

Nr. 11  
03. March 2014

Leopold-Franzens-Universität Innsbruck



Preprint-Series: Department of Mathematics - Applied Mathematics

## Multi-penalty regularization for detecting relevant variables

Katerina Hlavackova-Schindler, Valeriya Naumova and  
Sergiy Pereverzyev Jr.



Technikerstraße 21a - 6020 Innsbruck - Austria  
Tel.: +43 512 507 53803 Fax: +43 512 507 53898  
<https://www.applied-math.uibk.ac.at>

# Multi-penalty regularization for detecting relevant variables

Kateřina Hlaváčková-Schindler<sup>1</sup>, Valeriya Naumova<sup>2</sup> and  
Sergiy Pereverzyev Jr.<sup>3</sup>

March 2014

<sup>1</sup> *Chair of Bioinformatics, University of Natural Resources and Life Sciences,  
Vienna, Austria*  
e-mail: katerina.schindler@gmail.com

<sup>2</sup> *Simula Research Laboratory,  
Lysaker, Norway*  
e-mail: valeriya@simula.no

<sup>3</sup> *Applied Mathematics Group, Department of Mathematics,  
University of Innsbruck, Austria*  
e-mail: sergiy.pereverzyev@uibk.ac.at

## Abstract

In this paper we propose a new method for detecting relevant variables from a priori given high-dimensional data under the assumption that input-output dependence is described by a nonlinear function depending on a few variables. The method is based on the inspection of the behavior of discrepancies of a multi-penalty regularization with a component-wise penalization for small and large values of regularization parameters. We provide the justification of the proposed method under a certain condition on sampling operators. The effectiveness of the method is demonstrated in the example with synthetic data and in the reconstruction of gene regulatory networks. In the latter example, the obtained results provide a clear evidence of the competitiveness of the proposed method.

**Keywords:** multi-penalty regularization, variables detection, causality networks, gene regulatory networks.

# 1 Introduction and description of approach

Natural and social phenomena usually emerge from the behavior of complex systems consisting of interacting components or variables. In practice, we do not have a direct access to the “laws” governing the underlying relationships between them; instead, we are faced with a dataset recorded from the possibly interacting variables. How can we tell from these given data whether there exists any relationship between two or more variables?

This question can be made precise by considering a dataset

$$Z_N = \{ (x_1^i, x_2^i, \dots, x_p^i, y^i) \}_{i=1}^N$$

of observed values  $y^i$ ,  $i = 1, 2, \dots, N$ , of a variable of interest  $y$  paired with simultaneously observed values  $x_\nu^i$ ,  $\nu = 1, 2, \dots, p$ , of the variables  $x_1, x_2, \dots, x_p$  that possibly interact with  $y$ . Then the set  $Z_N$  is used to quantify how strong is the effect of  $x = (x_1, x_2, \dots, x_p)$  on  $y$ . An instance of this situation is the problem of reconstructing from the set  $Z_N$  a multivariate function  $y = f(x_{\nu_1}, x_{\nu_2}, \dots, x_{\nu_l})$  that depends only on a subset  $\{x_{\nu_i}\}_{i=1}^l$  of the variables  $\{x_\nu\}_{\nu=1}^p$  (very often,  $l$  is much smaller than  $p$ ). In this work, we are interested in detecting such relevant variables  $x_{\nu_i}$  from given data  $Z_N$ .

Note that the above problem has been extensively studied under the assumption that the target function  $f$  depends linearly on the relevant variables such that it admits the representation

$$f(x) = \sum_{j=1}^p \beta_j x_j$$

with only a few non-zero coefficients  $\beta_j$  for  $j = \nu_1, \nu_2, \dots, \nu_l$ . Under such assumption the problem of detecting the relevant variables from the data set  $Z_N$  can be reduced to the linear regression with a sparsity constraint. The latter one is now fairly well understood and can be solved efficiently by means of  $l_1$ -regularization. We refer the reader to the classical work by [5] and the more recent one [7] (see also references therein) for comprehensive treatments of this subject.

Despite the computational benefit of the linear regression, it should be noted that this model is too simple to be always appropriately matched to the underlying dynamics and may sometimes lead to a misspecification (see the discussion in the last section). A suitable alternative is to adopt the situation where the target function  $f$  depends nonlinearly on the relevant variables. This situation is much less understood, and in the literature it is mostly restricted to the so-called *additive models* [9, 20, 16] in which the target function is assumed to be the sum

$$f(x) = \sum_{j=1}^p f_j(x_j) \tag{1}$$

of nonlinear univariate functions  $f_j$  in some Reproducing Kernel Hilbert Spaces (RKHS)  $\mathcal{H}_j$  such that  $f_j \equiv 0$  for  $j \notin \{\nu_i\}_{i=1}^l$ .

In [2] it has been observed that detection of the relevant variables in the model (1) can be performed by using a technique from the *multiple kernel learning* [3, 13]. Then an estimator of the target function (1) can be constructed as the sum  $\sum_{j=1}^p f_j^\lambda(x_j)$  of

the minimizers of the functional

$$T_\lambda^q(f_1, f_2, \dots, f_p; Z_N) = \frac{1}{N} \sum_{i=1}^N \left( y^i - \sum_{j=1}^p f_j(x_j^i) \right)^2 + \sum_{j=1}^p \lambda_j \|f_j\|_{\mathcal{H}_j}^q, \quad (2)$$

i.e.,

$$T_\lambda^q(f_1^\lambda, f_2^\lambda, \dots, f_p^\lambda; Z_N) = \min\{ T_\lambda^q(f_1, f_2, \dots, f_p; Z_N), f_j \in \mathcal{H}_j, j = 1, 2, \dots, p \},$$

where  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_p)$  is a vector of the regularization parameters, and  $q > 0$ .

A different approach has been recently proposed in [16]. This approach is based on the idea that the importance of a variable can be captured by partial derivatives. Then in [16] the target function is estimated as the minimizer of the functional

$$\hat{T}_\lambda(f; Z_N) = \frac{1}{N} \sum_{i=1}^N (y^i - f(x^i))^2 + \lambda_1 \|f\|_{\mathcal{H}}^2 + \lambda_2 \sum_{j=1}^p \left( \frac{1}{N} \sum_{i=1}^N \left( \frac{\partial f(x^i)}{\partial x_j} \right)^2 \right)^{1/2}, \quad (3)$$

where  $x^i = (x_1^i, x_2^i, \dots, x_p^i)$ , and  $\mathcal{H}$  is some RKHS of functions  $f = f(x_1, x_2, \dots, x_p)$ .

Note that the choice of the regularization parameters  $\lambda_j$  is an open issue in the both above mentioned approaches. For the multiple kernel learning scheme based on (2), an a priori parameter choice strategy has been proposed in [13]. In this strategy the choice of  $\lambda_j$  depends only on kernels generating RKHS  $\mathcal{H}_j$  and on distribution of the points  $x_j^i$  involved in  $Z_N$ . It is clear that such a strategy may not be suitable for detecting relevant variables, because the functions (1) depending on different variables  $x_j$  may be associated with the same  $\mathcal{H}_j$  and  $x_j^i$ . As to the scheme based on (3), no recipe for choosing the parameters  $\lambda_1, \lambda_2$  was given.

Observe also that a numerical implementation of the above mentioned approaches can be non trivial. For example, the functional (3), as well as the functional (2) with  $q \in (0, 1]$ , is not differentiable and, hence, its minimization cannot be done by simple gradient methods. Moreover, the minimizers of the functionals can only be computed in an iterative fashion requiring the solution of a system of  $M = O(Np)$  equations at each step, and this can be computationally expensive for large  $N$  and/or  $p$ .

In the present paper we propose a new approach attempting to detect relevant variables one by one such that the dimension of the corresponding system of equations increases only when it is necessary. The first step of our approach consists in constructing the minimizers  $f_j = f_j^{\lambda_j}(x_j)$  of the functionals  $T_{\lambda_j}^q(f_j; Z_N)$  defined by (2) with  $q = 2$ ,  $p = 1$ ,  $\lambda_1 = \lambda_j$ ,  $x_1^i = x_j^i$ ,  $\mathcal{H}_1 = \mathcal{H}_j$ ,  $j = 1, 2, \dots$ . From the representer theorem [11, 22], it follows that such minimization is reduced to solving systems of  $N$  linear equations. Then, the minimizers  $f_j^{\lambda_j}(x_j)$  are used to rank the variables  $x_j$  according to the values of the discrepancies

$$\mathcal{D}(f_j^{\lambda_j}(x_j); Z_N) = \left( \frac{1}{N} \sum_{i=1}^N \left( y^i - f_j^{\lambda_j}(x_j^i) \right)^2 \right)^{1/2}, \quad j = 1, 2, \dots,$$

as follows: the smaller the value of  $\mathcal{D}(f_j^{\lambda_j}(x_j); Z_N)$ , the higher the rank of  $x_j$ . This step can be seen as an attempt to interpret the data  $Z_N$  by using only a univariate function, and the variable with the highest rank is considered as the first relevant variable  $x_{\nu_1}$ .

The next step consists in testing the hypothesis that a variable with the second highest rank, say  $x_\mu$ , is also the relevant one. For such a testing we consider the minimizers  $f_{\nu_1}^{\lambda_{\nu_1}}$ ,  $f_\mu^{\lambda_\mu}$  of the functional

$$T_\lambda^2(f_{\nu_1}, f_\mu; Z_N) = \frac{1}{N} \sum_{i=1}^N (y^i - f_{\nu_1}(x_{\nu_1}^i) - f_\mu(x_\mu^i))^2 + \lambda_{\nu_1} \|f_{\nu_1}\|_{\mathcal{H}_{\nu_1}}^2 + \lambda_\mu \|f_\mu\|_{\mathcal{H}_\mu}^2. \quad (4)$$

Our idea is based on the observation [17] that in multi-penalty regularization with a component-wise penalization, such as (4), one needs to use small as well as large values of the regularization parameters  $\lambda_{\nu_1}$ ,  $\lambda_\mu$ . So, in the proposed approach the variable  $x_\mu$  is considered as the relevant one if for  $\{\lambda_{\nu_1}, \lambda_\mu\} \subset (0, 1)$ , the values of the discrepancy

$$\mathcal{D}(f_{\nu_1}^{\lambda_{\nu_1}}, f_\mu^{\lambda_\mu}; Z_N) = \left( \frac{1}{N} \sum_{i=1}^N \left( y^i - f_{\nu_1}^{\lambda_{\nu_1}}(x_{\nu_1}^i) - f_\mu^{\lambda_\mu}(x_\mu^i) \right)^2 \right)^{1/2} \quad (5)$$

are smaller than the ones for  $\lambda_{\nu_1} \in (0, 1)$ ,  $\lambda_\mu > 1$ . If it is not the case, then the above mentioned hypothesis is rejected, and in the same way we test the variable with the third highest rank, and so on. In the next section we provide a theoretical justification for the use of the values of the discrepancies corresponding to regularization parameters from different intervals for detecting the relevant variables.

On the other hand, if the variable  $x_\mu$  has been accepted as the second relevant variable, i.e.,  $x_\mu = x_{\nu_2}$ , then to test whether or not the variable with the third highest rank, say  $x_\nu$ , can be taken as the third relevant variable, i.e., whether or not  $x_\nu = x_{\nu_3}$ , we use the minimizers  $f_{\nu_1}^{\lambda_{\nu_1}}$ ,  $f_{\nu_2}^{\lambda_{\nu_2}}$ ,  $f_\nu^{\lambda_\nu}$  of the functional

$$T_\lambda^2(f_{\nu_1}, f_{\nu_2}, f_\nu; Z_N) = \frac{1}{N} \sum_{i=1}^N (y^i - f_{\nu_1}(x_{\nu_1}^i) - f_{\nu_2}(x_{\nu_2}^i) - f_\nu(x_\nu^i))^2 + \lambda_{\nu_1} \|f_{\nu_1}\|_{\mathcal{H}_{\nu_1}}^2 + \lambda_{\nu_2} \|f_{\nu_2}\|_{\mathcal{H}_{\nu_2}}^2 + \lambda_\nu \|f_\nu\|_{\mathcal{H}_\nu}^2, \quad (6)$$

where, with a little abuse of notation, we use the same symbols  $f_{\nu_1}$ ,  $f_{\nu_1}^{\lambda_{\nu_1}}$  as in (4),(5). Then as above the variable  $x_\nu$  is considered as the relevant one if for  $\{\lambda_{\nu_1}, \lambda_{\nu_2}, \lambda_\nu\} \subset (0, 1)$ , the values of the discrepancy

$$\mathcal{D}(f_{\nu_1}^{\lambda_{\nu_1}}, f_{\nu_2}^{\lambda_{\nu_2}}, f_\nu^{\lambda_\nu}; Z_N) = \left( \frac{1}{N} \sum_{i=1}^N \left( y^i - f_{\nu_1}^{\lambda_{\nu_1}}(x_{\nu_1}^i) - f_{\nu_2}^{\lambda_{\nu_2}}(x_{\nu_2}^i) - f_\nu^{\lambda_\nu}(x_\nu^i) \right)^2 \right)^{1/2} \quad (7)$$

are smaller than for  $\{\lambda_{\nu_1}, \lambda_{\nu_2}\} \subset (0, 1)$ ,  $\lambda_\nu > 1$ . If it is not the case, then the variable with the next highest rank is tested in the same way.

If the discrepancy (7) does exhibit the above mentioned behavior, then for testing the variable with the next highest rank in accordance with the proposed approach, we need to add to (6) one more penalty term corresponding to that variable, so that the functional  $T_\lambda^2(f_1, f_2, \dots, f_p; Z_N)$  of the form (2) containing the whole set of penalties may appear only at the end of the testing procedure.

In the next sections after presenting the theoretical background, we will illustrate the application of the proposed approach to the recovery of causal relationships in a gene regulatory network, and compare it with the results known from the literature.

## 2 Theoretical background

At first, we will write a system of necessary conditions for the minimizers of the functional (2), where, according to the proposed approach,  $p$  may take values  $1, 2, \dots$ , and  $q = 2$ .

Let  $\mathbb{R}^N$  be the  $N$ -dimensional Euclidean space of vectors  $u = (u^1, u^2, \dots, u^N)$  equipped with the norm  $\|u\|_{\mathbb{R}^N} := \left( \frac{1}{N} \sum_{i=1}^N (u^i)^2 \right)^{1/2}$  and the corresponding inner product  $\langle \cdot, \cdot \rangle_{\mathbb{R}^N}$ .

Consider the sampling operators  $S_{N,j}$  mapping the RKHS  $\mathcal{H}_j$  generated by the kernels  $K_j = K_j(x_j, v_j)$ ,  $j = 1, 2, \dots, p$ , into  $\mathbb{R}^N$  such that for  $f \in \mathcal{H}_j$ ,

$$S_{N,j}f = (f(x_j^1), f(x_j^2), \dots, f(x_j^N)) \in \mathbb{R}^N.$$

In view of the reproducing property  $f(x_j^i) = \langle K_j(x_j^i, \cdot), f(\cdot) \rangle_{\mathcal{H}_j}$  of the kernels  $K_j$ , we can write the adjoints  $S_{N,j}^* : \mathbb{R}^N \rightarrow \mathcal{H}_j$  of the sampling operators as follows

$$(S_{N,j}^*u)(x_j) = \frac{1}{N} \sum_{i=1}^N K_j(x_j^i, x_j) u^i. \quad (8)$$

In terms of  $S_{N,j}$ , the functional (2) has the form

$$T_\lambda^2(f_1, f_2, \dots, f_p; Z_N) = \left\| Y - \sum_{j=1}^p S_{N,j} f_j \right\|_{\mathbb{R}^N}^2 + \sum_{j=1}^p \lambda_j \|f_j\|_{\mathcal{H}_j}^2, \quad (9)$$

where  $Y = (y^1, y^2, \dots, y^N)$ . Then, using the standard technique of the calculus of variations, we obtain the following system of equations for the minimizers  $f_j^{\lambda_j}$

$$\lambda_j f_j^{\lambda_j} + \sum_{\nu=1}^p S_{N,j}^* S_{N,\nu} f_\nu^{\lambda_\nu} = S_{N,j}^* Y, \quad j = 1, 2, \dots, p. \quad (10)$$

From (8) and (10), it is clear that  $f_j^{\lambda_j}$  can be represented as

$$f_j^{\lambda_j}(x_j) = \sum_{i=1}^N \gamma_i^j K_j(x_j^i, x_j). \quad (11)$$

Note that (11) can be seen as an analog of the well-known representer theorem [11, 22] for the case of the regularization with a component-wise penalization in RKHS. This allows the reduction of the minimization of (9) to solving systems of  $Np$  linear equations with respect to  $\gamma_i^j$ . Recall that in the approach described above,  $p$  will successively take the values  $1, 2, \dots$ , such that the dimension of the corresponding system (10) increases only when it is necessary.

Now for the sake of definiteness and simplicity of the presentation, suppose that

$$Y = S_{N,1}f_1 + S_{N,2}f_2 + \varepsilon, \quad (12)$$

where  $f_1 = f_1(x_1)$ ,  $f_2 = f_2(x_2)$ , and the vector  $\varepsilon \in \mathbb{R}^N$  may represent a noise in measurements, as well as a contribution to the data  $Y$  coming from functions of other relevant variables. Note that (12) means that  $x_1, x_2$  are relevant variables.

Below we analyze the behavior of the discrepancy (5) for  $\nu_1 = 1$ ,  $\mu = 2$ , and  $Y = (y^1, y^2, \dots, y^N)$  given by (12). This means that we consider the second step of the proposed approach when the variables  $x_1, x_2$  have already received the ranks 1 and 2 respectively. The analysis of other steps and possibilities can be done similarly, but it is too technical and is omitted here for brevity.

It is easy to check that for  $p = 2$ , the solutions of the system (10) can be written as follows:

$$\begin{aligned} f_1^{\lambda_1} &= \left( \frac{\lambda_1}{\lambda_2} \mathbb{I}_{K_1} + S_{N,1}^* (\lambda_2 \mathbb{I}_N + S_{N,2} S_{N,2}^*)^{-1} S_{N,1} \right)^{-1} S_{N,1}^* (\lambda_2 \mathbb{I}_N + S_{N,2} S_{N,2}^*)^{-1} Y, \\ f_2^{\lambda_2} &= \left( \frac{\lambda_2}{\lambda_1} \mathbb{I}_{K_2} + S_{N,2}^* (\lambda_1 \mathbb{I}_N + S_{N,1} S_{N,1}^*)^{-1} S_{N,2} \right)^{-1} S_{N,2}^* (\lambda_1 \mathbb{I}_N + S_{N,1} S_{N,1}^*)^{-1} Y, \end{aligned} \quad (13)$$

where  $\mathbb{I}_N$  is the identity matrix of size  $N \times N$ , and  $\mathbb{I}_{K_j}$  is the identity operator on RKHS  $\mathcal{H}_j$  generated by the kernel  $K_j(x_j, v_j)$ ,  $j = 1, 2$ .

Now, we introduce the main assumption used in our theoretical analysis. This assumption is formulated in terms of the elements of the singular-value decomposition of the sampling operators

$$S_{N,j} = \sum_{i=1}^N a_{ij} h_{ij} \langle \kappa_{ij}, \cdot \rangle_{\mathcal{H}_j}, \quad j = 1, 2, \quad (14)$$

where  $\{h_{ij}\}$ ,  $\{\kappa_{ij}\}$  are some orthonormal systems in  $\mathbb{R}^N$  and  $\mathcal{H}_j$  respectively, and  $a_{ij} \geq 0$ . Our assumption is that  $S_{N,j}$  share the same system of  $\{h_{ij}\}$ , i.e.

$$\{h_{i,1}\} = \{h_{i,2}\} = \{h_i\}. \quad (15)$$

The assumption (15) is in fact an assumption on the distribution of the sampling points  $\{x_j^i\}$ . We illustrate it in the following simple example.

**Example 1.** Let  $N = 2$ , and  $x_1^1 = x_1^2 = t$ ,  $x_2^1 = \tau_1$ ,  $x_2^2 = \tau_2$ . This means that the sampling points belong to a line parallel to the  $x_2$ -axis. If  $x_1$  is already accepted as the relevant variable, then such sampling points allow an easy test whether or not  $x_2$  should be accepted as the relevant variable. Indeed, if  $y^1 \neq y^2$ , then one really needs one more variable to explain the given data  $Y = (y^1, y^2)$ .

In the considered case, the sampling operators look as follows

$$S_{N,1}f = (f(t), f(t)), \quad S_{N,2}f = (f(\tau_1), f(\tau_2)).$$

Assume that both RKHS are generated by the same Gaussian kernel  $K(x, v) = e^{-(x-v)^2}$ . Then

$$\begin{aligned} S_{N,2}f &= (1, 1) \langle (K(\tau_1, \cdot) + K(\tau_2, \cdot))/2, f \rangle_{\mathcal{H}_2} + (1, -1) \langle (K(\tau_1, \cdot) - K(\tau_2, \cdot))/2, f \rangle_{\mathcal{H}_2}, \\ S_{N,1}f &= (1, 1) \langle K(t, \cdot), f \rangle_{\mathcal{H}_1}, \end{aligned}$$

and it is easy to check that these operators admit the decomposition (14) with

$$\begin{aligned} h_{1,1} &= h_{1,2} = (1, 1)/\sqrt{2}, \quad h_{2,1} = h_{2,2} = (1, -1)/\sqrt{2}, \\ \kappa_{1,2} &= \frac{1}{\sqrt{2}} (K(\tau_1, \cdot) + K(\tau_2, \cdot)) / (1 + e^{-(\tau_1 - \tau_2)^2})^{1/2}, \\ \kappa_{2,2} &= \frac{1}{\sqrt{2}} (K(\tau_1, \cdot) - K(\tau_2, \cdot)) / (1 - e^{-(\tau_1 - \tau_2)^2})^{1/2}, \quad \kappa_{1,1} = K(t, \cdot), \\ a_{1,2} &= (1 + e^{-(\tau_1 - \tau_2)^2})^{1/2}, \quad a_{2,2} = (1 - e^{-(\tau_1 - \tau_2)^2})^{1/2}, \quad a_{1,1} = \sqrt{2}, \quad a_{2,1} = 0. \end{aligned}$$

Thus, in the considered case the assumption (15) is satisfied.  $\square$

We would like to stress that the assumption (15) is only of the theoretical nature. At the same time, it is clear that a successful detection of relevant variables cannot be done from the data sampled at poorly distributed points  $\{x_j^i\}$ . Therefore, some restrictions on the sampling operators are unavoidable, and the condition (15) is just one of them.

**Theorem 1.** *If (12) and (15) hold true, then*

$$\|Y - S_{N,1}f_1^{\lambda_1} - S_{N,2}f_2^{\lambda_2}\|_{\mathbb{R}^N} \leq \frac{1}{2} \left( \sqrt{\lambda_1} \|f_1\|_{\mathcal{H}_1} + \sqrt{\lambda_2} \|f_2\|_{\mathcal{H}_2} \right) + \|\varepsilon\|_{\mathbb{R}^N}.$$

*Proof.* From (13)-(15), it follows that

$$S_{N,1}f_1^{\lambda_1} + S_{N,2}f_2^{\lambda_2} = \sum_{i=1}^N \frac{\lambda_1 a_{i,1}^2 + \lambda_2 a_{i,2}^2}{\lambda_1 \lambda_2 + \lambda_1 a_{i,2}^2 + \lambda_2 a_{i,1}^2} h_i \langle h_i, Y \rangle_{\mathbb{R}^N},$$

Then, in view of (12), we have

$$Y - S_{N,1}f_1^{\lambda_1} - S_{N,2}f_2^{\lambda_2} = \Sigma_1 + \Sigma_2 + \Sigma_3, \quad (16)$$

where

$$\Sigma_1 = \sum_{i=1}^N \frac{\lambda_1 \lambda_2 a_{i,1}}{\lambda_1 \lambda_2 + \lambda_1 a_{i,2}^2 + \lambda_2 a_{i,1}^2} h_i \langle \kappa_{i,1}, f_1 \rangle_{\mathcal{H}_1}, \quad (17)$$

$$\Sigma_2 = \sum_{i=1}^N \frac{\lambda_1 \lambda_2 a_{i,2}}{\lambda_1 \lambda_2 + \lambda_1 a_{i,2}^2 + \lambda_2 a_{i,1}^2} h_i \langle \kappa_{i,2}, f_2 \rangle_{\mathcal{H}_2}, \quad (18)$$

$$\Sigma_3 = \sum_{i=1}^N \frac{\lambda_1 \lambda_2}{\lambda_1 \lambda_2 + \lambda_1 a_{i,2}^2 + \lambda_2 a_{i,1}^2} h_i \langle h_i, \varepsilon \rangle_{\mathbb{R}^N}.$$

Observe now that

$$\begin{aligned} \|\Sigma_1\|_{\mathbb{R}^N} &= \left( \sum_{i=1}^N \left( \frac{\lambda_1 \lambda_2 a_{i,1}}{\lambda_1 \lambda_2 + \lambda_1 a_{i,2}^2 + \lambda_2 a_{i,1}^2} \right)^2 \langle \kappa_{i,1}, f_1 \rangle_{\mathcal{H}_1}^2 \right)^{1/2} \\ &\leq \left( \sum_{i=1}^N \left( \frac{\lambda_1 a_{i,1}}{\lambda_1 + a_{i,1}^2} \right)^2 \langle \kappa_{i,1}, f_1 \rangle_{\mathcal{H}_1}^2 \right)^{1/2} \\ &\leq \sup_t \left| \frac{\lambda_1 t}{\lambda_1 + t^2} \right| \|f_1\|_{\mathcal{H}_1} = \frac{\sqrt{\lambda_1}}{2} \|f_1\|_{\mathcal{H}_1}. \end{aligned}$$

Moreover, in the same way, we obtain that

$$\begin{aligned} \|\Sigma_2\|_{\mathbb{R}^N} &\leq \frac{\sqrt{\lambda_2}}{2} \|f_2\|_{\mathcal{H}_2}, \\ \|\Sigma_3\|_{\mathbb{R}^N} &= \left( \sum_{i=1}^N \left( \frac{\lambda_1 \lambda_2}{\lambda_1 \lambda_2 + \lambda_1 a_{i,2}^2 + \lambda_2 a_{i,1}^2} \right)^2 \langle h_i, \varepsilon \rangle_{\mathbb{R}^N}^2 \right)^{1/2} \\ &\leq \left( \sum_{i=1}^N \langle h_i, \varepsilon \rangle_{\mathbb{R}^N}^2 \right)^{1/2} = \|\varepsilon\|_{\mathbb{R}^N}. \end{aligned}$$

Combining these bounds with (16), we obtain the asserted statement.  $\square$



The above theorem allows for the conclusion that if there is a contribution to the data  $Y$  that comes from the functions of variables, say  $x_1, x_2$ , then the values of the discrepancy corresponding to the small values of the regularization parameters  $\{\lambda_1, \lambda_2\} \subset (0, 1)$  are expected to be dominated by the ones corresponding to at least one large parameter. It can be also seen from (17)-(18) that if  $a_{i,1} \neq 0$  and  $a_{i,2} \neq 0$ , then the terms in  $\Sigma_1$  and  $\Sigma_2$  are monotonically increasing functions of  $\lambda_1, \lambda_2$ .

Using the similar argument, we can extend the statement of the theorem to any number of variables, provided that corresponding sampling operators share a common singular system in  $\mathbb{R}^N$ . Then the above conclusion can also be made for more than two variables, and it is the reason behind the use of the values of the discrepancies corresponding to large and small values of the regularization parameters for detecting relevant variables as it has been described in Introduction. Thus, if the discrepancy

$$\mathcal{D} \left( f_{\nu_1}^{\lambda_{\nu_1}}, f_{\nu_2}^{\lambda_{\nu_2}}, \dots, f_{\nu_l}^{\lambda_{\nu_l}}; Z_N \right) = \left\| Y - \sum_{j=1}^l S_{N,\nu_j} f_{\nu_j}^{\lambda_{\nu_j}} \right\|_{\mathbb{R}^N} \quad (19)$$

as a function of  $(\lambda_{\nu_1}, \lambda_{\nu_2}, \dots, \lambda_{\nu_l})$  exhibits the above mentioned behavior, then the variables  $x_{\nu_1}, x_{\nu_2}, \dots, x_{\nu_l}$  are considered as the relevant ones.

Since in applications it is usually difficult to check the values of (19) for all  $\lambda_{\nu_1}, \lambda_{\nu_2}, \dots, \lambda_{\nu_l}$ , one can realize the above mentioned approach by using Monte-Carlo-type simulations. Namely, if  $x_{\nu_1}, x_{\nu_2}, \dots, x_{\nu_{l-1}}$  have been already accepted as relevant variables, then the values of (19) for the randomly chosen  $(\lambda_{\nu_1}, \lambda_{\nu_2}, \dots, \lambda_{\nu_l}) \in (0, 1)^l$  are compared to the ones for the randomly chosen  $(\lambda_{\nu_1}, \lambda_{\nu_2}, \dots, \lambda_{\nu_l}) \in (0, 1)^{l-1} \times [1, B]$ ,  $B > 1$ , and  $x_{\nu_l}$  is accepted as the relevant variable if in the above simulations the values of (19) for  $(\lambda_{\nu_1}, \lambda_{\nu_2}, \dots, \lambda_{\nu_l}) \in (0, 1)^l$  are dominated by the ones for  $(\lambda_{\nu_1}, \lambda_{\nu_2}, \dots, \lambda_{\nu_l}) \in (0, 1)^{l-1} \times [1, B]$ .

*Remark 1.* Note that the conclusion about the ordered behavior of the discrepancy made on the basis of Theorem 1 can be seen as an extension of the following interpretation of the values of discrepancies  $\|S_{N,1} f_j^{\lambda_j} - Y\|_{\mathbb{R}^N}$  for the single penalty regularization. From [26, Lemma 3.1], it follows that

$$\begin{aligned} \lim_{\lambda_j \rightarrow 0} \left\| S_{N,j} f_j^{\lambda_j} - Y \right\|_{\mathbb{R}^N} &= \inf_{f \in \mathcal{H}_j} \|S_{N,j} f - Y\|_{\mathbb{R}^N}, \\ \lim_{\lambda_j \rightarrow \infty} \left\| S_{N,j} f_j^{\lambda_j} - Y \right\|_{\mathbb{R}^N} &= \|Y\|_{\mathbb{R}^N}. \end{aligned}$$

Then it is clear that if  $\mathcal{H}_j$  is dense in the corresponding space of continuous functions, and

$$Y = S_{N,j} f_j + \varepsilon, \quad \|\varepsilon\|_{\mathbb{R}^N} < \|Y\|_{\mathbb{R}^N},$$

then for small  $\lambda_j$  and large  $\bar{\lambda}_j$ , one can expect

$$\left\| S_{N,j} f_j^{\lambda_j} - Y \right\|_{\mathbb{R}^N} < \left\| S_{N,j} f_j^{\bar{\lambda}_j} - Y \right\|_{\mathbb{R}^N}.$$

On the other hand, if  $Y \in (\text{Range}(S_{N,j}))^\perp$  such that there is no contribution to  $Y$  allowing a representation in terms of the values of  $f_j \in \mathcal{H}_j$  at the points  $\{x_j^i\}_{i=1}^N$ , then the discrepancies  $\left\| S_{N,j} f_j^{\lambda_j} - Y \right\|_{\mathbb{R}^N}$  do not behave in the ordered way.

Of course, in the case of the single variable and penalty, no additional assumptions like, for example, (15) are needed to justify the ordered behavior of the discrepancy  $\left\| S_{N,j} f_j^{\lambda_j} - Y \right\|_{\mathbb{R}^N}$  for  $Y = S_{N,j} f_j + \varepsilon$ .  $\square$

At the end of this theoretical section, we illustrate the above approach on the example from [16], where for  $p = 40$  and  $N = 100$ , the data set  $Z_N = \{ (x_1^i, x_2^i, \dots, x_p^i; y^i) \}_{i=1}^N$  is simulated in such a way that the values  $x_j^i$  are sampled uniformly at random from the interval  $[-2, 2]$ , and

$$y^i = \sum_{j=1}^4 (x_j^i)^2 + \varepsilon^i, \quad (20)$$

where  $\varepsilon^i$  are zero-mean Gaussian random variables with variances chosen so that the signal-to-noise ratio is 15 : 1.

The input (20) means that in this example the target function (1) depends on the first 4 variables. Recall that in our approach we, at first, need to rank the variables  $x_1, x_2, \dots, x_{40}$  according to the values of the discrepancies  $\mathcal{D}(f_j^{\lambda_j}(x_j); Z_N)$ ,  $j = 1, 2, \dots, 40$ , where  $f_j^{\lambda_j}$  is the minimizer of the Tikhonov functional

$$T_\lambda(f; Z_N) = \frac{1}{N} \sum_{i=1}^N (y^i - f(x_j^i))^2 + \lambda \|f\|_{\mathcal{H}}^2. \quad (21)$$

In our experiments, we choose in (21)  $\lambda = \lambda_j = \lambda^{(k_j)}$  from the set

$$\Lambda_{50} = \{ \lambda = \lambda^{(k)} = 10^{-4} \cdot (1.3)^k, k = 1, 2, \dots, 50 \}$$

according to the quasi-optimality criterion (see, e.g. [25, 4, 12]). Moreover, in (21) the space  $\mathcal{H}$  is chosen to be the RKHS generated by the polynomial kernel of degree 2. This choice is made according to [16], where the same kernel has been used in the approach (3) for dealing with the data (20).

For the considered simulation of the data (20) the sequence of the variables ordered according to their ranks looks as follows:

$$x_2, x_4, x_3, x_1, x_{33}, x_6, \dots, x_{18}. \quad (22)$$

Then as it is described above, the next step consists in testing whether the values of the discrepancy

$$\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}; Z_N) = \left( \frac{1}{N} \sum_{i=1}^N (y^i - f_2^{\lambda_2}(x_2^i) - f_4^{\lambda_4}(x_4^i))^2 \right)^{1/2}$$

corresponding to the small values  $\lambda_2, \lambda_4$  are dominated by the ones corresponding to the small  $\lambda_2$  and the large  $\lambda_4$ . Here and below we use the convention that in the notation  $\mathcal{D}(f_{\mu_1}^{\lambda_{\mu_1}}, f_{\mu_2}^{\lambda_{\mu_2}}, \dots, f_{\mu_l}^{\lambda_{\mu_l}}; Z_N)$ , the symbols  $f_{\mu_j}^{\lambda_{\mu_j}}$  mean the minimizers of the functional

$$T_\lambda^2(f_{\mu_1}, f_{\mu_2}, \dots, f_{\mu_l}; Z_N) = \frac{1}{N} \sum_{i=1}^N \left( y^i - \sum_{j=1}^l f_{\mu_j}(x_{\mu_j}^i) \right)^2 + \sum_{j=1}^l \lambda_{\mu_j} \|f_{\mu_j}\|_{\mathcal{H}_{\mu_j}}^2.$$

In our experiments the small values of the regularization parameters are randomly chosen within the set

$$\Lambda_{50}^{\text{small}} = \{ \lambda = \lambda^{(k)} = 10^{-4} \cdot (1.3)^k, k = 1, 2, \dots, 15 \},$$

while the large values are selected at random from

$$\Lambda_{50}^{\text{large}} = \{ \lambda = \lambda^{(k)} = 10^{-4} \cdot (1.3)^k, k = 40, 41, \dots, 50 \}.$$

Moreover, in all experiments the random choice of the regularization parameters from  $\Lambda_{50}^{\text{small}}$  and  $\Lambda_{50}^{\text{large}}$  is performed 15 times.

For the considered simulations of the data (20) and randomly chosen  $\lambda_2, \lambda_4$ , the values of the discrepancy  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}; Z_N)$  are displayed in Figure 1. In this Figure, as well as in all figures below, the horizontal axis represents the run number of the simulation, while the vertical axis represents the observed value of the discrepancy. Furthermore, the values corresponding to the simulations with the parameters from  $\Lambda_{50}^{\text{small}}$  and  $\Lambda_{50}^{\text{large}}$  are connected by red solid lines and blue dashed lines respectively.

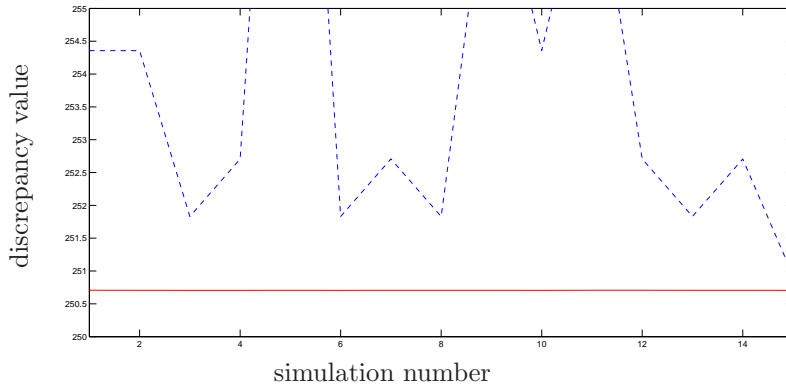


Figure 1: The experiment with the data (20). The behavior of the discrepancy  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}; Z_N)$  for  $\{\lambda_2, \lambda_4\} \subset \Lambda_{50}^{\text{small}}$  (red solid line), and  $\lambda_2 \in \Lambda_{50}^{\text{small}}, \lambda_4 \in \Lambda_{50}^{\text{large}}$  (blue dashed line).

Note that in Figure 1 and in some other figures below, the curves displaying the values of the discrepancy for the regularization parameters from  $\Lambda_{50}^{\text{small}}$  look like straight lines. In view of Theorem 1, the fluctuations in the values of the discrepancy corresponding to the small values of the regularization parameters are indeed small. They are not so much visible because of the vertical axis scaling used in the figures.

According to our approach, the behavior of the discrepancy displayed in Figure 1 means that the corresponding variables  $x_2, x_4$  have to be accepted as the relevant ones. Then taking into account the ranking (22), we need to check the behavior of the discrepancy  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}, f_3^{\lambda_3}; Z_N)$  for  $\{\lambda_2, \lambda_4, \lambda_3\} \subset \Lambda_{50}^{\text{small}}$ , and  $\{\lambda_2, \lambda_4\} \subset \Lambda_{50}^{\text{small}}, \lambda_3 \in \Lambda_{50}^{\text{large}}$ . This behavior is displayed in Figure 2 and it allows the acceptance of  $x_3$  as the next relevant variable.

In view of Figure 3 displaying the behavior of the discrepancy

$$\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}, f_3^{\lambda_3}, f_1^{\lambda_1}; Z_N),$$

the same conclusion can be made regarding the variable  $x_1$ .

At the same time, further testing along the ranking list (22) shows that the discrepancies  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}, f_3^{\lambda_3}, f_1^{\lambda_1}, f_j^{\lambda_j}; Z_N)$  with  $j = 33, 6, \dots, 18$  do not exhibit the ordered behavior for  $\lambda_j \in \Lambda_{50}^{\text{small}}$  and  $\lambda_j \in \Lambda_{50}^{\text{large}}$ . Typical examples are displayed in

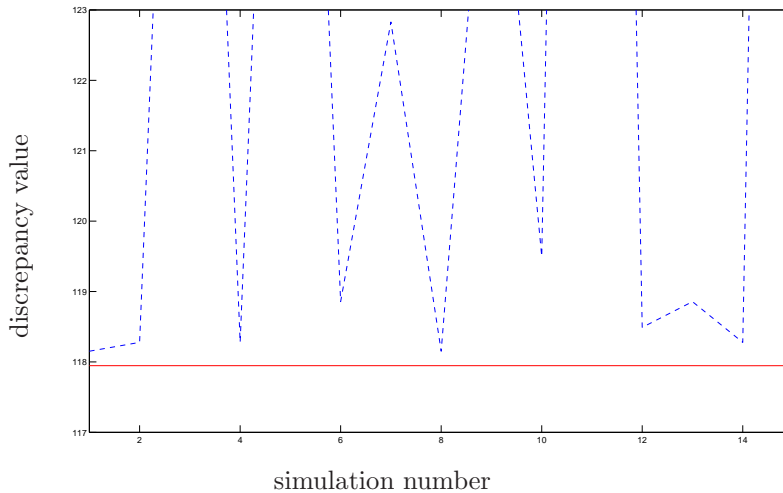


Figure 2: The experiment with the data (20). The behavior of the discrepancy  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}, f_3^{\lambda_3}; Z_N)$  for  $\{\lambda_2, \lambda_4, \lambda_3\} \subset \Lambda_{50}^{\text{small}}$  (red solid line), and  $\{\lambda_2, \lambda_4\} \subset \Lambda_{50}^{\text{small}}$ ,  $\lambda_3 \in \Lambda_{50}^{\text{large}}$  (blue dashed line).

Figure 4 and Figure 5. Therefore, our approach does not allow the acceptance of  $x_{33}, x_6, \dots, x_{18}$  as the relevant variables.

Thus, for the considered simulation of the data (20) all relevant variables are correctly detected by the proposed approach.

### 3 Application to the reconstruction of a causality network

In this section we discuss the application of our approach based on multi-penalty regularization to the inverse problem of detecting causal relationships between genes from the time series of their expression levels.

Viewing each gene in a genome as a distinct variable, say  $u_\nu$ , associated to the rate of gene expression, the value  $u_\nu^t = u_\nu(t)$  of this variable at time moment  $t$  can be influenced by the values  $u_j^\tau = u_j(\tau)$ ,  $j = 1, \dots, p$ , at the time moments preceding  $t$ , i.e.,  $\tau < t$ . This influence is realized through the regulatory proteins produced by genes. Moreover, gene expression levels  $u_j^\tau$  are often interpreted and measured in terms of levels or amounts of such proteins. Therefore, time series gene expression data can be used for detecting causal relationships between genes and constructing gene regulatory networks allowing better insights into the underlying cellular mechanisms.

A gene regulatory network or, more generally, a causality network is a directed graph with nodes that are variables  $u_\nu$ ,  $\nu = 1, 2, \dots, p$ , and directed edges representing causal relations between variables. We write  $u_\nu \leftarrow u_j$  if the variable  $u_j$  has the causal influence on the variable  $u_\nu$ . An example of such a network is presented in Figure 6. This network contains genes that are active in the human cancer cell line HeLa [27]. It was derived from the biological experiments in [14], and then, it was used for testing several algorithms devoted to the causality detection [21, 15, 23, 19]. Using the same data as in the above papers, we discuss an applicability of our ap-

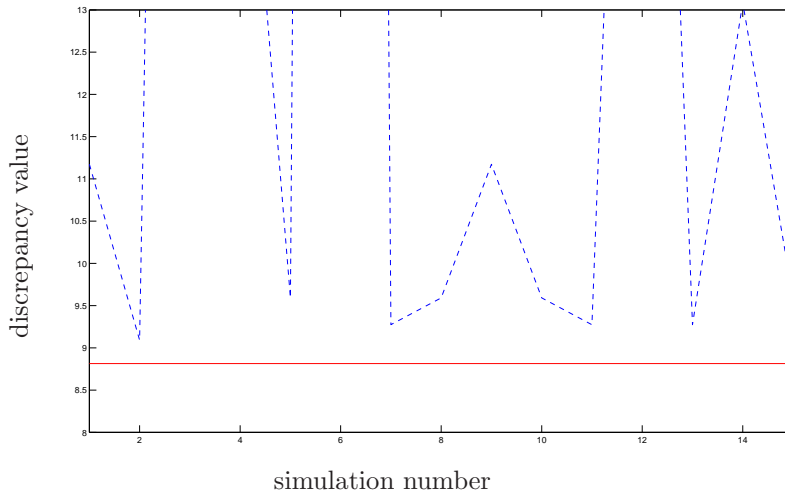


Figure 3: The experiment with the data (20). The behavior of the discrepancy  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}, f_3^{\lambda_3}, f_1^{\lambda_1}; Z_N)$  for  $\{\lambda_2, \lambda_4, \lambda_3, \lambda_1\} \subset \Lambda_{50}^{\text{small}}$  (red solid line), and  $\{\lambda_2, \lambda_4, \lambda_3\} \subset \Lambda_{50}^{\text{small}}, \lambda_1 \in \Lambda_{50}^{\text{large}}$  (blue dashed line).

proach in reconstructing the causalities within this network.

A causality network can be characterized by the so-called adjacency matrix  $A = \{A_{\nu,j}\}_{\nu,j=1}^p$  with the following elements  $A_{\nu,j} = 1$  if  $u_\nu \leftarrow u_j$ , otherwise,  $A_{\nu,j} = 0$ . In Figure 7 we present the adjacency matrix  $A = A^{\text{true}}$  corresponding to the causality network displayed in Figure 6. Adjacency matrices allow a convenient comparison of different reconstruction methods of causality networks.

Note that causality networks arise in various scientific contexts. A detailed overview of the approaches for measuring a causal influence can be found in [10], where it is mentioned that the introduction of the concept of causality into the analysis of data observed in time series is due to Clive W. J. Granger [8], who was awarded the Nobel Prize in Economic Sciences in 2003.

The concept of causality in the Granger approach is based on the assumption that (i) the cause should precede its effect, and (ii) the cause contains an information about the effect that is in no other variable. A consequence of these assumptions is that the causal variable  $u_j$  can help to forecast the effect variable  $u_\nu$ . In this restricted sense of causality, referred to as Granger causality, the variable  $u_j$  is said to cause another variable  $u_\nu$  if future values  $u_\nu^t$ ,  $t = L + 1, L + 2, \dots, T$ , of  $u_\nu$  can be better predicted using the past values  $u_j^\tau$ ,  $u_\nu^\tau$ ,  $\tau = t - 1, t - 2, \dots, t - L$ , of  $u_j$  and  $u_\nu$  rather than using only the past values of  $u_\nu$ . Here  $L$  is the maximum lag allowed in the past observations, and we assume that the available time series data are  $\{u_j^t\}_{t=1}^T$ ,  $\{u_\nu^t\}_{t=1}^T$ .

The notion of Granger causality was originally defined for a pair of time series and was based on linear regression models. If we are interested in cases in which  $p$  time series variables are presented and we wish to determine a causal relationships between them, then we naturally turn to the Graphical Granger modeling based on the linear multivariate regression of the form

$$u_\nu^t \approx \sum_{j=1}^p \sum_{l=1}^L \beta_j^l u_j^{t-l}, \quad t = L + 1, L + 2, \dots, T. \quad (23)$$

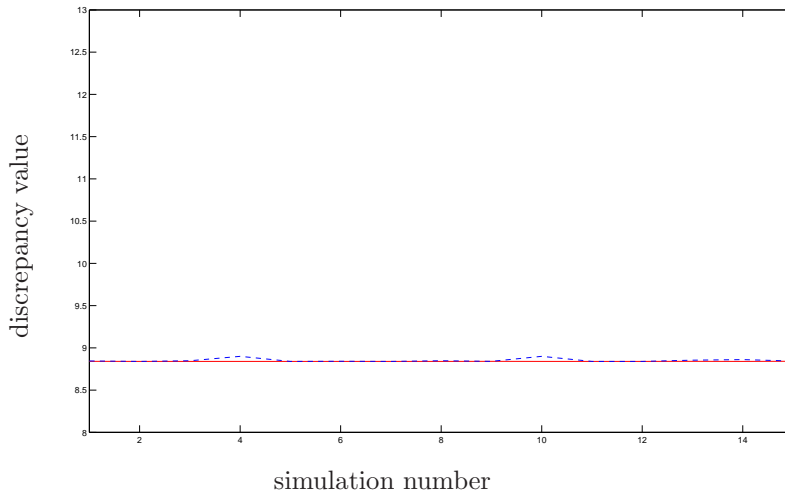


Figure 4: The experiment with the data (20). The behavior of the discrepancy  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}, f_3^{\lambda_3}, f_1^{\lambda_1}, f_{33}^{\lambda_{33}}; Z_N)$  for  $\{\lambda_2, \lambda_4, \lambda_3, \lambda_1, \lambda_{33}\} \subset \Lambda_{50}^{\text{small}}$  (red solid line), and  $\{\lambda_2, \lambda_4, \lambda_3, \lambda_1\} \subset \Lambda_{50}^{\text{small}}, \lambda_{33} \in \Lambda_{50}^{\text{large}}$  (blue dashed line).

Then,  $u_j$  is said to be Granger-causal for  $u_\nu$  if the corresponding coefficients  $\beta_j^l$ ,  $l = 1, 2, \dots, L$ , are in some sense significant. Thus, we are interested in selecting the most important coefficients. For this purpose, a particular relevant class of methodologies is those that combine regression with variable selection, such as the Lasso [24, 31], which minimizes the squared discrepancy plus a penalty on the sum, or the weighted sum of the absolute values of the regression coefficients  $\beta_j^l$ .

Lasso-type estimates have been used for discovering graphical Granger causality by a number of researchers, including [1, 23, 19]. Note that in regularization theory Lasso is known as the  $l_1$ -Tikhonov regularization. It has been extensively studied in the framework of the reconstruction of the sparse structure of an unknown signal. It should be also mentioned that the sparsity enforcing regularization techniques, such as Lasso, are viewed now as a methodology for the quantitative inverse problems in systems biology [6].

At the same time, as it is mentioned in [15], the Lasso estimate of the graphical Granger causality may result in a model (23) in which the large (significant) coefficients  $\beta_j^l$  appear in many sums  $\sum_{l=1}^L \beta_j^l u_j^{t-l}$ . Such a model is hard to interpret, because of natural groupings existing between time series variables  $\{u_j^{t-l}\}_{l=1}^L$ ,  $j = 1, 2, \dots, p$ . We mean that the time series variables  $\{u_j^{t-l}\}_{l=1}^L$  with the same index, say  $j = j_1$ , should be either selected or eliminated as a whole. The group Lasso procedure [29, 30] was invented to address this issue and it was used in [15] in order to obtain the corresponding Granger graphical model of gene regulatory networks. According to this model, a gene  $u_{j_1}$  causes a gene  $u_\nu$  if in (23) the coefficients  $\beta_{j_1}^l$ ,  $l = 1, 2, \dots, L$ , are significant components of the vector  $\beta = (\beta_j^l)$  solving the minimization problem

$$\sum_{t=L+1}^T \left( u_\nu^t - \sum_{j=1}^p \sum_{l=1}^L \beta_j^l u_j^{t-l} \right)^2 + \lambda \sum_{j=1}^p \left( \sum_{l=1}^L (\beta_j^l)^2 \right)^{1/2} \rightarrow \min_{\beta}. \quad (24)$$

Here note that similarly to (3) by using the square root  $(\cdot)^{1/2}$  in the penalty term,

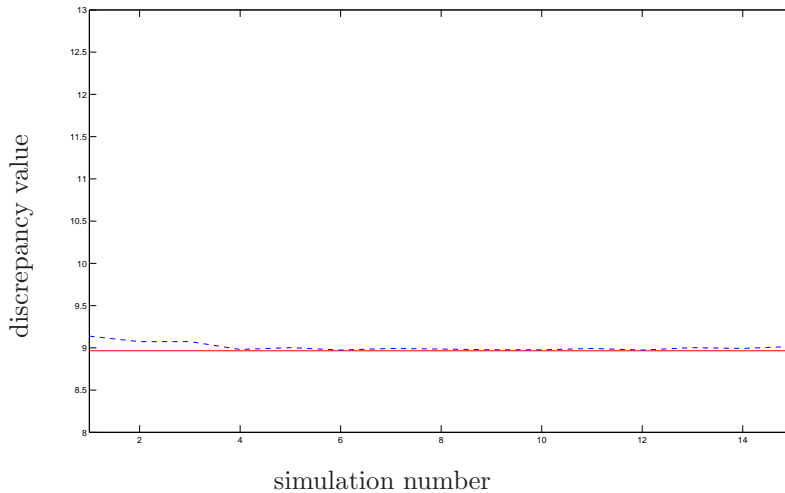


Figure 5: The experiment with the data (20). The behavior of the discrepancy  $\mathcal{D}(f_2^{\lambda_2}, f_4^{\lambda_4}, f_3^{\lambda_3}, f_1^{\lambda_1}, f_6^{\lambda_6}; Z_N)$  for  $\{\lambda_2, \lambda_4, \lambda_3, \lambda_1, \lambda_6\} \subset \Lambda_{50}^{\text{small}}$  (red solid line), and  $\{\lambda_2, \lambda_4, \lambda_3, \lambda_1\} \subset \Lambda_{50}^{\text{small}}, \lambda_6 \in \Lambda_{50}^{\text{large}}$  (blue dashed line).

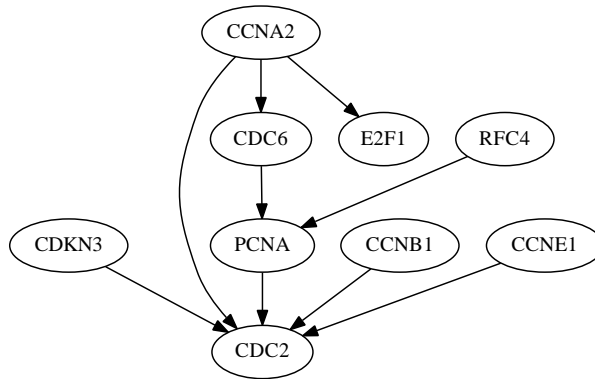


Figure 6: Causality network of the human cancer cell line HeLa from the BioGRID database ([www.thebiogrid.org](http://www.thebiogrid.org)).

one encourages the coefficients associated with each particular gene to be similar in amplitude, as contrary to using the  $l_1$ -norm, for example. The opposite side of this is that the procedures of minimizing (24) are nonlinear and require the solution of  $O(pL)$  equations on each iteration step. This can be computationally expensive for large number  $p$  of genes.

On the other hand, the above mentioned natural groupings between the values  $u_j^t$  of variables  $u_j$  can be introduced already in the multivariate regression by considering instead of (23) the following form

$$u_\nu^t \approx \sum_{j=1}^p f_j \left( \sum_{l=1}^L \beta_j^l u_j^{t-l} \right), \quad t = L + 1, L + 2, \dots, T, \quad (25)$$

where  $f_j$  are univariate functions in some Reproducing Kernel Hilbert Spaces  $\mathcal{H}_j$ . Note that (25) can be seen as a particular form of structural equation models discussed in [18]. Then a conclusion that the gene  $u_{j_1}$  causes the gene  $u_\nu$  can be

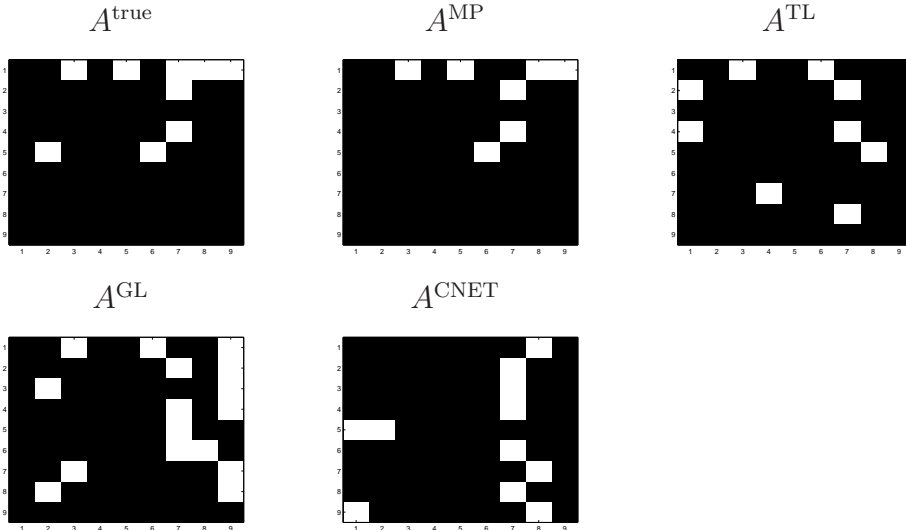


Figure 7: The adjacency matrix  $A^{\text{true}}$  for the causality network in Figure 6 and its various estimations. The white squares correspond to  $A_{i,j} = 1$ ; the black squares — to the zero-elements. The genes are numbered in the following order: CDC2, CDC6, CDKN3, E2F1, PCNA, RFC4, CCNA2, CCNB1, CCNE1.

drawn by determining that the variable  $x_{j_1}$  is a relevant variable of a function of the form (1) whose values at the points

$$x_j^i = \sum_{l=1}^L \beta_j^l u_j^{L+i-l}, \quad i = 1, 2, \dots, T - L, \quad j = 1, 2, \dots, p, \quad (26)$$

are equal to

$$y^i = u_\nu^{L+i}, \quad i = 1, 2, \dots, T - L. \quad (27)$$

Of course, the latter conclusion can be drawn only when in (26) some values of the coefficients  $\beta_j^l$  have been already set. For example, these regression coefficients can be precomputed in (23) by some inexpensive algorithm such as the ordinary or regularized least squares (OLS or RLS). Note that such a precomputation step is also required in Adaptive Lasso [31] that has been discussed in the context of the regulatory networks discovery in [15], and where an auxiliary vector estimator of the coefficients in (23) is usually obtained by OLS or Ridge Regression.

Another possibility of determining the coefficients  $\beta_j^l$  in (26) is to use the output vector of any of the graphical Granger models based on (23) such as [15, 23]. In this case, the discussed approach provides an opportunity of additional evaluation of these models in the sense that causal relationships detected by them and confirmed in the discussed approach can be considered as more certain.

After specifying the coefficients  $\beta_j^l$  in (26), the values (26), (27) can form the data set  $Z_N = \{(x_1^i, x_2^i, \dots, x_p^i; y^i)\}_{i=1}^N$ ,  $N = T - L$ . Then, the detection of the relevant variables from the data  $Z_N$  follows the approach described in Section 1 and analyzed in Section 2. The only adjustment is that in view of the idea of Granger causality (comparison of the accuracy of regressing for  $u_\nu$  in terms of its own past values with that of regressing in terms of the values  $u_\nu$  and the values of a possible cause), we start the ranking list of variables with the variable  $x_\nu$  when looking for the genes causing the gene  $u_\nu$ .



	$P$	$R$	$F_1$
$A^{\text{MP}}$	1	0.78	0.88
$A^{\text{CNET}}$	0.36	0.44	0.4
$A^{\text{GL}}$	0.24	0.44	0.3
$A^{\text{TL}}$	0.3	0.33	0.32

Table 1: The values of the performance measures for the adjacency matrices in Figure 7.

Below we present the results of the application of the proposed approach to the data of the gene expressions for the network of genes displayed in Figure 6. These data is taken as in [21, 15, 23]. In (9), (10), (25) all univariate functions  $f_j$  are assumed to be in the same RKHS generated by the Gaussian kernel  $K(x, v) = e^{-(x-v)^2}$ . Moreover, the standard RLS-algorithm has been used for precomputing the regression coefficients in (25), (26). The regularization parameter in RLS has been chosen according to the quasi-optimality criterion. As in [21, 15, 23] the gene expressions  $\{u_j^t\}$  are observed for  $t = 1, 2, \dots, 47$ , and, as it is suggested in [19], the maximum lag was chosen as  $L = 4$ . Then, we follow the same steps as in the illustrating example in Section 2. In particular, we use the same sets  $\Lambda_{50}^{\text{small}}, \Lambda_{50}^{\text{large}}$ .

The application of the proposed approach to the above mentioned data results in the adjacency matrix  $A^{\text{MP}}$  displayed in Figure 7.

As it has been already mentioned, the data corresponding to the causality network in Figure 6 was used for testing several methods devoted to the regulatory networks modeling. First, it was used in [21], where the authors developed a search-based algorithm, called CNET, and applied it to this set of data. Then, the same set of nine genes was also analyzed in [15] by means of group Lasso (GL) algorithm based on the minimization of the functionals of the form (24). In [23] the authors pointed out some limitations of GL-algorithm and proposed to overcome them by means of the so-called truncating Lasso (TL) penalty algorithm. Figure 7 presents the adjacency matrices  $A^{\text{CNET}}, A^{\text{GL}}, A^{\text{TL}}$  of the estimated causality network with the genes from Figure 6 obtained respectively by the algorithms from [21, 15, 23].

As in [23] to assess the performance of the discussed algorithms, we use three well-known performance measures: precision ( $P$ ), recall ( $R$ ), and their harmonic mean ( $F_1$ ) (see, e.g., [28]). Table 1 contains the values of these measures for the adjacency matrices given by the discussed methods and displayed in Figure 7. This table shows that the best performance is achieved by our approach.

To illustrate the steps of our approach in reconstructing the network from Figure 6, we present Figures 8-10 displaying the behavior of the discrepancies, which in the present context play the role of the indicators for the causal relationships. These figures are related to CDC2 gene numbered as  $x_1$ . We take this gene as an example because its causing genes are poorly detected by the CNET, GL, and TL algorithms.

Using the data for this gene and transforming them into (26),(27) with  $\nu = 1$ , we receive the following sequence of the variables ordered according to their ranks

$$x_1, x_3, x_7, x_5, x_4, x_9, x_6, x_2, x_8.$$

Figures 8-10 display the behavior of the discrepancies

$$\begin{aligned} & \mathcal{D}(f_1^{\lambda_1}, f_3^{\lambda_3}, f_5^{\lambda_5}; Z_N), \\ & \mathcal{D}(f_1^{\lambda_1}, f_3^{\lambda_3}, f_5^{\lambda_5}, f_9^{\lambda_9}, f_6^{\lambda_6}; Z_N), \\ & \mathcal{D}(f_1^{\lambda_1}, f_3^{\lambda_3}, f_5^{\lambda_5}, f_9^{\lambda_9}, f_8^{\lambda_8}; Z_N) \end{aligned}$$

considered respectively at 3th, 6th and 8th steps of our approach. The reason to present these steps as examples is explained below.

The behavior of the discrepancy displayed in Figure 8 indicates that according to our approach, the variable  $x_5$ , which corresponds to PCNA gene, should be considered as the cause for CDC2. From Figure 6 one can see that this causal relationship is true, but it has not been detected by any other considered algorithms.

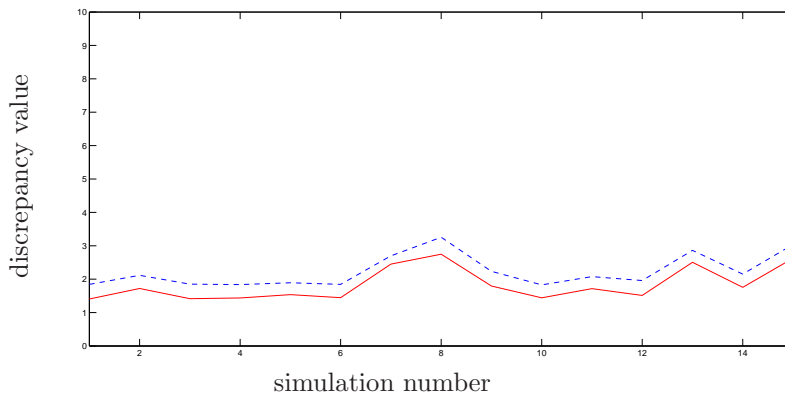


Figure 8: The experiment with the gene expressions data. The behavior of the discrepancy  $\mathcal{D}(f_1^{\lambda_1}, f_3^{\lambda_3}, f_5^{\lambda_5}; Z_N)$  for small values of the regularization parameters (red solid line) and for small  $\lambda_1, \lambda_3$  and large  $\lambda_5$  (blue dashed line).

According to our approach, the interpretation of the erratic behavior of the discrepancies in Figure 9 is that  $x_6$  is not the relevant variable, and therefore, the corresponding gene RFC4 does not cause CDC2. This conclusion is also in agreement with Figure 6. At the same time, the relationship  $\text{RFC4} \rightarrow \text{CDC2}$  is wrongly detected by both Lasso-based algorithms GL and TL.

The situation in Figure 10 is opposite. According to our approach, the behavior displayed in this Figure means that  $x_8$  is the relevant variable and, thus,  $\text{CCNB1} \rightarrow \text{CDC2}$ . This relationship is true, but it was not detected by the Lasso-based algorithms.

Therefore, in our opinion, Table 1 and Figures 8-10 can be seen as an evidence of the reliability of the proposed approach in the application to the real data.

## 4 Conclusion

We have proposed a new method for detecting the relevant variables. The method is based on the inspection of the behavior of discrepancies of multi-penalty regularization with a component-wise penalization for small and large values of the

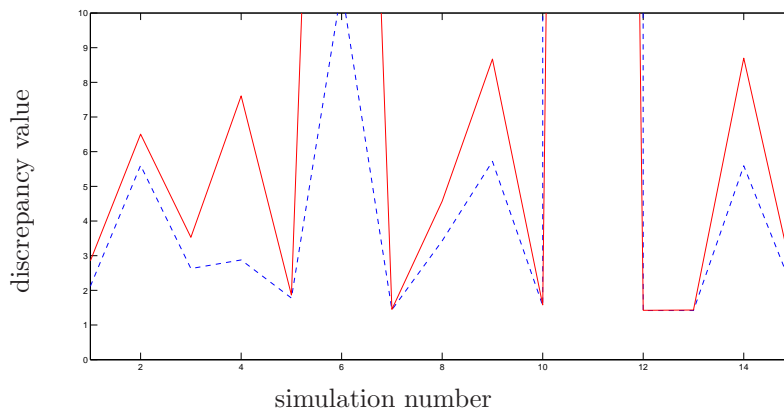


Figure 9: The experiment with the gene expressions data. The behavior of the discrepancy  $\mathcal{D}(f_1^{\lambda_1}, f_3^{\lambda_3}, f_5^{\lambda_5}, f_9^{\lambda_9}, f_6^{\lambda_6}; Z_N)$  for small values of the regularization parameters (red solid line) and for small  $\lambda_1, \lambda_3, \lambda_5, \lambda_9$  and large  $\lambda_6$  (blue dashed line).

regularization parameters. An ordered behavior suggests the acceptance of the hypothesis that the corresponding variable is the relevant one, while an erratic behavior of discrepancies is the signal for the rejection of the hypothesis.

We provided justification of the proposed method under the condition that the corresponding sampling operators share a common singular system in  $\mathbb{R}^n$ . Then we demonstrated the method in the application to the inverse problem of the reconstruction of the gene regulatory networks.

A promising performance of the method in the mentioned application calls for its further investigation. In particular, it is interesting to study the conditions on the sampling points/operators guaranteeing or preventing the detection of the relevant variables. It is also interesting to study the application of the proposed approach to the detection of the cause-effect relationships in various scientific contexts. As it was mentioned, the approach can be realized on the top of different techniques for discovering Granger causality. So, the coupling of the known techniques with the presented approach is a further interesting point for detailed investigations.

## 5 Acknowledgements

The major part of this work has been prepared, when the second author was staying at RICAM as a PostDoc. She gratefully acknowledges the partial support by the Austrian Fonds zur Förderung der Wissenschaftlichen Forschung (FWF), grant P25424 “Data-driven and problem-oriented choice of the regularization space.”

## References

- [1] A. Arnold, Y. Liu, and N. Abe. Temporal causal modeling with graphical Granger methods. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 66–75, USA, 2007. ACM New York.

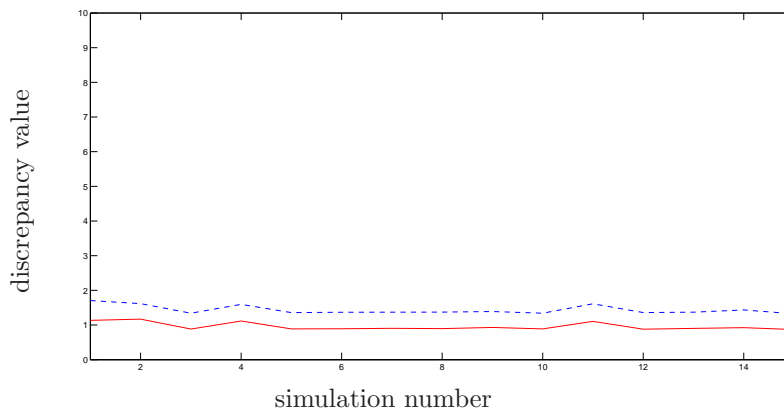


Figure 10: The experiment with the gene expressions data. The behavior of the discrepancy  $\mathcal{D}(f_1^{\lambda_1}, f_3^{\lambda_3}, f_5^{\lambda_5}, f_9^{\lambda_9}, f_8^{\lambda_8}; Z_N)$  for small values of the regularization parameters (red solid line) and for small  $\lambda_1, \lambda_3, \lambda_5, \lambda_9$  and large  $\lambda_8$  (blue dashed line).

- [2] F. Bach. Exploring large feature spaces with hierarchical multiple kernel learning. In D. Koller and et al, editors, *Advances in Neural Information Processing Systems 21*, pages 105–112. 2009.
- [3] F. Bach, G. R. G. Lanckriet, and M. I. Jordan. Multiple kernel learning, conic duality, and the SMO algorithm. In *Proceedings of the Twenty-first International Conference on Machine Learning*, 2004.
- [4] F. Bauer and M. Reiß. Regularization independent of the noise level: an analysis of quasi-optimality. *Inverse Probl.*, 24(5):16 p., 2008.
- [5] I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Commun. Pure Appl. Math.*, 57(11):1413–1457, 2004.
- [6] H. W. Engl, C. Flamm, J. Lu, P. Kügler, S. Müller, and P. Schuster. Inverse problems in systems biology. *Inverse Probl.*, 25(12):123014, 2009.
- [7] S. Foucart and H. Rauhut. *A Mathematical Introduction to Compressive Sensing*. Springer, New York, 2013.
- [8] C. Granger. Investigating causal relations by econometric models and crossspectral methods. *Econometrica*, 37:424–438, 1969.
- [9] T. Hastie and R. Tibshirani. *Generalized additive models*. Chapman and Hall, 1990.
- [10] K. Hlavackova-Schindler, M. Palus, M. Vejmelka, and J. Bhattacharya. Causality detection based on information-theoretic approaches in time series analysis. *Physics Reports*, 441:1–46, 2007.
- [11] G. S. Kimeldorf and G. Wahba. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2):495–502, 1970.

- [12] S. Kindermann and A. Neubauer. On the convergence of the quasioptimality criterion for (iterated) Tikhonov regularization. *Inverse Probl. Imaging*, 2(2):291–299, 2008.
- [13] V. Koltchinskii and M. Yuan. Sparsity in multiple kernel learning. *Ann. Stat.*, 38(6):3660–3695, 2010.
- [14] X. Li and et al. Discovery of time-delayed gene regulatory networks based on temporal gene expression profiling. *BMC Bioinformatics*, 7(26), 2006.
- [15] A. C. Lozano, N. Abe, Y. Liu, and S. Rosset. Grouped graphical Granger modeling for gene expression regulatory networks discovery. *Bioinformatics*, 25:110–118, 2009.
- [16] S. Mosci, L. Rosasco, M. Santoro, A. Verri, and S. Villa. Nonparametric sparsity and regularization. Technical Report 41, MIT, CSAIL, Cambridge, USA, 2011.
- [17] V. Naumova and S. Pereverzyev. Multi-penalty regularization with a component-wise penalization. *Inverse Probl.*, 29(7):15, 2013.
- [18] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge, 2000.
- [19] S. Pereverzyev Jr and K. Hlaváčková-Schindler. Graphical Lasso Granger method with 2-levels-thresholding for recovering causality networks. Technical report, University of Innsbruck, Department of Mathematics, Applied Mathematics Group, 2013.
- [20] P. Ravikumar, J. Lafferty, H. Liu, and L. Wasserman. Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(5):1009–1030, 2009.
- [21] F. Sambo, B. D. Camillo, and G. Toffolo. CNET: an algorithm for reverse engineering of causal gene networks. In *NETTAB2008*, Varenna, Italy, 2008.
- [22] B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In *Computational learning theory*, pages 416–426. Springer: Lecture Notes in Computer Science 2111, 2001.
- [23] A. Shojaie and G. Michailidis. Discovering graphical Granger causality using the truncating lasso penalty. *Bioinformatics*, 26:i517–i523, 2010.
- [24] R. Tibshirani. Regression shrinkage and selection via the Lasso. *J. R. Stat. Soc. B*, 58:267–288, 1996.
- [25] A. N. Tikhonov and V. B. Glasko. Use of the regularization method in non-linear problems. *USSR Comp. Math. Math. Phys.*, 5:93–107, 1965.
- [26] G. M. Vainikko and A. Y. Veretennikov. *Iteration Procedures in Ill-Posed Problems*. Moscow: Nauka, 1986. In Russian.
- [27] M. L. Whitfield and et al. Identification of genes periodically expressed in the human cell cycle and their expression in tumors. *Mol. Biol. Cell*, 13:1977–2000, 2002.

- [28] Wikipedia. Precision and recall — Wikipedia, The Free Encyclopedia, 2014. [Online; accessed 3-January-2014].
- [29] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *J. R. Stat. Soc. B*, 68:49–67, 2006.
- [30] P. Zhao, G. Rocha, and B. Yu. The composite absolute penalties family for grouped and hierarchical variable selection. *Ann. Stat.*, 37(6A):3468–3497, 2009.
- [31] H. Zou. The adaptive Lasso and its oracle properties. *J. Am. Stat. Ass.*, 101(476):1418–1429, 2006.