic distance in the hyper-torus  $\mathbb{T}^p$ , but it is rather complicated to compute. To overcome these drawbacks, we define a distance between two points  $\mathbf{a}, \mathbf{b} \in [0, 2\pi)^p$  in Equation (19) that is easy to compute and takes into account the periodicity of circular data.

$$d(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}^*\|_2 \quad (19)$$





**Figure 4.** Circular distance computation for  $p = 2$  defined by Equation (19).

where

$$\mathbf{b}^* = (b_i^*)_{i=1}^p = \begin{cases} b_i & \text{if } |a_i - b_i| \leq \pi \\ b_i + 2\pi & \text{if } |a_i + b_i| > \pi \text{ and } a_i > \pi \\ b_i - 2\pi & \text{if } |a_i - b_i| > \pi \text{ and } a_i \leq \pi \end{cases}$$

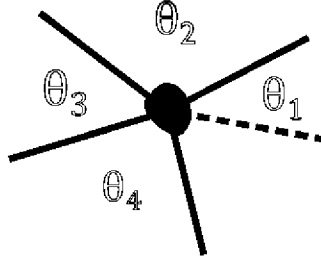
Figure 4 provides a graphical explanation on how this distance works for  $p = 2$ .

## 4.2. Data Extraction and Modeling

Two sets of 3D reconstructions of pyramidal neurons were downloaded from NeuroMorpho.Org,<sup>35</sup> a public repository of neural reconstructions. One contains 1915 reconstructions of human pyramidal cells from different brain regions,<sup>36</sup> whereas the other has 90 reconstructions of mouse pyramidal cells from the neocortex.<sup>37</sup> A third set of 38 reconstructions of rat pyramidal neurons from the somatosensory cortex was downloaded from a different source.<sup>38,39</sup> Table IV contains the count of reconstructions per data set aggregated by the number of basal dendrites. We only pick those cases where the number of samples is at least six; otherwise, we are not able to estimate the KL divergence properly as with smaller samples the test set would be too small to provide a reliable KL estimation. If a neuron has  $p$  basal dendrites, then we have  $p - 1$  variables between dendrites since the last angle is fully determined by the rest. Note that there is a linear restriction over the angles as they must total less than  $2\pi$ , although it is not considered in the model. The restriction should be taken into account when computing the conditional distribution of one angle given the others, as the range of accepted values is restricted. However, the model is still valid as a first approach, and these issues will be considered as future improvements.

**Table IV.** Number of reconstructions per data set aggregated by the number of basal dendrites.

Basal count	Data set		
	Human	Mouse	Rat
1	8	0	0
2	11	0	0
3	133	1	1
4	475	11	4
5	598	31	13
6	420	32	6
7	143	9	4
8	87	5	7
9	18	1	1
10	11	0	2
11	8	0	0
12	2	0	0
13	1	0	0
Total	1915	90	38

**Figure 5.** Angles between dendrites.

Since basal dendritic roots do not exactly lie in a plane, we need to project them onto the basal plane, defined as the plane that minimizes the distance of the basal roots to their projection onto the plane, which is exactly what principal component analysis solves. Then, we need to establish a consistent criterion to select the basal tree from which we start numbering in counter-clockwise order as shown in Figure 5. The data extraction procedure is detailed in Algorithm 2, which is applied to measure the angles from each reconstruction.

ALGORITHM 2.

Input: 3D reconstruction

Output:  $\theta_1, \dots, \theta_p$  angles between basal dendrites

Steps:

1. Extract Cartesian coordinates of the basal roots
2. Find the basal plane using principal component analysis on the coordinates of the basal roots
3. Project the basal roots onto the basal plane
4. Measure the total length of each basal tree and designate the longest as the first
5. Measure planar angles between the projected basal roots in counter-clockwise order

**Table V.** Mean KL divergence results for angles between basal dendrites  $\pm$  standard deviation (lower is better).

$p$	Samples	Multivariate von Mises				MV Normal	
		Nonpenalized	$h_{j>i} = 0.5$	$\mathbf{H} = \mathbf{H}_{dist}$	$L_1(1)$	Linear	Wrapped
<i>Rat</i>							
4	12	1.56 $\pm$ 1.26	0.62 $\pm$ 0.42	0.68 $\pm$ 0.45	<b>0.51 <math>\pm</math> 0.36</b>	0.79 $\pm$ 0.76	1.09 $\pm$ 0.90
7	8	9.31 $\pm$ 1.86	2.72 $\pm$ 0.85	<b>1.96 <math>\pm</math> 0.64</b>	2.60 $\pm$ 0.85	2.69 $\pm$ 0.85	2.69 $\pm$ 0.94
<i>Mouse</i>							
3	11	2.04 $\pm$ 1.16	1.01 $\pm$ 0.48	<b>0.56 <math>\pm</math> 0.44</b>	0.75 $\pm$ 0.48	1.33 $\pm$ 0.85	1.78 $\pm$ 1.28
4	31	0.37 $\pm$ 0.35	0.72 $\pm$ 0.40	<b>0.34 <math>\pm</math> 0.26</b>	0.55 $\pm$ 0.31	<b>0.34 <math>\pm</math> 0.27</b>	0.35 $\pm$ 0.30
5	32	0.44 $\pm$ 0.35	0.61 $\pm$ 0.40	0.39 $\pm$ 0.28	0.42 $\pm$ 0.28	<b>0.35 <math>\pm</math> 0.29</b>	<b>0.35 <math>\pm</math> 0.26</b>
6	9	8.95 $\pm$ 1.99	3.07 $\pm$ 0.64	<b>0.94 <math>\pm</math> 0.63</b>	1.78 $\pm$ 1.11	2.29 $\pm$ 1.12	2.65 $\pm$ 1.30
<i>Human</i>							
3	475	<b>0.51 <math>\pm</math> 0.10</b>	0.53 $\pm$ 0.10	0.52 $\pm$ 0.09	0.52 $\pm$ 0.09	0.55 $\pm$ 0.09	1.16 $\pm$ 0.11
4	598	<b>0.74 <math>\pm</math> 0.10</b>	<b>0.74 <math>\pm</math> 0.10</b>	0.76 $\pm$ 0.10	0.76 $\pm$ 0.11	0.85 $\pm$ 0.11	0.94 $\pm$ 0.11
5	420	0.87 $\pm$ 0.13	0.88 $\pm$ 0.11	0.85 $\pm$ 0.12	0.85 $\pm$ 0.11	0.80 $\pm$ 0.13	<b>0.68 <math>\pm</math> 0.12</b>
6	143	0.67 $\pm$ 0.21	0.70 $\pm$ 0.21	0.63 $\pm$ 0.23	0.69 $\pm$ 0.24	0.51 $\pm$ 0.17	<b>0.47 <math>\pm</math> 0.22</b>
7	87	0.33 $\pm$ 0.22	0.56 $\pm$ 0.26	0.35 $\pm$ 0.22	0.45 $\pm$ 0.28	<b>0.27 <math>\pm</math> 0.24</b>	0.29 $\pm$ 0.23
8	18	4.69 $\pm$ 1.96	1.71 $\pm$ 0.70	0.80 $\pm$ 0.69	<b>0.63 <math>\pm</math> 0.43</b>	0.84 $\pm$ 0.68	0.84 $\pm$ 0.68
9	11	8.96 $\pm$ 1.60	3.10 $\pm$ 0.93	<b>1.41 <math>\pm</math> 0.93</b>	2.88 $\pm$ 1.04	2.26 $\pm$ 1.19	2.25 $\pm$ 1.09
Average rank		4.62	4.31	2.23	3.38	3.00	3.46

For the F-norm penalized multivariate von Mises estimation, the prior matrix is  $\Phi = \mathbf{0}$  in both cases, whereas the penalization parameter is 1 in the  $L_1$  case. Best results are in bold characters.

6. Return  $\theta_1, \dots, \theta_p$

Probably, the weakest point in the procedure described in Algorithm 2 is the determination of the order of the angles, i.e., selecting the first angle based on the longest basal tree. In the absence of a common coordinate system for all reconstructions in a data set, any local criterion would be (up to some point) arbitrary. Indeed, we have verified that the null hypothesis that all marginal distributions are equal is not rejected by the  $q$ -sample uniform-scores test<sup>12</sup> in all cases included in Table V. This result suggests that either the criterion is totally arbitrary, the reconstruction angles are truly equidistributed, or both are true. Despite these shortcomings, this method is valuable as, to the best of our knowledge, there are no other multivariate approaches in the literature in this regard.

In our experiments, we compare the KL divergence estimate between the underlying distribution and the multivariate von Mises distributions with parameters estimated with the unpenalized maximum pseudo-likelihood, penalized with  $\Phi = \mathbf{0}$  and  $h_{j>i} = 0.5$  and penalized with  $\Phi = \mathbf{0}$  and  $\mathbf{H} = \mathbf{H}_{dist}$ . Here,  $\mathbf{H}_{dist}$  is defined in such a way that  $h_{ij}$  is linearly proportional to the distance between the  $i$ th and  $j$ th angles with  $\max h_{j>i} = 1$  and diagonal elements  $h_{ii} = \min h_{j>i}$ . For example,

in the case depicted in Figure 5 the  $\mathbf{H}_{dist}$  matrix is

$$\begin{pmatrix} 0.5 & 0.5 & 1 & 1 \\ 0 & 0.5 & 0.5 & 1 \\ 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0 & 0.5 \end{pmatrix}$$

In addition, we include in the comparison the  $L_1$  penalized multivariate von Mises distribution,<sup>8</sup> the multivariate normal and the wrapped-normal multivariate distributions with parameters estimated using the method of moments but shrinking the covariance matrix if it is not positive definite. To get an honest estimation of the KL divergence, we performed repeated train and test validation with 100 repetitions. In each trial, the 50% of the original samples are selected at random without replacement as training set, leaving the remaining portion of samples as test set to compute the KL divergence estimation.

### 4.3. Results

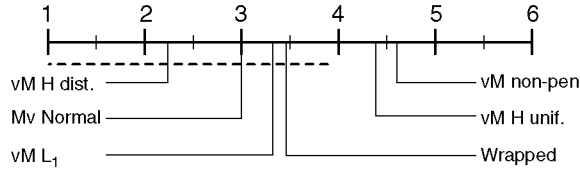
Results are summarized in Table V. The penalized von Mises distribution obtains similar or better results than the unregularized estimation, specially in the cases where the number of samples is low with respect to the number of variables. Also, we can see how the penalization with the structure-aware confidence matrix  $\mathbf{H}_{dist}$  provides, in general terms, better estimations than the ones using a uniform confidence matrix.

To assess the statistical significance of the results, we first apply a variation of the Friedman test, a nonparametric equivalent of the repeated ANOVA. The test compares the ranks of each method in different data sets under the null hypothesis that all methods are equivalent.<sup>40</sup> The statistic

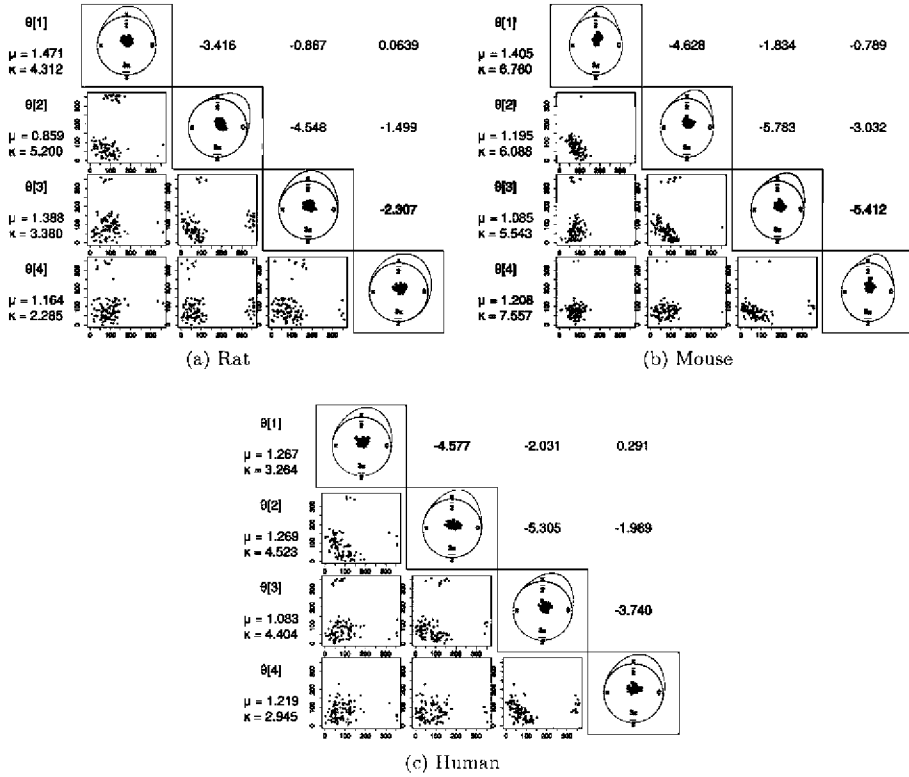
$$F_F = \frac{(N - 1)\chi_F^2}{N(k - 1) - \chi_F^2} \quad (20)$$

is distributed as an F-distribution with  $k - 1$  and  $(k - 1)(N - 1)$  degrees of freedom under the null hypothesis. In Equation (20),  $k$  is the number of methods that we are comparing (six in our case),  $N$  is the number of data sets on which we evaluate our algorithms (13), and  $\chi_F^2$  is the Friedman statistic. The test rejects the null hypothesis with a significance level  $\alpha = 0.05$ , which means that there are significant performance differences between the methods. Then, we compare all methods with the F-norm penalized von Mises estimator with structure-aware confidence matrix  $\mathbf{H}$  using the Bonferroni–Dunn test. The test finds significant differences ( $\alpha = 0.05$ ) with the nonpenalized von Mises estimation and with the F-norm penalized estimator with uniform confidence matrix. Figure 6 summarizes these results.

The plots in Figure 7 summarize the von Mises distribution fitted for pyramidal neurons with five basal dendrites from the three species: human, rat, and mouse. The boxes in the diagonal show the marginal distribution of the original data for each variable. Each diagonal box contains a rose plot (i.e., a circular histogram) and



**Figure 6.** Comparison of F-norm penalized von Mises estimator with structure-aware confidence matrix  $H$  against the others with the Bonferroni–Dunn test. All methods with ranks outside the dashed interval are significantly different ( $\alpha = 0.05$ ) from the control.



**Figure 7.** Fitted distributions for the angles between basal dendrites from rat, mouse and human with five basal dendrites ( $p = 4$ ).

the circular density approximated with the von Mises circular kernel. The boxes in the upper triangle are the  $\lambda_{ij}$  parameters of the fitted distributions, whereas the boxes in the lower triangle contain a scatterplot of the variables in that row/column. Finally, the  $\kappa_i$  and  $\mu_i$  parameters are represented in the first column. We do not observe significant differences between species, except for the anomalous low mean for the second angle for the rat in Figure 7a.

## 5. CONCLUSION AND FUTURE WORK

In this paper, we have introduced a simplified formulation of the log pseudo-likelihood for the multivariate von Mises distribution that it is also faster to compute, reducing significantly the time needed to obtain the maximum pseudo-likelihood estimator. We have also proposed a penalization term based on the Frobenius norm of the  $\mathbf{P}$  matrix that allows to penalize each parameter independently. We have proved the consistency of the penalized estimator, and we have also provided some recommendations to set the penalization parameters. Additionally, we have defined a multivariate circular distance that is easy to compute and takes into account the periodicity of circular data and then used it to compute a nonparametric approximation of the KL divergence between two circular samples. The penalized estimator has been tested in both synthetic and real-world experiments using the KL approximation with the circular distance in the latter case. The results prove that the penalized estimator provides similar or better estimations, specially when the number of samples is low, and we have a reasonable prior knowledge of the dependencies between the variables.

An analytical study of the properties of the maximum pseudo-likelihood estimator could provide consistent rules to set the penalization parameters  $\Phi$  and  $H$ , whereas the study on efficient approximations of the normalization term of the multivariate von Mises distribution will help to apply this same penalization to the maximum likelihood instead. Another major work line in the future could be to analyze the convergence rate of the penalized estimator compared to the nonregularized estimator. Additionally, it would be interesting to work on a more efficient sampling method based on rejection sampling in high-dimensional settings, for example, by using a mixture of multivariate normal distributions as reference.

All methods described in this paper have been included in the R package *mvCircular*<sup>a</sup> that implements sampling and fitting of the multivariate von Mises distribution as well as multivariate circular plots and statistics.

### Acknowledgments

We thank the referees for their review and comments, which significantly contributed to improving the quality of this paper.

This work was partially supported by the Spanish Ministry of Science and Innovation (MCINN) via a doctoral grant to the first author (FPU014/04818), the Spanish Ministry of Economy and Competitiveness (MINECO) through the Cajal Blue Brain (C080020-09; the Spanish partner of the Blue Brain initiative from EPFL) and TIN2013-41592-P projects, by the Regional Government of Madrid through the S2013/ICE-2845-CASI-CAM-CM project, and by the European Union's Seventh Framework Programme (FP7/2007-2013) under grant agreement no. 604102 (Human Brain Project).

<sup>a</sup><https://github.com/lrodriguezlujan/mvcircular>.

## References

1. Zhang J, Chowdhury S, Messac A, Castillo L. A multivariate and multimodal wind distribution model. *Renew Energy* 2013;51:436–447.
2. Masseran N, Razali AM, Ibrahim K, Latif MT. Fitting a mixture of von Mises distributions in order to model data on wind direction in peninsular Malaysia. *Energy Convers Manage* 2013;72:94–102.
3. Bahlmann C. Directional features in online handwriting recognition. *Pattern Recognit* 2006;39(1):115–125.
4. Baltieri D, Vezzani R, Cucchiara R. People orientation recognition by mixtures of wrapped distributions on random trees. In: *Lecture Notes in Computer Science, Vol 7576*. Berlin:Springer;2012. pp 270–283.
5. Masden EA, Reeve R, Desholm M, Fox AD, Furness RW, Haydon DT. Assessing the impact of marine wind farms on birds through movement modelling. *J R Soc Interface* 2012;9(74):2120–2130.
6. Lagona F. Regression analysis of correlated circular data based on the multivariate von Mises distribution. *Environ Ecol Stat* 2015;23:1–25.
7. Mardia KV, Hughes G, Taylor CC, Singh H. A multivariate von Mises distribution with applications to bioinformatics. *Can J Stat* 2008;36(1):99–109.
8. Razavian N, Kamisetty H, Langmead CJ. The von Mises graphical model: regularized structure and parameter learning. Tech. Rep. CMU-CS-11-108, Department of Computer Science, Carnegie Mellon University, Pittsburgh, PA; 2011.
9. Mardia KV, Kent JT, Zhang Z, Taylor CC, Hamelryck T. Mixtures of concentrated multivariate sine distributions with applications to bioinformatics. *J Appl Stat* 2012;39(11):2475–2492.
10. Banerjee A, Dhillon IS, Ghosh J, Sra S. Clustering on the unit hypersphere using von Mises–Fisher distributions. *J Mach Learn Res* 2005;6:1345–1382.
11. Fisher NI. *Statistical analysis of circular data*. Cambridge, UK: Cambridge University Press; 1993.
12. Mardia KV, Jupp PE. *Directional statistics*. Hoboken, NJ: Wiley; 2009.
13. Gatto R, Jammalamadaka SR. The generalized von Mises distribution. *Stat Methodol* 2007;4(3):341–353.
14. Kato S, Jones M. An extended family of circular distributions related to wrapped Cauchy distributions via Brownian motion. *Bernoulli* 2013;19(1):154–171.
15. Wang F, Gelfand AE. Directional data analysis under the general projected normal distribution. *Stat Methodol* 2013;10(1):113–127.
16. Hernandez-Stumpfhauser D, Breidt FJ, van der Woerd MJ. The general projected normal distribution of arbitrary dimension: modeling and Bayesian inference. *Bayesian Anal Advance Publication*, 19 January 2016.
17. Singh H, Hnizdo V, Demchuk E. Probabilistic model for two dependent circular variables. *Biometrika* 2002;89(3):719–723.
18. Tan KM, London P, Mohan K, Lee SI, Fazel M, Witten D. Learning graphical models with hubs. *J Mach Learn Res* 2014;15(1):3297–3331.
19. Bickel PJ, Levina E. Regularized estimation of large covariance matrices. *Annals Stat* 2008;36(1):199–227.
20. Rodriguez-Lujan L, Bielza C, Larrañaga P. Regularized multivariate von Mises distribution. In: *Lecture Notes in Computer Science, Vol 9422*. Berlin: Springer; 2015. pp 25–35.
21. Presnell B, Morrison SP, Littell RC. Projected multivariate linear models for directional data. *J Am Stat Assoc* 1998;93(443):1068–1077.
22. Fisher R. Dispersion on a sphere. *Proc R S London A: Math Phys Eng Sci* 1953;217:295–305.
23. Mardia KV, Voss J. Some fundamental properties of a multivariate von Mises distribution. *Commun Stat Theory Methods* 2014;43(6):1132–1144.

24. Best DJ, Fisher NI. Efficient simulation of the von Mises distribution. *Appl Stat* 1979;28:152–157.
25. Besag J. Statistical analysis of non-lattice data. *Statistician* 1975;24(3):179–195.
26. Liang P, Jordan MI. An asymptotic analysis of generative, discriminative, and pseudolikelihood estimators. In: *Proc 25th Int Conf on Machine Learning, Helsinki, FI, 5–9 July*. New York: ACM; 2008. pp 584–591.
27. Zhu C, Byrd RH, Lu P, Nocedal J. Algorithm 778: L-BFGS-B: Fortran subroutines for large-scale bound-constrained optimization. *ACM Trans Math Softw* 1997;23(4):550–560.
28. Anderson E, Bai Z, Bischof C, Blackford S, Dongarra J, Du Croz J, Greenbaum A, Hammarling S, McKenney A, Sorensen DC. *LAPACK users' guide*, 3rd ed. Philadelphia, PA: Society for Industrial and Applied Mathematics; 1999.
29. Fessler JA. Mean and variance of implicitly defined biased estimators (such as penalized maximum likelihood): applications to tomography. *IEEE Trans Image Process* 1996;5(3):493–506.
30. Anderson TW. *An introduction to multivariate statistical analysis*, 3rd ed. Hoboken, NJ: Wiley; 2003.
31. Gelman A, Carlin JB, Stern HS. *Bayesian data analysis*, 2nd ed. Chapman & Hall/CRC: Boca Raton, FL, USA; 2004.
32. Qi H, Sun D. An augmented Lagrangian dual approach for the h-weighted nearest correlation matrix problem. *IMA J Num Anal* 2011;31(2):491–511.
33. Rall W, Segev I, Rinzel J, Shepherd GM. *The theoretical foundation of dendritic function*. Cambridge, MA: MIT Press; 1995.
34. Pérez-Cruz F. Kullback-Leibler divergence estimation of continuous distributions. In: *IEEE Int Symp on Information Theory, Toronto, ON, 6–11 July*. IEEE; 2008. pp 1666–1670.
35. Ascoli GA, Donohue DE, Halavi M. *Neuromorpho.org: a central resource for neuronal morphologies*. *J Neurosci* 2007;27(35):9247–9251.
36. Jacobs B, Schall M, Prather M, Kapler E, Driscoll L, Baca S, Jacobs J, Ford K, Wainwright M, Trembl M. Regional dendritic and spine variation in human cerebral cortex: a quantitative Golgi study. *Cerebral Cortex* 2001;11(6):558–571.
37. Ballesteros-Yáñez I, Valverde O, Ledent C, Maldonado R, DeFelipe J. Chronic cocaine treatment alters dendritic arborization in the adult motor cortex through a CB1 cannabinoid receptor-dependent mechanism. *Neuroscience* 2007;146(4):1536–1545.
38. Markram H, Muller E, Ramaswamy S, Reimann MW, Abdellah M, Sanchez CA, Ailamaki A, Alonso-Nanclares L, Antille N, Arsever S, Kahou GA, Berger TK, Bilgili A, Buncic N, Chalimourda A, Chindemi G, Courcol JD, Delalondre F, Delattre V, Druckmann S, Dumusc R, Dynes J, Eilemann S, Gal E, Gevaert ME, Ghobril JP, Gidon A, Graham JW, Gupta A, Haenel V, Hay E, Heinis T, Hernando JB, Hines M, Kanari L, Keller D, Kenyon J, Khazen G, Kim Y, King JG, Kisvarday Z, Kumbhar P, Lasserre S, Le Bé JV, Magalhães BR, Merchán-Pérez A, Meystre J, Morrice BR, Muller J, Muñoz-Céspedes A, Muralidhar S, Muthurasa K, Nachbaur D, Newton TH, Nolte M, Ovcharenko A, Palacios J, Pastor L, Perin R, Ranjan R, Riachi I, Rodríguez JR, Riquelme JL, Rössert C, Sfyarakis K, Shi Y, Shillcock JC, Silberberg G, Silva R, Tauheed F, Telefont M, Toledo-Rodriguez M, Tränkler T, Van Geit W, Díaz JV, Walker R, Wang Y, Zaninetta SM, DeFelipe J, Hill SL, Segev I, Schürmann F. Reconstruction and simulation of neocortical microcircuitry. *Cell* 2015;163(2):456–492.
39. Ramaswamy S, Courcol JD, Abdellah M, Adaszewski SR, Antille N, Arsever S, Atenekeng G, Bilgili A, Brukay Y, Chalimourda A, Chindemi G, Delalondre F, Dumusc R, Eilemann S, Gevaert ME, Gleeson P, Graham JW, Hernando JB, Kanari L, Katkov Y, Keller D, King JG, Ranjan R, Reimann MW, Rössert C, Shi Y, Shillcock JC, Telefont M, Van Geit W, Diaz JV, Walker R, Wang Y, Zaninetta SM, DeFelipe J, Hill SL, Muller J, Segev I, Schürmann F, Muller EB, Markram H. The neocortical microcircuit collaboration portal: a resource for rat somatosensory cortex. *Front Neural Circuits* 2015;9(44).
40. Demšar J. Statistical comparisons of classifiers over multiple data sets. *J Mach Learn Res* 2006;7:1–30.