

Bayesian optimization of the PC algorithm for learning Gaussian Bayesian networks

Irene Córdoba¹, Eduardo C. Garrido-Merchán², Daniel Hernández-Lobato²,
Concha Bielza¹, and Pedro Larrañaga¹

¹ Universidad Politécnica de Madrid, Departamento de Inteligencia Artificial

² Universidad Autónoma de Madrid, Departamento de Ingeniería Informática

Abstract. The PC algorithm is a popular method for learning the structure of Gaussian Bayesian networks. It carries out statistical tests to determine absent edges in the network. It is hence governed by two parameters: (i) The type of test, and (ii) its significance level. These parameters are usually set to values recommended by an expert. Nevertheless, such an approach can suffer from human bias, leading to suboptimal reconstruction results. In this paper we consider a more principled approach for choosing these parameters in an automatic way. For this we optimize a reconstruction score evaluated on a set of different Gaussian Bayesian networks. This objective is expensive to evaluate and lacks a closed-form expression, which means that Bayesian optimization (BO) is a natural choice. BO methods use a model to guide the search and are hence able to exploit smoothness properties of the objective surface. We show that the parameters found by a BO method outperform those found by a random search strategy and the expert recommendation. Importantly, we have found that an often overlooked statistical test provides the best over-all reconstruction results.

1 Introduction

Graphical models serve as a compact representation of the relationships between variables in a domain. An important subclass is the Bayesian network, where conditional independences are encoded by missing edges in a directed graph with no cycles. By exploiting these independences, Bayesian networks yield a modular factorization of the joint probability distribution underlying the data. Of particular interest are Gaussian Bayesian networks for modelling variables in a continuous domain, which have been widely applied in real scenarios such as gene network discovery [11] and neuroscience [1].

When learning graphical models from data, two main tasks are usually differentiated: structure and parameter learning. The former consists in recovering the graph structure, and the latter amounts to fitting the numerical quantities in the model. In Gaussian Bayesian networks, parameter learning involves using standard linear regression theory, whereas structure learning is not an easy task in general, given the combinatorial search space of acyclic digraphs. There

are two main approaches one can find in the literature for structure discovery in Bayesian networks: score-and-search heuristics, where the search space is explored looking for the network which optimizes a given score function, and constraint-based approaches, where statistical tests are performed in order to include or exclude dependencies between variables.

A popular constraint-based method with consistency guarantees is the PC algorithm [6,16]. In this method, a backward stepwise testing procedure is performed for determining absent edges in the resulting graph. Thus, of critical importance are the choice of the statistical test to be performed, and the significance level at which the potential edges are going to be tested. However, both are usually fixed after a grid search or directly set by expert knowledge [2,6]. In the literature on Bayesian network structure learning some empirical studies explore exact structure recovery [9], the behavior of score-and-search algorithms [8], and the impact of the significance level in the PC algorithm for high dimensional sparse scenarios [2,6]. We are not aware, however, of any research work using elaborated methods for hyper-parameters selection in this context.

In this paper we show that Bayesian optimization (BO) can be used as an alternative methodology for choosing the significance level and the statistical test in the PC algorithm. BO has been recently applied successfully in different optimization problems [14,15]. We consider here a structure learning scenario in moderately sparse settings that is representative of those considered in [6]. We show that BO outperforms, in terms of structure recovery error, in a relatively small number of iterations, both a baseline approach based on a grid search and specific values set by expert knowledge obtained from previous results on this problem [6]. Furthermore, we also analyze what values for the statistical test and the significance level are recommended by the BO approach, and compare them with those often used by the relevant literature on the subject.

This article is organized as follows. In Section 2, we introduce the main concepts relative to Gaussian Bayesian networks that will be used throughout the rest of the paper. Then, in Section 3, we describe the PC algorithm, emphasizing its hyper-parameters and how they may affect its performance. Black box BO is outlined in Section 4, with emphasis on the particular characteristics of our problem. The experimental setting as well as the results we have obtained are described in Section 5. Finally, we conclude the paper in Section 6, where we also point out the main planned lines of future work.

2 Preliminaries on Gaussian Bayesian networks

Throughout the remainder of the paper, X_1, \dots, X_p will denote p random variables, and \mathbf{X} the random vector they form. For a subset of indices $I \subseteq \{1, \dots, p\}$, \mathbf{X}_I will denote the random vector corresponding only to the variables indexed by I . We will use $\mathbf{X}_I \perp\!\!\!\perp \mathbf{X}_J \mid \mathbf{X}_K$ for denoting that \mathbf{X}_I is conditionally independent of \mathbf{X}_J given the values of \mathbf{X}_K , being I, K, J disjoint subsets of $\{1, \dots, p\}$. Let $G = (V, E)$ be an acyclic digraph, where $V = \{1, \dots, p\}$ is the vertex set and $E \subseteq V \times V$ is the edge set. When G is part of a graphical model, its vertex

set V can be thought of as indexing a random vector $\mathbf{X} = (X_1, \dots, X_p)$. In a Bayesian network, the graph G is constrained to be acyclic directed and with no multiple edges.

A common interpretation of edges in a Bayesian network is the ordered Markov property, although many more exist, which can be shown to be equivalent [7]. This property is stated as follows. For a vertex $i \in V$, the set of parents of i is defined as $\text{pa}(i) := \{j : (j, i) \in E\}$. In every acyclic digraph, an ancestral order \prec can be found between the nodes where it is satisfied that if $j \in \text{pa}(i)$, then $j \prec i$, that is, the parents of a vertex come before it in the ancestral order. For notational simplicity, in the remainder we will assume that the vertex set $V = \{1, \dots, p\}$ is already ancestrally ordered. The ordered Markov property of Bayesian networks can be written in this context as

$$X_i \perp\!\!\!\perp \mathbf{X}_{\{1, \dots, i-1\} \setminus \text{pa}(i)} \mid \mathbf{X}_{\text{pa}(i)}$$

for all $i \in V$.

The above conditional independences, together with the properties of the multivariate Gaussian distribution, allow to express a Gaussian Bayesian network as a system of recursive linear regressions. Indeed, if for each $i \in V = \{1, \dots, p\}$, we consider the regression of X_i on its predecessors in the ancestral order, X_1, \dots, X_{i-1} , then from the results regarding conditioning on multivariate Gaussian random variables we obtain

$$X_i = \sum_{j=1}^{i-1} \beta^{ji|1, \dots, i-1} X_j + \epsilon_i, \quad (1)$$

where the regression coefficient $\beta^{ji|1, \dots, i-1} = 0$ when $j \notin \text{pa}(i)$, and ϵ_i are independent Gaussian random variables with zero mean and variance equal to the conditional variance of X_i on X_1, \dots, X_{i-1} . Therefore, both the structure and parameters of a Gaussian Bayesian network can be directly read off from the system of linear regressions in Equation (1).

3 Structure learning with the PC algorithm

The PC algorithm for learning Gaussian Bayesian networks proceeds by first estimating the skeleton, that is, the underlying undirected graph, of the acyclic digraph, and then orienting it. That is, for each vertex $i \in V = \{1, \dots, p\}$, it looks through the set of its neighbors, which we will denote as $\text{ne}(i)$, and selects a node $j \in \text{ne}(i)$ and subset $C \subseteq \text{ne}(i) \setminus \{j\}$. Then, the conditional independence $X_i \perp\!\!\!\perp X_j \mid \mathbf{X}_C$ is tested on the available data. It is a backward stepwise elimination method, in the sense that it starts with the complete undirected graph, and then proceeds by testing conditional independences in order to remove edges, doing so incrementally in the size of the neighbor subset C .

The PC main phase pseudocode can be found in Algorithm 1. The output of Algorithm 1 is the skeleton, or undirected version, of the estimated Gaussian Bayesian network, which is later oriented. Algorithm 1 is typically called

Algorithm 1 The PC algorithm in its population version

Input: Conditional independence information about $\mathbf{X} = (X_1, \dots, X_p)$ **Output:** Skeleton of the Gaussian Bayesian network

```

1:  $G \leftarrow$  complete undirected graph on  $\{1, \dots, p\}$ 
2:  $l \leftarrow -1$ 
3: repeat
4:    $l \leftarrow l + 1$ 
5:   repeat
6:     Select  $i$  such that  $(i, j) \in E$  and  $|\text{ne}(i) \setminus \{j\}| \geq l$ 
7:     repeat
8:       Choose new  $C \subseteq \text{ne}(i) \setminus \{j\}$  with  $|C| = l$ 
9:       if  $X_i \perp\!\!\!\perp X_j \mid \mathbf{X}_C$  then
10:         $E \leftarrow E \setminus \{(i, j), (j, i)\}$ 
11:       end if
12:     until  $(i, j)$  has been deleted or all neighbor subsets of size  $l$  have been tested
13:   until All  $(i, j) \in E$  such that  $|\text{ne}(i) \setminus \{j\}| \geq l$  have been tested
14: until  $|\text{ne}(i) \setminus \{j\}| < l$  for all  $(i, j) \in E$ 

```

the *population* version of the PC algorithm [6], since it assumes that perfect information is available about the conditional independence relationships present in the data. This is useful for illustrating the behavior and main properties of the algorithm; however, in real scenarios this is unrealistic, and statistical tests must be performed on the data in order to determine which variable pairs, with respect to different node subsets, are conditionally independent.

3.1 Significance level and statistical test

The criteria for removing edges is related to the ordered Markov property and Equation (1). In particular, from multivariate Gaussian analysis we know that for $i \in V$ and $j < i$,

$$\beta^{ji|1, \dots, i-1} = 0 \iff \rho^{ji|1, \dots, i-1} = 0,$$

where $\rho^{ji|1, \dots, i-1}$ denotes the partial correlation coefficient between X_i and X_j with respect to X_1, \dots, X_{i-1} . In the PC algorithm, at iteration l , the null hypothesis $H_0 : \rho^{ji|C} = 0$ is tested against the alternative hypothesis $H_1 : \rho^{ji|C} \neq 0$, where C is a subset of the neighbors of i (excluding j) in the current estimator of the skeleton such that $|C| = l$.

The significance level at which H_0 will be tested, which we will denote in the remainder as α , is typically smaller or equal than 0.05, and serves to control the type I error. The other parameter of importance is the statistical test itself. The usual choice for this is a Gaussian test based on the Fisher's Z transform of the partial correlation coefficient [2,6], which is asymptotically normal. However, there are other choices available in the literature that could be considered and can be found in standard implementations of the algorithm. For example, the

`bnlearn` R package [13] provides the standard Student’s t test for the untransformed partial correlation coefficient, and the χ^2 test and a test based on the shrinkage James-Stein estimator, for the mutual information [4].

3.2 Evaluating the quality of the learned structure

When performing structure discovery in graphical models, there are several ways of evaluating the results obtained by an algorithm. As a starting point, one could use standard error rates, such as the true positive and false positive rates. These rates simply take into account the original acyclic digraph $G = (V, E)$, and the estimated one \hat{G} , with edge set \hat{E} . Then, E with \hat{E} are compared element-wise. This is a common approach in Bayesian networks.

We have preferred however to use the Structural Hamming Distance (SHD) [17]. This measure is motivated as follows. In Bayesian networks, there is not a unique correspondence between the model and the acyclic digraph that represents it. That is, if we denote as $\mathcal{M}(G)$ the set of multivariate Gaussian distributions whose conditional independence model is compatible (in the sense of the pairwise Markov property and Equation (1)) with the acyclic digraph G , then we may have two distinct acyclic digraphs G_1 and G_2 such that $\mathcal{M}(G_1) = \mathcal{M}(G_2)$. In such case, G_1 and G_2 are said to be Markov equivalent.

The SHD measure between two acyclic digraph structures G_1 and G_2 takes into account this issue of non unique correspondence. In particular, it counts the number of operations that have to be performed in order to transform the Markov equivalence class of one graph into the other. Thus, given two acyclic digraphs that are distinct but Markov equivalent, their true positive and false positive rates could be nonzero, while their SHD is guaranteed to be zero.

4 Black-box Bayesian optimization

Denote the SHD objective function as $f(\boldsymbol{\theta})$, which depends of the parameters in the PC algorithm, $\boldsymbol{\theta} = (\alpha, T)$, that are going to be optimized, α , the significance level, and T , the independence test. We can view $f(\boldsymbol{\theta})$ as a black-box objective function with noisy evaluations $y_i = f(\boldsymbol{\theta}) + \epsilon_i$, with ϵ_i being a, typically, Gaussian noise term. With BO the number of evaluations of f needed to solve the optimization problem are drastically reduced. Let the observed data until step $t - 1$ of the algorithm be $\mathcal{D}_{t-1} = \{(\boldsymbol{\theta}_i, y_i)\}_{i=1}^{t-1}$. At iteration t of BO, a probabilistic model $p(f(\boldsymbol{\theta}) | \mathcal{D}_{t-1})$, typically a Gaussian process (GP) [12], is fitted to the data collected so far. The uncertainty about f provided by the probabilistic model is then used to generate an acquisition function $a_t(\boldsymbol{\theta})$, whose value at each input location indicates the expected utility of evaluating f there. Therefore, at iteration t , $\boldsymbol{\theta}_t$ is chosen as the one that maximizes the acquisition function. The described process is repeated until enough data about the objective has been collected. When this is the case, the GP predictive mean for $f(\cdot)$ can be optimized to find the solution of the optimization problem, or we can provide as a solution the best observation made so far.

The key for BO success is that evaluating the acquisition function is very cheap compared to the evaluation of the actual objective, because it only depends on the GP predictive distribution for the objective at any candidate point. The GP predictive distribution for $f(\boldsymbol{\theta}_t)$, the candidate location for next iteration, is given by a Gaussian distribution characterized by a mean $\boldsymbol{\mu}$ and a variance σ^2 with values

$$\begin{aligned}\boldsymbol{\mu} &= \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}, \\ \sigma^2 &= k(\boldsymbol{\theta}_t, \boldsymbol{\theta}_t) - \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}_*.\end{aligned}\tag{2}$$

where $\mathbf{y} = (y_1, \dots, y_{t-1})^t$ is a vector with the objective values observed so far; σ_n^2 is the variance of the additive Gaussian noise; \mathbf{k}_* is a vector with the prior covariances between $f(\boldsymbol{\theta}_t)$ and each y_i ; \mathbf{K} is a matrix with the prior covariances among each y_i ; and $k(\boldsymbol{\theta}_t, \boldsymbol{\theta}_t)$ is the prior variance at the candidate location $\boldsymbol{\theta}_t$. The covariance function $k(\cdot, \cdot)$ is pre-specified; for further details about GPs and example of covariance functions we refer the reader to [12]. Four steps of the BO process are illustrated graphically in Fig. 1 for a toy minimization problem.

In BO methods the acquisition function balances between exploration and exploitation in an automatic way. A typical choice for this function is the information-theoretic method Predictive Entropy Search (PES) [5]. In PES, we are interested in maximizing information about the location of the optimum value, $\boldsymbol{\theta}^*$, whose posterior distribution is $p(\boldsymbol{\theta}^* | \mathcal{D}_{t-1})$. This can be done through the negative differential entropy measure of $p(\boldsymbol{\theta}^* | \mathcal{D}_{t-1})$. Through several operations, an approximation to PES is given by

$$a(\boldsymbol{\theta}) = H[p(y | \mathcal{D}_{t-1}, \boldsymbol{\theta})] - \mathbb{E}_{p(\boldsymbol{\theta}^* | \mathcal{D}_{t-1})}[H[p(y | \mathcal{D}_{t-1}, \boldsymbol{\theta}, \boldsymbol{\theta}^*)]],$$

where $p(y | \mathcal{D}_{t-1}, \boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is the posterior predictive distribution of y given \mathcal{D}_{t-1} and the minimizer $\boldsymbol{\theta}^*$ of f , and $H[\cdot]$ is the differential entropy. The first term of the previous equation can be analytically solved as it is the entropy of the predictive distribution and the second term is approximated by Expectation Propagation [10]. We can see an example of the PES acquisition function in Fig. 1.

5 Numerical experiments

Since we will consider networks of different node size p , we will use in our experimental setting as the validation measure a normalized version of SHD with respect to the maximum edge number $p(p-1)/2$. The significance level α will be represented for the BO algorithm as a real variable whose range lies in the decimal logarithmic space $[-5, -1]$. The statistical test will be represented using a categorical variable whose value indicates one of the above mentioned four tests. Namely, two test based on the partial correlation coefficient: a Gaussian test based on the Fisher's Z transform and the Student's T test; and two test based on the mutual information: the χ^2 test, and a test based on the shrinkage James-Stein estimator. As outlined before, this problem is specially suitable for BO, since we do not have access to gradients, the objective evaluations may be expensive and they may be contaminated with noise.

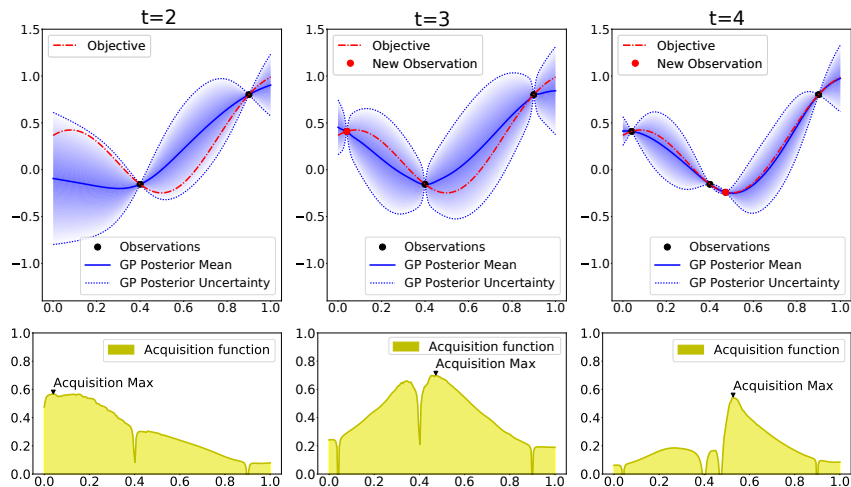


Fig. 1. An example of BO on a toy 1D noiseless problem, where a function (shown in red, dotted) sampled from a GP prior is optimized. The top figures show the GP (mean and standard deviation) estimation of the objective $f(\cdot)$ in blue. The acquisition function PES is shown in the bottom figures in yellow. The acquisition is high where the GP predicts a low objective and where the uncertainty is high. From left to right, iterations (t) of BO are executed and the points (black) suggested by the acquisition function are evaluated in the objective, we show in red the new point for every iteration. In a small number of iterations, the GP is able to almost exactly approximate the objective.

We have employed `Spearmint` (<https://github.com/HIPS/Spearmint>) for BO and the `pc.stable` function from the `bnlearn` R package [13] for the PC algorithm execution. We have run BO with the PES acquisition function over a set of Gaussian Bayesian networks generated following the simulation methodology of [6]. That is, the absent edges in the acyclic digraph G are sampled by using independent Bernoulli random variables with probability of success $d = n/(p-1)$, where p is the vertex number of G and n is the average neighbor size. The probability d can be thought of as an indicator of the density of the network: smaller d values mean sparser networks. The node size p is obtained from a grid of values $\{25, 50, 75, 100\}$, while the average neighbor size is $n \in \{2, 8\}$. Finally, we consider different sample sizes $N \in \{25, 50, 75, 100\}$. Therefore, we have a total of 32 different network learning scenarios, that are representative of those that can be found in [6]. We create 40 different replicas of the experiment and report average results across them, in order to provide more robust results. In each of these replicas, the nonzero regression coefficients in Equation (1) are sampled from a uniform distribution on $[0.1, 1]$, following [6].

For BO, we have used the PES acquisition function and 10 Monte Carlo iterations for sampling the parameters of the GP. The acquisition function is averaged across these 10 samples. We have used the Matérn covariance function

for $k(\cdot, \cdot)$ (Equation (2)) and the transformation described in [3] so that the GP can deal with the categorical variable (the test type). We compare BO with a random search (RS) strategy of the average normalized SHD error surface and with the expert criterion (EC), taken from [6]. These authors recommend a value of $\alpha = 0.01$ and use the Fisher’s Z partial correlation test. At each iteration, BO provides a candidate solution which corresponds to the best observation made so far. We stop the search in BO and RS after 30 evaluations of the objective.

The average normalized SHD results obtained are shown in Fig. 2. We show the relative difference in log-scale with respect to the best observed result. Therefore, the lower the values obtained, the better. We show the mean and standard deviation of this measure along the 40 replicas of the experiment, for each of the three methods compared (BO, RS and EC). We can see that EC is easily improved after only 10 iterations of BO and RS. Furthermore, BO outperforms RS providing significantly better results as more evaluations are performed. Importantly, the standard deviation of the results of BO are fairly small in the last iterations. This means that BO is very robust to the different replicas of the experiments.

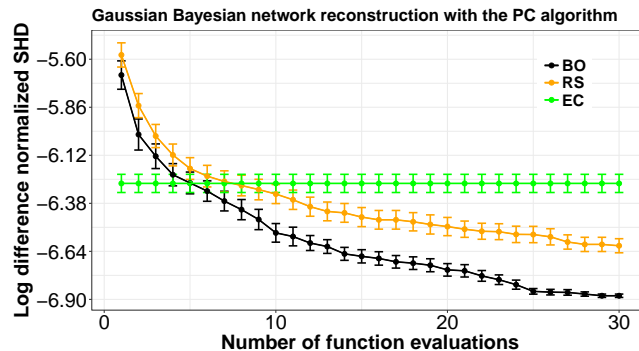


Fig. 2. Logarithmic difference with respect to the best observed average normalized SHD obtained in 40 replicas of the 32 considered Gaussian Bayesian networks.

Since the expert criterion is outperformed, we are interested in the parameter suggestions delivered by BO. In order to explore these results, we have generated two histograms that summarize the suggested parameters by BO in the last iteration, shown in Fig. 3. We observe that the most frequently recommended test is the James-Stein shrinkage estimator of the mutual information [4], while the most frequent recommendation for the significance level is concentrated at values lower than 0.025.

These results are very interesting from the viewpoint of graphical models learning. The first observation is that the optimal value obtained for the significance level is fairly close to the one suggested in [6]. However, the SHD results are arguably better for the BO than for the human expert. This may be explained

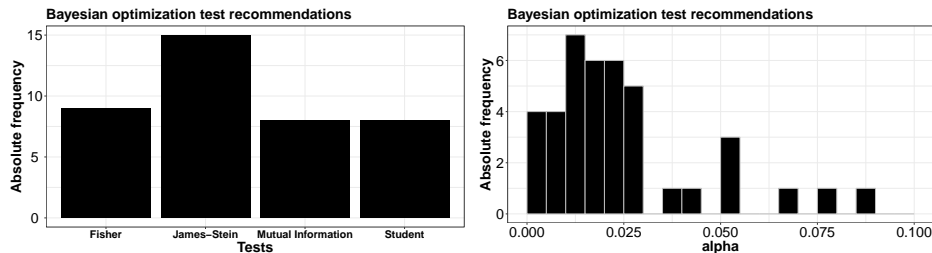


Fig. 3. Histograms with the recommended parameters by BO in the last iteration.

by the second interesting result we have obtained. Namely, the shrinkage James-Stein estimator of the mutual information is suggested more times than the that the extended Fisher’s Z partial correlation test. Therefore, in the context of sparse, high-dimensional networks, where we may have $p > N$ (such as in our experimental set-up and the one in [6]), it may be better to focus on the selection of the statistical test, rather than on carefully adjusting the significance level. In the literature, however, it is often done the other-way-around, and more effort is put on carefully adjusting the significance level.

6 Conclusions and future work

In this paper we have proposed the use of BO for selecting the optimal parameters of PC algorithm for structure recovery in Gaussian Bayesian networks. We have observed that, in a small number of iterations, the expert suggestion is outperformed by the recommendations provided by a BO method. Furthermore, an analysis of the recommendations made by the BO algorithm shows interesting results about the relative importance of the selection of the statistical test, as opposed to the selection of the significance level. In the literature, however, it is often that the selection of the significance level receives more attention.

For future work, we would like to apply BO in higher dimensional settings, where the number of nodes increases exponentially, whereas the number of samples increases linearly. This is also a typical scenario in Gaussian Bayesian network real applications. We would also like to explore how different error measures, such as the true positive and false positive rates, affect the obtained results when they are optimized using BO. Finally, we plan to extend this methodology to consider multi-objective optimization scenarios and also several constraints, since current BO methods are able to handle these problems too.

Acknowledgements: We acknowledge the use of the facilities of Centro de Computación Científica (CCC) at Universidad Autónoma de Madrid, and financial support from Comunidad de Madrid, grant S2013/ICE-2845; from the Spanish *Ministerio de Economía, Industria y Competitividad*, grants TIN2016-79684-P, TIN2016-76406-P, TEC2016-81900-REDT; from the Cajal Blue Brain project (C080020-09, the Spanish partner of the EPFL Blue Brain initiative); and from Fundación BBVA (Scientific Re-

search Teams in Big Data 2016). Irene Córdoba is supported by grant FPU15/03797 from the Spanish *Ministerio de Educación, Cultura y Deporte*.

References

1. Bielza, C., Larrañaga, P.: Bayesian networks in neuroscience: A survey. *Frontiers in Computational Neuroscience* **8**, 131 (2014)
2. Colombo, D., Maathuis, M.H.: Order-independent constraint-based causal structure learning. *Journal of Machine Learning Research* **15**(1), 3741–3782 (2014)
3. Garrido-Merchán, E.C., Hernández-Lobato, D.: Dealing with categorical and integer-valued variables in Bayesian optimization with Gaussian processes (2018), arXiv:1805.03463
4. Hausser, J., Strimmer, K.: Entropy inference and the James-Stein estimator, with application to nonlinear gene association networks. *Journal of Machine Learning Research* **10**, 1469–1484 (2009)
5. Hernández-Lobato, J.M., Hoffman, M.W., Ghahramani, Z.: Predictive entropy search for efficient global optimization of black-box functions. In: *Advances in Neural Information Processing Systems*. pp. 918–926 (2014)
6. Kalisch, M., Bühlmann, P.: Estimating high-dimensional directed acyclic graphs with the PC-algorithm. *Journal of Machine Learning Research* **8**, 613–636 (2007)
7. Lauritzen, S.L., Dawid, A.P., Larsen, B.N., Leimer, H.G.: Independence properties of directed Markov fields. *Networks* **20**(5), 491–505 (1990)
8. Malone, B., Järvisalo, M., Myllymäki, P.: Impact of learning strategies on the quality of bayesian networks: An empirical evaluation. In: *Proceedings of the Thirty-First the conference on Uncertainty in Artificial Intelligence*. pp. 562–571 (2015)
9. Malone, B., Kangas, K., Järvisalo, M., Koivisto, M., Myllymäki, P.: Empirical hardness of finding optimal Bayesian network structures: Algorithm selection and runtime prediction. *Machine Learning* **107**(1), 247–283 (2018)
10. Minka, T.P.: Expectation propagation for approximate bayesian inference. In: *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*. pp. 362–369 (2001)
11. Ness, R.O., Sachs, K., Vitek, O.: From correlation to causality: Statistical approaches to learning regulatory relationships in large-scale biomolecular investigations. *Journal of Proteome Research* **15**(3), 683–690 (2016)
12. Rasmussen, C.E.: Gaussian processes in machine learning. In: *Advanced Lectures on Machine Learning*, pp. 63–71. Springer (2004)
13. Scutari, M.: Learning Bayesian networks with the bnlearn R package. *Journal of Statistical Software* **35**(3), 1–22 (2010)
14. Shahriari, B., Swersky, K., Wang, Z., Adams, R.P., de Freitas, N.: Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE* **104**(1), 148–175 (2016)
15. Snoek, J., Larochelle, H., Adams, R.P.: Practical Bayesian optimization of machine learning algorithms. In: *Advances in Neural Information Processing Systems*. pp. 2951–2959 (2012)
16. Spirtes, P., Glymour, C., Scheines, R.: *Causation, Prediction, and Search*. MIT Press (2000)
17. Tsamardinos, I., Brown, L.E., Aliferis, C.F.: The max-min hill-climbing Bayesian network structure learning algorithm. *Machine Learning* **65**(1), 31–78 (2006)