

# A PREDICTIVE BASED REGRESSION ALGORITHM FOR GENE NETWORK SELECTION

STÉPHANE GUERRIER<sup>1\*</sup>, NABIL MILI<sup>2\*</sup>, ROBERTO MOLINARI<sup>2</sup>,  
SAMUEL ORSO<sup>2</sup>, MARCO AVELLA-MEDINA<sup>2</sup> & YANYUAN MA<sup>3</sup>

<sup>1</sup>DEPARTMENT OF STATISTICS  
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN, USA  
EMAIL: [stephane@illinois.edu](mailto:stephane@illinois.edu)

<sup>2</sup>RESEARCH CENTER FOR STATISTICS  
GENEVA SCHOOL OF ECONOMICS AND MANAGEMENT  
UNIVERSITY OF GENEVA, SWITZERLAND

<sup>3</sup>DEPARTMENT OF STATISTICS  
UNIVERSITY OF SOUTH CAROLINA, USA

## Abstract

Gene selection has become a common task in most gene expression studies. The objective of such research is often to identify the smallest possible set of genes that can still achieve good predictive performance. To do so, many of the recently proposed classification methods require some form of dimension-reduction of the problem which finally provide a single model as an output and, in most cases, rely on the likelihood function in order to achieve variable selection. We propose a new prediction-based objective function that can be tailored to the requirements of practitioners and can be used to assess and interpret a given problem. Based on cross-validation techniques and the idea of importance sampling, our proposal scans low-dimensional models under the assumption of sparsity and, for each of them, estimates their objective function to assess their predictive power in order to select. Two applications on cancer data sets and a simulation study show that the proposal compares favorably with competing alternatives such as, for example, Elastic Net and Support Vector Machine. Indeed, the proposed method not only selects smaller models for better, or at least comparable, classification errors but also provides a set of selected models

---

<sup>0\*</sup> The first two authors are Joint First Authors.

instead of a single one, allowing to construct a network of possible models for a target prediction accuracy level.

**Keywords:** Biomarker selection, Genomic networks, Disease classification, Breast cancer, Acute leukemia, Model averaging

# 1 Introduction

1 Gene selection has become a common task in most gene expression studies.  
2 The problem of assigning tumours to a known class is an example that is of  
3 particular importance and has received considerable attention in the last ten years.  
4 Conventional class prediction methods of leukemia or other cancers are in general  
5 based on microscopical examination of stained tissue specimens. However, such  
6 methods require highly trained specialists and are subjective (Tibshirani *et al.*,  
7 2002).

8 To avoid these drawbacks, many automatic selection methods have been  
9 proposed recently. The goal of these methods is often to identify the smallest  
10 possible set of genes that can still achieve good predictive performance (Díaz-  
11 Uriarte and De Andres, 2006), although this is not necessarily the only criterion  
12 based on which model (gene) selection is carried out (see, for example Leng *et al.*,  
13 2006). However, these methods have the advantage of being objective and have  
14 improved the correct classification rate in various cases. Among the different  
15 methodologies brought forward in this context we can find those proposed by  
16 Tibshirani *et al.* (2002), Dudoit *et al.* (2002), Zhu and Hastie (2004), Zou and  
17 Hastie (2005). See also Díaz-Uriarte and De Andres (2006) and the references  
18 therein for other approaches.

19 Nonetheless, many of these methods do not necessarily respond to the needs of  
20 practitioners and researchers when they approach the gene selection process. First  
21 of all, many of them have to rely on some form of size reduction and often require  
22 a subjective input to determine the dimension of the problem. Also, many of these  
23 methods often provide a single model as an output whereas genes interact inside  
24 biological systems and can be interchangeable in explaining a specific response.  
25 The idea of interchangeability of genes in explaining responses appears for instance  
26 in Kristensen *et al.* (2012). These authors use the PARADIGM algorithm of Vaske  
27 *et al.* (2010) to combine mRNA expression and DNA copy number in order to  
28 construct clusters of patients that provide the best predictive value. The resulting  
29 clusters can be seen as being characterized by different significantly expressed  
30 genes and we can refer to their interactive structure as *paradigmatic* networks.

31 Another issue of most existing gene selection methods is their reliance on the  
32 likelihood function, or a penalized version of it, as a means to develop a selection  
33 criterion. However, the likelihood function may not necessarily be the quantity  
34 that users are interested in as they may want to target some other kind of loss  
35 function such as, for example, the classification error. Of course, maximizing

36 the likelihood function is not typically the same as minimizing a particular loss  
37 function. Moreover, adapting these methods to handle missing or contaminated  
38 data is not straightforward. This has limited the applicability and reliability of  
39 these methods in many practical cases.

40 To eliminate the limitations of the gene selection procedures described above,  
41 this paper proposes an objective function for out-of-sample predictions that can  
42 be tailored to the requirements of practitioners and researchers. This is achieved  
43 by enabling them to select a criterion according to which they would like to assess  
44 and/or interpret a given problem. However, the optimization of such a criterion  
45 function is typically not an easy task since the function can be discontinuous,  
46 non-convex and would require computationally intensive techniques. To tackle  
47 this issue, we propose a solution using a different approach based on a procedure  
48 that resembles *importance sampling*. This new approach provides a general and  
49 flexible framework for gene selection as well as for other model selection problems.

50 The advantages of this proposal are multiple:

- 51 • **Flexibility:** It allows the users to specify a criterion that can be tailored to  
52 the specific problem setting. It is able to handle different kinds of responses,  
53 problems of missing and contaminated data, multicollinearity, etc.
- 54 • **Prediction Power:** The result of the procedure is a set of models with  
55 high predictive power with respect to the specified criterion. It is especially  
56 suitable in selecting genes and models to achieve accurate predictions.
- 57 • **Dimension-reduction:** It can provide an assessment of the dimension of  
58 the problem because it greatly reduces the number of necessary covariates  
59 and eases the interpretation without requiring any preliminary size reduction.
- 60 • **Network-building:** With the reduced model size, it preserves the capacity  
61 to build gene-networks to provide a more general view of the potential  
62 paradigmatic structures of the genetic information.

63 This last aspect is of great interest for gene selection since this list can provide  
64 insight into the complex mechanisms behind different biological phenomena.  
65 Different cases, some of which can be found in Section 4, indicate that this method  
66 appears to outperform other methods in terms of criteria minimization while,  
67 at the same time, selects models of considerably smaller dimension which allow  
68 improved interpretation of the results. The set of selected models can naturally  
69 be viewed as a network of possible structures of genetic information. We call

70 this a paradigmatic network. In Section 4 we give an example of a graphical  
 71 representation of such networks based on the analysis of one of two cancer data  
 72 sets which are discussed therein.

73 In this paper we first describe and formalize the proposed approach within the  
 74 model selection statistical framework in Section 2. In Section 3 we illustrate the  
 75 techniques and algorithms used to address the criterion minimization problem  
 76 highlighted in Section 2. The performance of our approach is then illustrated on  
 77 two data sets concerning leukemia classification (Golub *et al.*, 1999) and breast  
 78 cancer classification (Chin *et al.*, 2006) in Section 4. We conclude the paper in  
 79 Section 6 by summarizing the benefits of the new approach and providing an  
 80 outlook on other potential applications that can benefit from this methodology.

## 81 2 Approach

82 To introduce the proposed method, let us first define some notation which will be  
 83 used throughout this paper:

- 84 1. Let  $\mathcal{J}_f = \{1, 2, \dots, p\}$  be the set of indices for  $p$  potential covariates included  
 85 in the  $n \times p$  matrix  $\mathbf{X}$ . We allow  $\mathbf{X}$  to include a vector of 1s.
- 86 2. Let  $\mathcal{J} = \mathcal{P}(\mathcal{J}_f) \setminus \emptyset$ ,  $|\mathcal{J}| = 2^p - 1$ , be the power set including all possible  
 87 models that can be constructed with the  $p$  covariates excluding the empty  
 88 set.
- 89 3. Let  $j \in \mathcal{J}$  be a model belonging to the above mentioned power set.
4. Let  $\beta^j \in \mathbb{R}^p$  be the parameter vector for model  $j$ , i.e.

$$\beta_k^j = \begin{cases} \beta_k & \text{if } k \in j \\ 0 & \text{if } k \notin j \end{cases}$$

90 where  $\beta_k^j$ ,  $\beta_k$  are respectively the  $k$ th element of  $\beta^j$  and  $\beta$ , with  $\beta =$   
 91  $(\beta_1, \dots, \beta_p)^T \in B \subseteq \mathbb{R}^p$ .

92 Keeping this notation in mind, for a given model  $j \in \mathcal{J}$  we have that

$$\mathbb{E}[Y|\mathbf{X}] = g(\mathbf{X}, \beta^j), \tag{1}$$

93 where  $\mathbb{E}[\cdot]$  is the expectation operator and  $g(\cdot, \cdot)$  is a link function known up to the  
 94 parameter vector  $\beta^j \in \mathbb{R}^p$ . Models of the form (1) are very general and include  
 95 all parametric models and a large class of semiparametric models when  $g(\cdot, \cdot)$  is  
 96 not completely known.

97 We assume that for a fixed  $j$ , based on a specific choice for model (1) with  
 98 corresponding parameter vector  $\beta^j$  and given a new covariate vector  $\mathbf{X}_0$ , the user  
 99 can construct a prediction  $\widehat{Y}(\mathbf{X}_0, \beta^j)$ . To assess the quality of this prediction  
 100 we assume that we have a divergence measure available which we denote as  
 101  $D\{\widehat{Y}(\mathbf{X}_0, \beta^j), Y_0\}$ . The only requirement imposed on the divergence measure is  
 102 that it satisfies the property of positiveness, i.e.

$$D(u, v) > 0 \text{ for } u \neq v$$

$$D(u, v) = 0 \text{ for } u = v.$$

103 With this property being respected, the divergence measure can arbitrarily be  
 104 specified by the user according to the interest in the problem. Examples of such  
 105 divergence measures include the  $L_1$  loss function

$$D\{\widehat{Y}(\mathbf{X}_0, \beta^j), Y_0\} = |\widehat{Y}(\mathbf{X}_0, \beta^j) - Y_0|$$

106 or an asymmetric classification error

$$D\{\widehat{Y}(\mathbf{X}_0, \beta^j), Y_0\} = I\{\widehat{Y}(\mathbf{X}_0, \beta^j) = 1, Y_0 = 0\}w_1$$

$$+ I\{\widehat{Y}(\mathbf{X}_0, \beta^j) = 0, Y_0 = 1\}w_2.$$

107 where  $w_1, w_2 \geq 0$ . The latter is for a Bernoulli response and is typically an  
 108 interesting divergence measure when asymmetric classification errors have to be  
 109 considered. Indeed, in most clinical situations, the consequences of classification  
 110 errors are not equivalent with respect to the direction of the misclassification. For  
 111 instance, the prognosis and the treatment of Estrogen Receptor (ER) positive  
 112 Breast Cancers (BC) are quite different from those of ER negative ones. Indeed,  
 113 if a patient with ER negative is treated with therapies designed for patients with  
 114 ER positive, the consequence is much more severe than if this were done the  
 115 other way round because of the excessive toxicities and potentially severe side  
 116 effects. It therefore makes sense to give different values to  $w_1$  and  $w_2$ . By defining  
 117  $w_1 > w_2$  we would take these risks into account, where  $w_1$  would be the weight for  
 118 a misclassification from ER negative to ER positive BC and  $w_2$  for the opposite  
 119 direction. Weight values can be modulated according to the current medical  
 120 knowledge and the clinical intuition of the physicians.

121 Considering this divergence measure  $D(\cdot, \cdot)$ , we are consequently interested in  
 122 finding the best models within the general class given in (1). To do so, we would  
 123 ideally aim at solving the following risk minimization problem :

$$\widehat{\boldsymbol{\beta}}^j \in \mathcal{B} \equiv \underset{j \in \mathcal{J}}{\operatorname{argmin}} \underset{\boldsymbol{\beta}^j}{\operatorname{argmin}} \mathbb{E}_0 \left[ D \left\{ \widehat{Y}(\mathbf{X}_0, \boldsymbol{\beta}^j), Y_0 \right\} \right], \quad (2)$$

124 where  $\mathbb{E}_0$  denotes the expectation on the new observation  $(Y_0, \mathbf{X}_0)$ . Let  $j_0$  denote  
 125 the models with the smallest cardinality among all  $\widehat{\boldsymbol{\beta}}^j \in \mathcal{B}$ . Note that there  
 126 could be more than one model with the same prediction property and of the  
 127 same size, hence  $j_0$  could contain more than one model. Let us define the models  
 128 corresponding to  $j_0$  as the “true” models. Thus, our “true” models are essentially  
 129 the most parsimonious models that minimize the expected prediction error.

130 The optimization problem in (2) is typically very difficult to solve. First of all,  
 131 supposing we do not consider interaction terms, the outer minimization would  
 132 require to compare a total of  $2^p - 1$  results, each a result of the inner minimization  
 133 problem. In addition, each of the  $2^p - 1$  inner minimization problems is also very  
 134 hard to solve, even if the risk  $\mathbb{E}_0[D\{\widehat{Y}(\mathbf{X}_0, \boldsymbol{\beta}^j), Y_0\}]$  were a known function of  $\boldsymbol{\beta}^j$ .  
 135 Indeed, the inner minimization problem is in general non-convex and could be  
 136 combinatorial, implying that the minimizer might not be unique. For example,  
 137 when  $D(\cdot, \cdot)$  is the classification error, this problem is combinatorial by nature. In  
 138 practice, the computational challenge is even greater because the risk function  
 139  $\mathbb{E}_0[D\{\widehat{Y}(\mathbf{X}_0, \boldsymbol{\beta}^j), Y_0\}]$  is a function of  $\boldsymbol{\beta}^j$  without explicit form and needs to be  
 140 approximated.

141 We propose to estimate  $\mathbb{E}_0[D\{\widehat{Y}(\mathbf{X}_0, \boldsymbol{\beta}^j), Y_0\}]$  via an  $m$ -fold cross-validation  
 142 (typically  $m = 10$ ) repeated  $K$  times. More specifically, for a sample of size  $n$ ,  
 143 we repeat the following procedure  $K$  times. At the  $k$ th repetition, we randomly  
 144 permute the sets  $(\mathbf{X}_r, Y_r)$ , with  $r$  indexing a row of the data (i.e.  $r = 1, \dots, n$ ),  
 145 and then select  $\lfloor n/m \rfloor$  observations from the permuted data to form “test” data  
 146 sets, subindexed  $(i, l)$  (with  $i = 1, \dots, \lfloor n/m \rfloor$  and  $l = 1, \dots, m$ ) and superindexed  
 147  $k = 1, \dots, K$ , i.e.  $(\mathbf{X}_{(i,l)}^k, Y_{(i,l)}^k)$ . Therefore  $l$  indicates the test set,  $i$  indicates the  
 148 observation within this test set and  $k$  represents the repetition (associated with  
 149 a certain permutation of the data). The classical 10-fold cross-validation, for  
 150 example, is obtained by defining  $K = 1$  and  $m = 10$ . Given this, the estimated  
 151 risk is

$$\widehat{\mathbb{E}}_0 \left[ D \left\{ \widehat{Y}(\mathbf{X}_0, \boldsymbol{\beta}^j), Y_0 \right\} \right] = \frac{1}{\lfloor n/m \rfloor m K} \sum_{k=1}^K \sum_{l=1}^m \sum_{i=1}^{\lfloor n/m \rfloor} D \left\{ \widehat{Y}(\mathbf{X}_{(i,l)}^k, \boldsymbol{\beta}^j), Y_{(i,l)}^k \right\}. \quad (3)$$

152 Having approximated the expectation  $\mathbb{E}_0$ , the minimization problem in (2)  
 153 becomes

$$\operatorname{argmin}_{j \in \mathcal{J}} \operatorname{argmin}_{\beta^j} \widehat{\mathbb{E}}_0 \left[ D \left\{ \widehat{Y}(\mathbf{X}_0, \beta^j), Y_0 \right\} \right]. \quad (4)$$

154 Despite the above approximation, the minimization problem remains compli-  
 155 cated for the reasons mentioned earlier. Thus, we further eliminate the inner  
 156 minimization problem in (4) by inserting an estimator  $\widehat{\beta}^j$  obtained indepen-  
 157 dently from the minimization procedure. More specifically, we assume that an  
 158 estimator of  $\beta^j$ , say  $\widehat{\beta}^{j,k}$ , is available based on model (1) and “training” observa-  
 159 tions  $(\mathbf{X}_{\lfloor n/m \rfloor + 1}^k, Y_{\lfloor n/m \rfloor + 1}^k), \dots, (\mathbf{X}_n^k, Y_n^k)$  (i.e. those observations excluded from the  
 160 above mentioned “test” data sets). This estimator can be any available estimator,  
 161 for example, the maximum likelihood estimator (MLE), a moment based estimator,  
 162 or a quantile regression based estimator, etc. (see, for example, [Azzalini, 1996](#);  
 163 [Hall et al., 2005](#); [Koenker, 2005](#)). We then replace the inner minimization in (4)  
 164 directly with the approximate expectation evaluated at  $\widehat{\beta}^{j,k}$ ’s and simplify (4) to

$$\operatorname{argmin}_{j \in \mathcal{J}} \frac{1}{\lfloor n/m \rfloor m K} \sum_{k=1}^K \sum_{l=1}^m \sum_{i=1}^{\lfloor n/m \rfloor} D \{ \widehat{Y}(\mathbf{X}_{(i,l)}^k, \widehat{\beta}^{j,k}), Y_{(i,l)}^k \}. \quad (5)$$

165 The intuition of replacing the inner minimization in (4) with a sample average  
 166 evaluated at an arbitrary estimator is due to the fact that this estimator, under a  
 167 fixed “true” model and regardless of whether this estimator is a standard MLE or  
 168 a minimizer of the divergence measure  $D(\cdot, \cdot)$ , is an approximation to the “true”  
 169 parameter. This means that, consequently, different estimators are “close” to each  
 170 other. As a consequence, the minimization problem in (5) can be considered to be  
 171 a close approximation to  $\min_{\beta^j} \mathbb{E}_0 \left[ D \left\{ \widehat{Y}(\mathbf{X}_0, \beta^j), Y_0 \right\} \right]$ . In fact, using an informal  
 172 law of large numbers argument, as  $n/m \rightarrow \infty$ , then we have that  $\widehat{\beta}^j \xrightarrow{p} \beta^j$ . If in  
 173 addition  $m \rightarrow \infty$  then, under some regularity conditions on  $D(\cdot, \cdot)$ , the averages  
 174 tend to the desired expectation. On the other hand, if instead we consider  $m$  as  
 175 fixed, we would have an unbiased estimator of the expected risk.

176 We now have an optimization problem in (5) which requires a comparison  
 177 of  $2^p - 1$  values and is much easier to solve. To further reduce the number of  
 178 comparisons, the following section describes some procedures and algorithms  
 179 allowing to solve this problem in a more efficient manner.

### 180 3 Heuristic procedure

181 To solve the optimization problem in (5), we propose an approach designed to  
182 have the following three features:

- 183 1. Identify a **set of models** that carry large predictive power instead of a  
184 single “best” model;
- 185 2. Find this set of models within a **reasonable time**, without having to  
186 explore all possible models;
- 187 3. This set achieves **sparsity**, i.e. most of the parameters in  $\beta$  will be fixed at  
188 zero in each of the models in the set.

189 Note that the last feature above reflects the belief that most of the covariates  
190 are irrelevant for the problem under consideration and should be excluded. Indeed,  
191 our method is designed to work effectively if such a sparsity assumption holds,  
192 putting it on the same level of almost all variable selection procedures in the  
193 literature. Moreover, we require the method to have the first feature in order to  
194 increase flexibility in terms of interpretation. Indeed, in many domains such as  
195 gene selection, for example, the aim may not be to find a single model but a set  
196 of variables (genes) that can be inserted in a paradigmatic structure to better  
197 understand the contribution of each of them via their interactions.

198 Given this goal, assume that we have at our disposal an estimate of the  
199 measure of interest  $D(\cdot, \cdot)$  for all possible  $2^p - 1$  models. In this case, our interest  
200 would be to select a set of “best” models by simply keeping the set of models  
201 that have a low discrepancy measure  $D(\cdot, \cdot)$ . It is of course unrealistic to obtain  
202 a discrepancy measure for all models in most practical cases because this would  
203 require a considerable amount of time for computation. Therefore, in order to  
204 achieve the second feature, instead of examining all possible models, we can  
205 randomly sample covariates from  $\mathcal{J}$ . The random sampling needs to be carefully  
206 devised because in practice, for example in gene selection problems, the number  
207 of covariates  $p$  can easily reach thousands or tens of thousands (see examples  
208 in Section 4, where  $p = 7, 129$  and  $p = 22, 215$  respectively). In such situations,  
209  $2^p - 1$  is an extremely large number and the probability of randomly sampling  
210 a “good” set of variables from the  $2^p - 1$  variables is very small. Using the  
211 sparsity property of the problem, we propose to start with the set of variables  $\mathcal{M}_0$   
212 (typically an empty set) and increase the model complexity stepwise. Throughout  
213 this procedure, we ensure that at step  $k$ , the most promising covariates based

214 on the evaluation at step  $k - 1$  are given higher probabilities of being randomly  
 215 drawn. The last idea is in the spirit of “importance sampling” in the sense that  
 216 covariates with more importance based on the previous step are “encouraged” to  
 217 be selected in the current step. Note that by construction we achieve sparsity if  
 218 we stop the stepwise search at models of size  $d_{\max} \ll p$ .

219 More formally, let us first define the set of all possible models of size  $d$  as

$$\mathcal{S}_d = \{(i_1, \dots, i_d) \mid i_1, \dots, i_d \in \mathcal{J}_f; i_1 < \dots < i_d\}.$$

220 We then define the set of promising models,  $\mathcal{S}_d^*$ , as the ones with an estimated  
 221 out-of-sample divergence measure  $D(\cdot, \cdot)$  below a certain estimated  $\alpha$ -quantile.  
 222 The value of  $\alpha$  is user-defined depending on the problem at hand, and is typically  
 223 a small value such as  $\alpha = 1\%$ . The formal definition of this set would then be

$$\mathcal{S}_d^* = \{j \mid j \in \mathcal{S}_d; \widehat{D}_j \leq \widehat{q}_d(\alpha)\},$$

224 where

$$\widehat{D}_j \equiv \frac{1}{\lfloor n/m \rfloor m K} \sum_{k=1}^K \sum_{l=1}^m \sum_{i=1}^{\lfloor n/m \rfloor} D\{\widehat{Y}(\mathbf{X}_{(i,l)}^k, \widehat{\boldsymbol{\beta}}^{j,k}), Y_{(i,l)}\}, \quad (6)$$

225 and  $\widehat{q}_d(\alpha)$  is the  $\alpha$ -quantile of the  $\widehat{D}_j$  ( $j \in \mathcal{S}_d$ ) values issued from  $B$  randomly  
 226 selected models. Finally, we define the set of indices of covariates that are in  $\mathcal{S}_d^*$   
 227 as

$$\mathcal{I}_d^* = \{i \mid i \in j, j \in \mathcal{S}_d^*\}$$

228 whose complement we define as  $\mathcal{I}_d^c$  (i.e. all those covariates that are not included  
 229 in  $\mathcal{I}_d^*$ ).

230 With this approach in mind and using the above notations, to start the  
 231 procedure we assume that we have  $p$  variables from which to select.

232 *A. Initial Step:* We start by adding the number of variables  $d = 1$  to our initial  
 233 variable set  $\mathcal{M}_0$  with the goal of finally obtaining the set  $\mathcal{I}_1^*$ .

234 1. Construct the  $p$  possible one variable models by augmenting  $\mathcal{M}_0$  with  
 235 each of the  $p$  available variables.

236 2. Compute  $\widehat{D}$  for every model obtained in Step A.1.

- 237           3. From Steps A.1 and A.2, construct the set  $\mathcal{I}_1^*$  using (3). Go to Step B  
 238           and let  $d = 2$ .
- 239    B. *General Step:* We define here the general procedure to construct  $\mathcal{I}_d^*$  for  
 240            $2 \leq d \leq d_{\max}$ .
- 241           1. Augment  $\mathcal{M}_0$  with  $d$  variables as follows:
- 242                   (i) Randomly select a set, either set  $\mathcal{I}_{d-1}^*$  with probability  $\pi$  or its  
 243                   complement  $\mathcal{I}_{d-1}^c$  with probability  $1 - \pi$ .
- 244                   (ii) Select one variable uniformly at random and without replacement  
 245                   from the set chosen in Step (i) and add this variable to  $\mathcal{M}_0$ .
- 246                   (iii) Repeat Steps (i) and (ii) until  $d$  variables are added to  $\mathcal{M}_0$ .
- 247           2. Construct a model of dimension  $d$  using the  $d$  variables selected in Step  
 248           B.1. Repeat Step B.1  $B$  times to construct  $B$  such models.
- 249           3. From Steps B.1 and B.2, construct the set  $\mathcal{I}_d^*$  according to (3). If  
 250            $d < d_{\max}$ , go to Step B and let  $d = d + 1$ , otherwise exit algorithm.

251           Once the algorithm is implemented, the user obtains an out-of-sample discrep-  
 252           ancy measure for all evaluated models. Given that the goal is to obtain a set  
 253           of models  $\mathcal{S}_d^*$  with high predictive power, the discrepancy measure delivers the  
 254           criterion based on which it is possible to determine the optimal model dimension  
 255           and the corresponding network structure.

## 256    3.1    Practical Considerations

257           The algorithm described above lays out the basic procedure to solve the problem  
 258           in (2). However, as many other heuristic selection procedures, there are a series  
 259           of “hyper-parameters” to be determined and certain aspects to be considered.  
 260           In the following paragraphs we will discuss some of these issues arising when  
 261           implementing our algorithm in practice.

### 262    3.1.1    Choice of algorithm inputs

263           The parameters  $d_{\max}$ ,  $B$ ,  $\alpha$  and  $\pi$  of the above algorithm are to be fixed by  
 264           the user. As mentioned earlier  $d_{\max}$  represents a reasonable upper bound for  
 265           the model dimension which is constrained to  $d_{\max} \leq l$ , where  $l$  depends on the  
 266           limitations of the estimation method and is commonly the sample size  $n$ . As

267 for the parameter  $B$ , a larger value is always preferable to better explore the  
 268 covariate space. However, a larger  $B$  implies heavier computations, hence a rule  
 269 of thumb that could be used is to choose this parameter such that  $p \leq B \leq \binom{p}{2}$ .  
 270 As mentioned earlier, the parameter  $\alpha$  should define a small quantile, typically  
 271 1%. Finally,  $\pi$  determines to what extent the user assigns importance to the  
 272 variables selected at the previous step. Given that  $d_{\max} \ll p$  and  $\alpha$  is small, we  
 273 will typically have that  $|\mathcal{I}_{d-1}^*| < |\mathcal{I}_{d-1}^c|$ . In this setting, a choice of  $\pi = 0.5$  for  
 274 example would deliver a higher probability for the variables in  $\mathcal{I}_{d-1}^*$  to be included  
 275 in  $\mathcal{I}_d^*$ . All other parameters being equal, increasing the value of  $\pi$  would decrease  
 276 the probability of choosing a variable in  $\mathcal{I}_{d-1}^c$  and vice versa. Moreover, we discuss  
 277 in Appendix A how the proposed algorithm can be adjusted to situations where  $p$   
 278 is either small or very large.

279 As a final note, it is also possible for the initial model  $\mathcal{M}_0$  to already contain  
 280 a set of  $p_0$  covariates which the user considers to be essential for the final output.  
 281 In this case, the procedure described above would remain exactly the same since  
 282 the procedure would simply select from the  $p$  covariates which are not in the  
 283 user-defined set and the final model dimension would simply be  $p_0 + d$ .

### 284 3.1.2 Model Dimension and Network Building

285 The final goal of the algorithm is to find a subset of models of dimension  $d^*$  that  
 286 in some way minimize the considered discrepancy. A possible solution would  
 287 be to select the set of models  $\mathcal{S}_{d^*}^*$  such that  $d^* = \min_{d \in \{1, \dots, d_{\max}\}} q_d(\alpha)$ . However,  
 288 the quantity  $q_d(\alpha)$  is unknown and replaced by its estimator  $\hat{q}_d(\alpha)$ . Due to  
 289 this, a solution that might be more appropriate would be to consider a testing  
 290 procedure to obtain  $d^*$  taking into account the variability of  $\hat{q}_d(\alpha)$ . For example,  
 291 we could find the dimension  $d^*$  such that we cannot reject the hypothesis that  
 292  $\hat{q}_{d^*}(\alpha) = \hat{q}_{d^*+1}(\alpha)$ . Thus we sequentially test whether  $\hat{q}_{j+1}$  is smaller than  $\hat{q}_j$   
 293 for  $j = 1, \dots, d_{\max}$ . As long as the difference is significant we increment  $j$  by one  
 294 unit, otherwise the minimum is reached and  $d^* = j$ .

295 The type of test and its corresponding rejection level are determined by the  
 296 user based on the nature of the divergence measure. For example, if we take the  
 297  $L_1$  loss function as a divergence, one could opt for the Mann-Whitney test or if  
 298 the loss function is a classification error (as in the applications in Section 4), one  
 299 could choose the binomial test or other tests for proportions. The rejection level  
 300 will depend, among others, on the number of tests that need to be run, typically  
 301 less than  $d_{\max} - 1$ , and need to be adjusted using, for example, the Bonferroni  
 302 correction. Finally, once the set  $\mathcal{S}_{d^*}^*$  is obtained, the user may still want to “filter”

303 the resulting models. Indeed, the number of models in the solution  $\mathcal{S}_{d^*}^*$  may be  
 304 large and the corresponding divergence estimates may vary considerably from  
 305 model to model. Since these divergence measures are estimators, we again propose  
 306 a multiple testing procedure to reduce the number of models in  $\mathcal{S}_{d^*}^*$ . Before doing  
 307 so, we eliminate redundant models, thereby making sure that every model is  
 308 included only once. Then, we start the testing procedure with an empty set  
 309  $\mathcal{S}_{d^*}^0 = \emptyset$  to which we add the model (or one of the models) that has the minimum  
 310 divergence measure estimate, denoted  $\widehat{D}_{j_{\min}}$ , where  $j_{\min} \in \mathcal{S}_{d^*}^*$  denotes this model.  
 311 Then for every model  $j \in \mathcal{S}_{d^*}^* \setminus j_{\min}$ , we test whether  $\widehat{D}_j$  is greater than  $\widehat{D}_{j_{\min}}$ . We  
 312 add the model to  $\mathcal{S}_{d^*}^0$  if the difference is not significant and stop adding models as  
 313 soon as the test deems that the divergence of the next model is indeed larger. By  
 314 doing so we finally obtain  $\mathcal{S}_{d^*}^0 \subseteq \mathcal{S}_{d^*}^*$  which is the set containing the models (and  
 315 hence covariates) which can be interpreted in a paradigmatic network. Generally  
 316 speaking, this network can be built starting from the most frequent covariate(s)  
 317 present in  $\mathcal{S}_{d^*}^0$  (we call these “hubs”) and, subsequently, connecting these with  
 318 the most frequent covariates included in the models with the previous hubs. This  
 319 can be continued until the number of connected hubs is equal to  $d^*$ .

## 320 3.2 Related literature

321 Some of the ideas put forth in this work have also been considered in the literature.  
 322 An extensive survey of the related works goes beyond the scope of this paper.  
 323 Here we briefly describe some of the connections to three main ideas that have  
 324 been explored to this point.

325 The first one is recognizing that practitioners might aim to minimize some  
 326 criterion that differs from likelihood-type losses. An interesting paper illustrating  
 327 this point is [Juang \*et al.\* \(1997\)](#) in the context of speech recognition. For their  
 328 classification problem, these authors propose to minimize a “smoothed” version of  
 329 the decision rule used for classification. The advantage of this procedure is that  
 330 it yields better misclassification errors than using pure likelihood based criteria  
 331 which intrinsically fit a distribution to the data. In the approach presented in this  
 332 work we also deliver an approximate solution but, as opposed to approximating  
 333 the problem and solving the latter in an exact manner as in [Juang \*et al.\* \(1997\)](#), we  
 334 define the exact problem and try to approximately minimize the misclassification  
 335 error through our algorithm.

336 Secondly, there is a large literature that uses stochastic search procedures to  
 337 explore the space of candidate models. Influential work in this direction includes

338 [George and McCulloch \(1993\)](#) and [George and McCulloch \(1997\)](#) who postulate  
339 hierarchical Bayesian models. In their set-up, subsets of promising predictors  
340 form models with higher posterior probabilities. An interesting application of  
341 this framework for disease classification using gene expression data is the work of  
342 [Yang and Song \(2010\)](#). [Cantoni \*et al.\* \(2007\)](#) also consider a random exploration  
343 of the space of possible models, but avoiding the Bayesian formulation of [George  
344 and McCulloch \(1993\)](#). Their approach defines a probability distribution for the  
345 various candidate models based on a cross-validated prediction error criterion and  
346 then uses a Markov Chain Monte-Carlo method to generate a sample from this  
347 probability distribution. An important feature of the stochastic search implied by  
348 our algorithm is that it is a greedy method, while the aforementioned methods are  
349 not. The typical forward/backward greedy algorithms proposed in the literature  
350 are not random, while existing stochastic procedures are not greedy. Thus, the  
351 combination of greedy approach and random search approach seems to be new  
352 (see for instance [Zhang, 2011](#), for some theory on greedy algorithms in sparse  
353 scenarios).

354 Finally, other authors have also considered providing a set of interesting models  
355 as opposed to a single “best” model. The stochastic search procedures mentioned  
356 in the above paragraph can naturally be used to obtain a group of interesting  
357 models. For example, [Cantoni \*et al.\* \(2007\)](#) consider a set of best indistinguishable  
358 models in terms of prediction. Random forests can be used to select variables and  
359 account for the stability of the chosen model as in [Díaz-Uriarte and De Andres  
360 \(2006\)](#). These methods can also be used to construct a set of interesting models.

## 361 4 Case Studies

362 In this section we provide an example of how the methodology proposed in this  
363 paper selects and groups genes to explain, describe and predict specific outcomes.  
364 We focus on the data-set (hereinafter *leukemia*) which collects information on  
365 Acute Myeloid Leukemia (AML) and Acute Lymphoblastic Leukemia (ALL) and  
366 is frequently used as an example for gene selection procedures. Indeed, [Golub  
367 \*et al.\* \(1999\)](#) were among the first to use this data to propose a gene selection  
368 procedure which was then followed up by other proposals that used the same  
369 data to compare their performance. We will use this data-set to underline the  
370 features and advantages of the proposed method. A second data-set concerning  
371 the research on breast cancer (presented in [Chin \*et al.\* \(2006\)](#)) is analysed in

372 Appendix C to show the outputs of the proposed method from another example<sup>1</sup>.

373 The analysis of these data-sets focuses both on the advantages of the proposed  
374 methodology and the biological interpretation of the outcomes. One of the goals  
375 of our method is to help decipher the complexity of biological systems. We will  
376 take on an overly simplified view of the cellular processes in which we will assume  
377 that one biomarker maps to only one gene that in turn has only one function.  
378 Although this assumption is not realistic, it allows us to give a straightforward  
379 interpretation of the selected models or “networks” which can therefore provide an  
380 approximate first insight into the relationships between variables and biomarkers  
381 (as well as between the biomarkers themselves). We clarify that we do not claim  
382 any causal nature in the conclusions we present in these analyses but we believe  
383 that the selected covariates can eventually be strongly linked to other covariates  
384 that may have a more obvious and direct interpretation for the problem at hand.  
385 Finally, the data-set has binary outcomes (as does the data-set in Appendix C),  
386 hence we will make use of the Classification Error (CE) as a measure of prediction  
387 performance and we will not assign weights to a given prediction error. This  
388 means that misclassification errors are given the same weight, in the sense that a  
389 false positive prediction (e.g. predicted “presence” when the truth is “absence”)  
390 is considered as undesirable as a false negative prediction. However, our method  
391 can consider also divergence measures based on unequal weights as highlighted in  
392 Section 2.

## 393 4.1 Acute Leukemia

394 Golub *et al.* (1999) were among the first to propose an automatic selection method  
395 for cancer classification and demonstrated the advantages of using such a method.  
396 One of the main applications of their method was on the *leukemia* data-set in  
397 which information regarding 72 patients is included, namely their type of leukemia  
398 (25 patients with AML and 47 patients with ALL) and 7,129 gene expressions  
399 used as explanatory variables to distinguish between two types of leukemia. As  
400 explained in Golub *et al.* (1999) this distinction is critical for successful treatment  
401 which substantially differs between classes. In fact, although remissions can be  
402 achieved using any of these therapies, cure rates are markedly increased and  
403 unwarranted toxicities are avoided when targeting the specific type of leukemia  
404 with the right therapy.

---

<sup>1</sup>The *Acute Leukemia* (Section 4.1) and the *Breast Cancer* (Appendix C) data-sets are made available in the R package “datamicroarray”.

### Quantile $\hat{q}_j$ for $\alpha = .01$

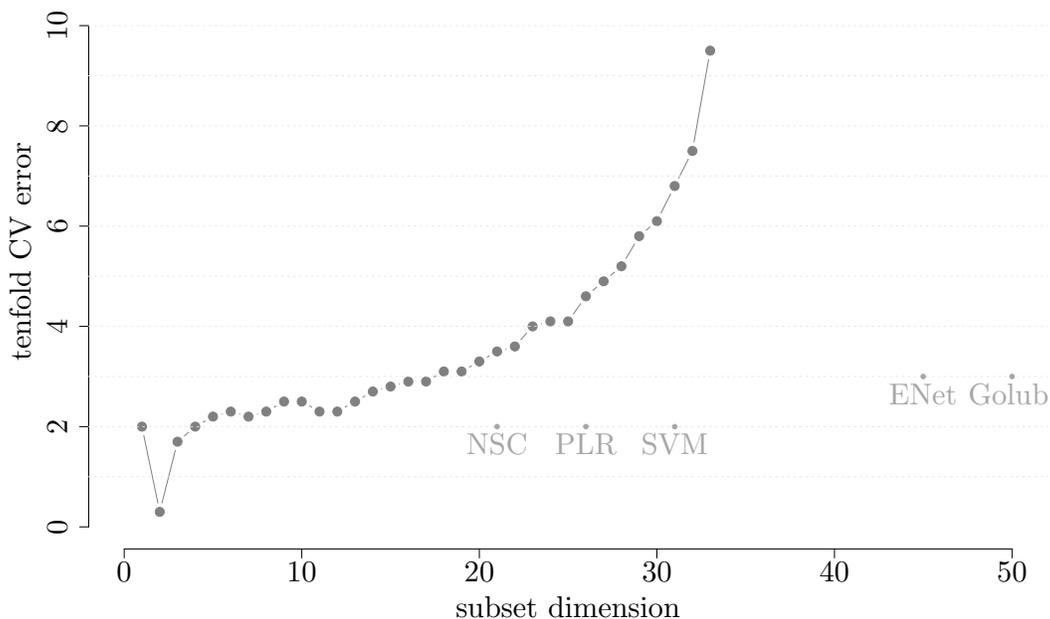


FIGURE 1: *Number of covariates vs.  $\hat{D}$  on leukemia cancer classification training set. The names are abbreviations for other selection method referred to in Table 1.*

#### 4.1.1 Statistical analysis

In order to understand how our proposed methodology performs compared to existing ones, we split the *leukemia* data into the same training set (38 patients) and test set (34 patients) as in the original work by Golub *et al.* (1999). We employ our method on the training set to understand the dimension of the model and to select the most relevant genes. Setting  $\alpha = 0.01$ , the corresponding observed quantile of the 10-fold cross-validation CE ( $\hat{D}$ ) is shown in Figure 1. It can be seen that the error immediately decreases to almost zero when using two covariates instead of one, after which it roughly monotonically increases, suggesting that the optimal model dimension is two.

In Figure 1 we also plotted the performance of the other selection methods used on this training data which are represented by labelled dots reporting the acronyms of these methods that are listed in Table 1. These cross-validation errors are taken from Zou and Hastie (2005) in the same setting in which we ran

419 the proposed method. However, another table in which the competing methods  
420 were ran using currently available software is presented in Appendix B where the  
421 conclusions in terms of comparison do not differ from those presented in Table  
422 1<sup>2</sup>. Indeed, the approach proposed in this work compares favourably to all other  
423 methods in terms of prediction power since they lie under the curve to the right  
424 of its minimum indicating that, compared to our method, they select models of  
425 considerably higher dimensions without achieving the same degree of performance  
426 in terms of CE. Therefore, for this particular case, our method outperforms the  
427 other methods. The sparsity and tenfold CV error are further illustrated in Table  
428 1, where we also present the average prediction error on the test data. Considering  
429 the latter, it can be seen how the performance of the different methods are similar  
430 but the proposed method (which we refer to as *Panning*) is able to achieve the  
431 same performance by selecting models of a considerably lower dimension. As a  
432 final note to the table, the last line reports the performance of model averaging.  
433 Indeed, if the interest lies in predicting, the algorithm of Section 3 provides a  
434 set of models whose CE is below a given quantile  $\alpha$ . The predictions of these  
435 models can be used in the spirit of model averaging where a general prediction  
436 can be obtained by taking the average of predictions of the selected set of models.  
437 The proposed methodology can therefore be potentially seen as a bridge between  
438 model selection and model averaging.

439

440 Once this procedure is completed, we can create a gene network to facilitate  
441 interpretation. This is a direct benefit of our method which does not deliver  
442 a single model after the selection process but provides a series of models that  
443 can be linked to each other and interpreted jointly. Indeed, the existence of a  
444 single model that links the covariates to the explained variable is probably not  
445 realistic in many settings, especially for gene classification. For this reason, the  
446 frequency with which each gene is included within the selected models and with  
447 which these genes are coupled with other genes provides the building block to  
448 create an easy-to-interpret gene network with powerful explanatory and predictive  
449 capacities. A graphical representation of this gene network can be found in  
450 Figure 2 where the size of a disk represents the frequency with which a particular  
451 biomarker is included in the selected models, and the line connecting the disks  
452 indicates the biomarkers that are included in the same model. Since the model  
453 dimension in this case is two, each biomarker is connected with only one other  
454 biomarker and, as can be observed, the proposed method identifies three main

---

<sup>2</sup>The use of the software making available the competing methods is described in Section 5.

Method	Tenfold CV error	Test error	Number of genes
Golub	3/38	4/34	50
Support vector machine (with recursive feature elimination)	2/38	1/34	31
Penalised logistic regression (with recursive feature elimination)	2/38	1/34	26
Nearest shrunken centroids	2/38	2/34	21
Elastic net	3/38	0/34	45
Panning Algorithm (107)			
Model a	0/38	2/34	2
Model b	0/38	2/34	2
Model c	0/38	2/34	2
[...]			
Model averaging		2/34	2

TABLE 1: *Summary of Leukemia classification results. The table is taken from Zou and Hastie (2005) where we added the Panning Algorithm. We obtained a total of 107 models of size 2 (109 different biomarkers) using a probability  $\alpha = 0.01$ ,  $B = 20'000$  bootstrap replicates, a selection probability  $\pi = 0.5$  with  $D(\cdot, \cdot)$  estimated through tenfold-CV repeated  $K = 10$  times. Models “a” to “c” are three examples out of the 107 models. All 107 models have a tenfold-CV error of 0. The best test error is 2 and the worst is 12. For model averaging all models are equally weighted.*

455 “hubs” for the networks (green disks) generating three networks. Appendix B also  
456 reports a related table where the biomarkers are listed according to their position  
457 in the model. These positions represent families of biomarkers (or genes) whose  
458 members are interchangeable. By the latter we mean that, given the presence of  
459 biomarkers from other families, specific biomarkers can be replaced by another  
460 biomarker from within the same family without losing predictive power. This is  
461 the idea behind finding a paradigmatic network for gene selection purposes. In  
462 the following paragraph we provide a summary biological interpretation of the  
463 the three main biomarkers (i.e. the most frequent in the selected models) which  
464 we call “hubs” from which the networks start.

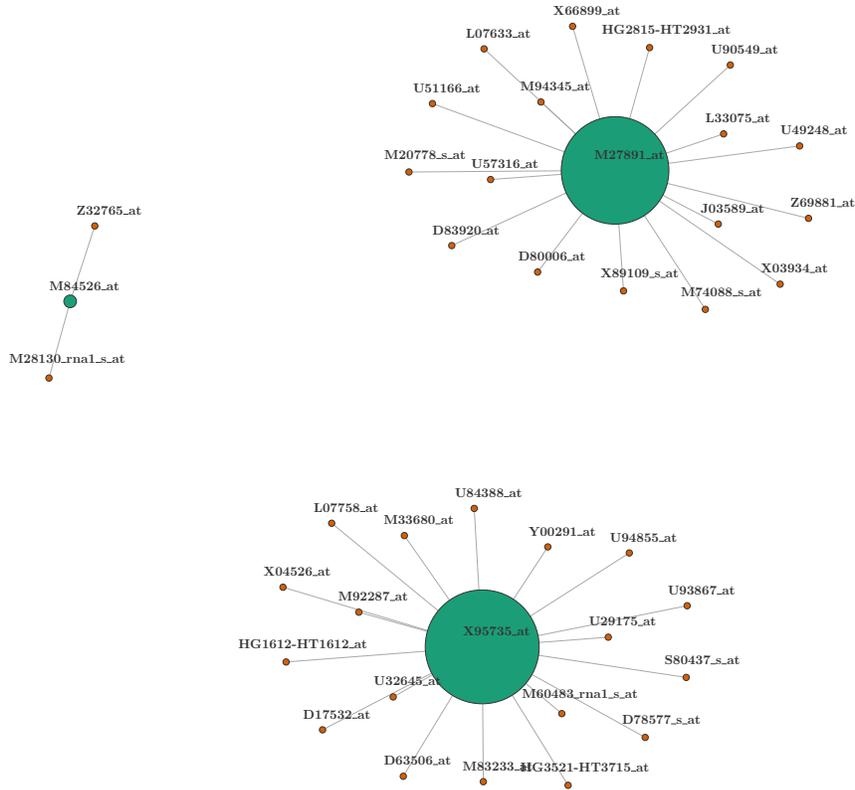


FIGURE 2: *Network representation of biomarkers selected from the leukemia data-set. Colors represent the position of covariates within the model: green for first position (hub) and orange for second. The width of the connecting lines is proportional to the frequency with which two biomarkers appear in the same model. The size of the disk is proportional to the frequency with which a biomarker is present within the selected set of models.*

#### 465 4.1.2 Biological interpretation

466 The three hubs that were identified are the following:

- 467 1. Cystatin C: a secreted cysteine protease inhibitor abundantly expressed in  
468 body fluids (see [Xu \*et al.\*, 2015](#));
- 469 2. Zyxin: a zinc-binding phosphoprotein that concentrates at focal adhesions  
470 and along the actin cytoskeleton;
- 471 3. Complement factor D: a rate-limiting enzyme in the alternative pathway of  
472 complement activation (see [White \*et al.\*, 1992](#)).

473 In the current state of knowledge about acute leukemia, these three hubs appear  
474 to make sense from a biological viewpoint. Cystatin C is directly linked to many  
475 pathologic processes through various mechanisms and recent studies indicate that  
476 the roles of Cystatin C in neuronal cell apoptosis induction include decreasing  
477 B-cell leukemia-2 (BCL-2) whose deregulation is known to be implicated in  
478 resistant AML (see [Sakamoto \*et al.\*, 2015](#)). Zyxin is a protein that interacts with  
479 Vasodilator-stimulated phosphoprotein (VASP) with both being involved in cellular  
480 adhesion and motility. VASP interacts with ABL (breakpoint cluster region-  
481 abelson) and is a substrate of the BcrAbl oncoprotein which drives oncogenesis in  
482 patients with chronic myeloid leukemia (CML) due to a constitutive activation  
483 of tyrosine kinase activity (see [Bernusso \*et al.\*, 2015](#)). Further results suggest  
484 that the phosphorylation and dephosphorylation cycle of VASP by the Abi-1-  
485 bridged mechanism regulates association of VASP with focal adhesions, which  
486 may regulate adhesion of Bcr-Abl-transformed leukaemic cells (see [Masahiro \*et al.\*,  
487 2012](#)). Finally, Complement factor D, together with several other components  
488 of both the classical and alternative complement cascade, is primarily expressed  
489 through both adipocytes and monocytes-macrophages in human subjects (see  
490 [White \*et al.\*, 1992](#); [Gabrielsson \*et al.\*, 2003](#)). A recent review in [Ratajczak \(2014\)](#)  
491 has stressed the role of the complement cascade as a trigger for hematopoietic  
492 stem cells from bone marrow into blood.

493 The interpretation of the network can be carried out through plots or tables  
494 such as those presented in Appendix B where the biomarkers can be grouped to-  
495 gether into clusters having the same biological traits, e.g. transcription/translation  
496 factor activity, DNA repair and catabolism, apoptotic activity. This grouping  
497 allows a more straightforward interpretation of the links between the different  
498 families thereby providing a more general overview of how the elements of the  
499 identified network interact.

## 500 5 Simulation study

501 In this section we present a simulation study whose goal is to highlight the  
502 practical benefits of the proposed method over competing methods frequently  
503 used in genomics. Considering the complexity of simulating from a gene network,  
504 in this setting we limit ourselves to considering the existence of a unique true  
505 model which therefore does not allow to assess one of the features of the proposed  
506 approach which is its network building capacities. Hence, this section specifically  
507 focuses on the prediction power and dimension-reduction ability of the method  
508 and, for the comparison with alternative methods to be fair, we only keep one  
509 model for each simulation replicate. This means that, once the dimension of the  
510 model has been identified, the model with the lowest estimated prediction error is  
511 kept (thereby discarding the other potential candidates).

In this optic, for the simulation study we mimicked the acute leukemia dataset seen in Section 4.1 where we set the true model to be generated by a combination of two gene expressions: Cystatin C ( $X_1$ ) and Thymine-DNA Glycosylase ( $X_2$ ) (see Section 4.1.2). Hence the response  $y^*$  in the simulations is a realization of a Bernoulli random variable with probability parameter  $\gamma$  which is obtained through a logit-link function applied to a linear combination of the two above-mentioned variables plus an intercept (with all  $\beta$  coefficients equaling one) i.e.:

$$\gamma = \frac{1}{1 + \exp^{-(1+X_1+X_2)}}.$$

512 Once the binary response variable  $y^*$  is generated, this is then separated into a  
513 training and a test set of the same size as that in the original data-set (i.e. 38  
514 and 34 respectively).

515 Using the implementation of the proposed algorithm available at the cor-  
516 responding GitHub repository<sup>3</sup>, the results of the simulations based on 100  
517 replications can be found in Table 2 where the median performances are reported.  
518 The proposed algorithm's hyper-parameters are  $\alpha = 0.01$ ,  $B = 20'000$ ,  $\pi = 0.5$   
519 and  $D(\cdot, \cdot)$  based on the classical tenfold-CV ( $K = 1$ ). To select the dimension  $d^*$ ,  
520 we ran the testing procedure described in Section 3.1.2 based on a  $p$ -value of 0.1.  
521 As mentioned earlier, unlike Table 1, we only kept one model of dimension  $d^*$   
522 instead of a set of models. This model was chosen such that it had the minimum  
523 training error and, if this minimum was not unique, then the model was randomly  
524 chosen among those achieving this minimum.

---

<sup>3</sup><https://github.com/SMAC-Group/panning>

525 Concerning the competing methods, these were implemented using existing  
526 R functions with default values. For the Elastic Net we used the R package  
527 “glmnet”, that implements the coordinate descent algorithm described in [Friedman](#)  
528 [et al. \(2010\)](#), using the `cv.glmnet()` function to select the lasso parameter. We  
529 performed a grid search over the values  $\{0.2, 0.4, 0.6, 0.8, 1\}$  for the parameter  
530  $\alpha$  of the Elastic Net and kept the value yielding the best deviance<sup>4</sup>. As for the  
531 Nearest Shrunken Centroids method of [Tibshirani et al. \(2002\)](#) we considered  
532 the R package “pamr”. We applied the function `pamr.train()` on the training  
533 data and took the value of the tuning parameter (threshold) yielding the best  
534 classification. The Support Vector Machines approach with recursive feature  
535 elimination was obtained through the function `fit.rfe()` in the “pathClass”  
536 R package . We used the function `crossval()` to select the soft-margin tuning  
537 parameter discussed in [Chapelle et al. \(2002\)](#). Finally, the penalized  $L_2$  logistic  
538 regression with greedy forward selection and backward deletion was implemented  
539 with the function `step.plr()` of the “stepPlr” R package. Note that this function  
540 also considers all possible interactions among the active variables and it is an  
541 implementation of the methodology proposed by [Park and Hastie \(2008\)](#). Finally,  
542 we used our own implementation for the logistic regression with greedy forward  
543 selection, selecting the model with the minimum BIC.

Method	Tenfold CV error	Test error	Number of genes
Panning Algorithm	0/38 (all)	1/34 (min: 0/34; max: 12/34)	2/7129 (all)
Elastic net	10/38 (min: 9/38; max: 12/38)	0/34 (all)	81/7129 (min: 1; max: 104)
Support vector machine	0/38 (all)	15/34 (all)	4/7129 (min: 4; max: 6)
Penalised logistic regression		12/34 (min: 8/34; max: 12/34)	5/7129 (all)
Nearest shrunken centroids	12/38 (min: 7/38; max: 18/38)	5/34 (min: 0/34; max: 5/34)	30/7129 (min: 3; max: 30)

TABLE 2: *Median performances of selection methods on 100 simulations based on a dataset of 7129 genes where only two are relevant.*

544

<sup>4</sup>Note that the special cases  $\alpha = 0$  and  $\alpha = 1$  correspond respectively to ridge regression and lasso.

545 Table 2 shows how the proposed method compares favorably in terms of  
546 median performance with the respect to the competing methods. Indeed, it is the  
547 best approach (or it is among the best) both in terms of cross-validation error as  
548 in terms test error. Even considering its maximum test error it is comparable to  
549 the other methods, keeping in mind that it selects models of extremely low (and  
550 above all correct) dimensions. For example, the Elastic Net is the without doubt  
551 the best in terms of test error but it selects a unique model of size 81 (in median)  
552 making its genetic interpretation much more complex. On the other hand, the  
553 proposed algorithm selects the correct dimension and, if considering the set of  
554 best models, would deliver a network which is more straightforward to interpret.

## 555 6 Conclusions

556 This paper has proposed a new model selection method with various advantages  
557 compared to existing approaches. Firstly, it allows the user to specify the criterion  
558 according to which they would like to assess the predictive quality of a model. In  
559 this setting, it gives an estimate of the dimension of the problem, allowing the user  
560 to understand how many gene expressions are needed in a model to well describe  
561 and predict the response of interest. Building on this, it provides a paradigmatic  
562 structure of the selected models where the selected covariates are considered as  
563 elements in an interconnected biological network. The approach can handle more  
564 variables than observations without going through dimension-reduction techniques  
565 such as pre-screening or penalization.

566 The problem definition of this method and the algorithmic structure used to  
567 solve it deliver further advantages such as the ability to cope with noisy inputs,  
568 missing data, multicollinearity and the capacity to deal with outliers within the  
569 response and the explanatory variables (robustness).

570 Some issues which must be taken into account concerning the proposed method  
571 are (i) its computational demand and (ii) its need for an external validation. As  
572 far as the first aspect goes, this can be considered indeed negligible compared to  
573 the time often required to collect the data it should analyse and can be greatly  
574 reduced according to the needs and requirements of the user. Concerning the  
575 second aspect, external validation is a crucial point which is often overlooked  
576 and is required for any model selection procedure. In this sense, the proposed  
577 method does not differ from any other existing approach in terms of additional  
578 requirements.

579 Having proposed a method with considerable advantages for gene selection  
580 using statistical ideas in model selection and machine learning, future research  
581 aims at studying the statistical properties of this approach to understand its  
582 asymptotic behaviour and develop the related inference tools.

## 583 **Acknowledgements**

584 We are very thankful to John Ramey (<http://ramhiser.com/>) for having pro-  
585 cessed the breast cancer and leukemia data set in Github and for having kindly  
586 answered our requests.

587 We thank Maria-Pia Victoria-Feser (Research Center for Statistics, University  
588 of Geneva, Switzerland) for her valuable comments and inputs as well as her  
589 institutional support.

## 590 **Funding and Conflict of interest**

591 No conflict of interest can be declared.

## 592 References

- 593 Andres, S. A. and Wittliff, J. L. (2012). Co-expression of genes with estrogen receptor- $\alpha$  and progesterone receptor in human breast carcinoma tissue.  
594  
595 *Hormone molecular biology and clinical investigation*, **12**(1), 377–390.
- 596 Azzalini, A. (1996). *Statistical inference based on the likelihood*, volume 68. CRC  
597 Press.
- 598 Bernusso, V. A., Machado-Neto, J. A., Pericole, F. V., Vieira, K. P., Duarte,  
599 A. S., Traina, F., Hansen, M. D., Saad, S. T. O., and Barcellos, K. S. (2015).  
600 Imatinib restores vasp activity and its interaction with zyxin in bcr–abl leukemic  
601 cells. *Biochimica et Biophysica Acta (BBA)-Molecular Cell Research*, **1853**(2),  
602 388–395.
- 603 Bohrer, L. R., Chuntova, P., Bade, L. K., Beadnell, T. C., Leon, R. P., Brady,  
604 N. J., Ryu, Y., Goldberg, J. E., Schmechel, S. C., Koopmeiners, J. S., *et al.*  
605 (2014). Activation of the fgfr–stat3 pathway in breast cancer cells induces  
606 a hyaluronan-rich microenvironment that licenses tumor formation. *Cancer*  
607 *research*, **74**(1), 374–386.
- 608 Cantoni, E., Field, C., Mills Flemming, J., and Ronchetti, E. (2007). Longitudinal  
609 variable selection by cross-validation in the case of many covariates. *Statistics*  
610 *in medicine*, **26**(4), 919–930.
- 611 Chapelle, O., Vapnik, V., Bousquet, O., and Mukherjee, S. (2002). Choosing  
612 multiple parameters for support vector machines. *Machine learning*, **46**(1-3),  
613 131–159.
- 614 Chin, K., DeVries, S., Fridlyand, J., Spellman, P. T., Roydasgupta, R., Kuo,  
615 W.-L., Lapuk, A., Neve, R. M., Qian, Z., Ryder, T., *et al.* (2006). Genomic and  
616 transcriptional aberrations linked to breast cancer pathophysiologies. *Cancer*  
617 *cell*, **10**(6), 529–541.
- 618 Chou, J., Provot, S., and Werb, Z. (2010). Gata3 in development and cancer  
619 differentiation: cells gata have it! *Journal of cellular physiology*, **222**(1), 42–49.
- 620 Christer, H., Peter, K., Margaret, L. A., Stephen, H., and Kathryn, M. T. (2013).  
621 A mechanism for epithelial-mesenchymal transition and anoikis resistance in  
622 breast cancer triggered by zinc channel zip6 and stat3 (signal transducer and  
623 activator of transcription 3). *Biochemical Journal*, **455**(2), 229–237.

- 624 Chung, S. S., Giehl, N., Wu, Y., and Vadgama, J. V. (2014). Stat3 activation in  
625 her2-overexpressing breast cancer promotes epithelial-mesenchymal transition  
626 and cancer stem cell traits. *International journal of oncology*, **44**(2), 403–411.
- 627 Díaz-Uriarte, R. and De Andres, S. A. (2006). Gene Selection and Classification  
628 of Microarray Data using Random Forest. *BMC Bioinformatics*, **7**(1), 3.
- 629 Dudoit, S., Fridlyand, J., and Speed, T. P. (2002). Comparison of Discrimination  
630 Methods for the Classification of Tumors using Gene Expression Data. *Journal  
631 of the American statistical association*, **97**(457), 77–87.
- 632 Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization paths for  
633 generalized linear models via coordinate descent. *Journal of statistical software*,  
634 **33**(1), 1.
- 635 Gabrielsson, B. G., Johansson, J. M., Lönn, M., Jernås, M., Olbers, T., Peltonen,  
636 M., Larsson, I., Lönn, L., Sjöström, L., Carlsson, B., *et al.* (2003). High  
637 expression of complement components in omental adipose tissue in obese men.  
638 *Obesity research*, **11**(6), 699–708.
- 639 George, E. and McCulloch, R. (1993). Variable selection via gibbs sampling.  
640 *Journal of the American Statistical Association*, **88**(423), 881–889.
- 641 George, E. I. and McCulloch, R. E. (1997). Approaches for bayesian variable  
642 selection. *Statistica sinica*, **7**(2), 339–373.
- 643 Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov,  
644 J. P., Coller, H., Loh, M. L., Downing, J. R., and Caligiuri, M. A. (1999).  
645 Molecular Classification of Cancer: Class Discovery and Class Prediction by  
646 Gene Expression Monitoring. *Science*, **286**(5439), 531–537.
- 647 Hall, A. R. *et al.* (2005). *Generalized method of moments*. Oxford University  
648 Press Oxford.
- 649 Juang, B., Hou, W., and Lee, C. (1997). Minimum classification error rate methods  
650 for speech recognition. *Speech and Audio Processing, IEEE Transactions on*,  
651 **5**(3), 257–265.
- 652 Koenker, R. (2005). *Quantile regression*. Number 38. Cambridge university press.
- 653 Kouros-Mehr, H., Slorach, E. M., Sternlicht, M. D., and Werb, Z. (2006). Gata-3  
654 maintains the differentiation of the luminal cell fate in the mammary gland.  
655 *Cell*, **127**(5), 1041–1055.

- 656 Kristensen, V. N., Vaske, C. J., Ursini-Siegel, J., Van Loo, P., Nordgard, S. H.,  
657 Sachidanandam, R., Sørli, T., Wärnberg, F., Haakensen, V. D., Helland, Å.,  
658 *et al.* (2012). Integrated molecular profiles of invasive breast tumors and ductal  
659 carcinoma in situ (dcis) reveal differential vascular and interleukin signaling.  
660 *Proceedings of the National Academy of Sciences*, **109**(8), 2802–2807.
- 661 Leng, C., Lin, Y., and Wahba, G. (2006). A note on the lasso and related  
662 procedures in model selection. *Statistica Sinica*, pages 1273–1284.
- 663 Masahiro, M., Mizuho, S., Yunfeng, Y., Masayoshi, I., Ryosuke, F., Takuya,  
664 O., Norihiro, I.-K., Tatsuo, T., and Naoki, W. (2012). Abi-1-bridged tyrosine  
665 phosphorylation of vasp by abelson kinase impairs association of vasp to focal  
666 adhesions and regulates leukaemic cell adhesion. *Biochemical Journal*, **441**(3),  
667 889–899.
- 668 Park, M. and Hastie, T. (2008). Penalized logistic regression for detecting gene  
669 interactions. *Biostatistics*, **9**(1), 30–50.
- 670 Ratajczak, M. (2014). A novel view of the adult bone marrow stem cell hierarchy  
671 and stem cell trafficking. *Leukemia*.
- 672 Sakamoto, K. M., Grant, S., Saleiro, D., Crispino, J. D., Hijjiya, N., Giles, F.,  
673 Plataniias, L., and Eklund, E. A. (2015). Targeting novel signaling pathways for  
674 resistant acute myeloid leukemia. *Molecular genetics and metabolism*, **114**(3),  
675 397–402.
- 676 Taniguchi, K. and Karin, M. (2014). Il-6 and related cytokines as the critical  
677 lynchpins between inflammation and cancer. In *Seminars in immunology*,  
678 volume 26, pages 54–74. Elsevier.
- 679 Tibshirani, R., Hastie, T., Narasimhan, B., and Chu, G. (2002). Diagnosis of  
680 Multiple Cancer Types by Shrunk Centroids of Gene Expression. *Proceedings*  
681 *of the National Academy of Sciences*, **99**(10), 6567–6572.
- 682 Vaske, C. J., Benz, S. C., Sanborn, J. Z., Earl, D., Szeto, C., Zhu, J., Haussler, D.,  
683 and Stuart, J. M. (2010). Inference of patient-specific pathway activities from  
684 multi-dimensional cancer genomics data using PARADIGM. *Bioinformatics*,  
685 **26**(12), i237–i245.
- 686 White, R. T., Damm, D., Hancock, N., Rosen, B., Lowell, B., Usher, P., Flier, J.,  
687 and Spiegelman, B. (1992). Human adiponin is identical to complement factor d  
688 and is expressed at high levels in adipose tissue. *Journal of Biological Chemistry*,  
689 **267**(13), 9210–9213.

- 690 Xu, Y., Ding, Y., Li, X., and Wu, X. (2015). Cystatin c is a disease-associated  
691 protein subject to multiple regulation. *Immunology and cell biology*.
- 692 Yang, A.-J. and Song, X.-Y. (2010). Bayesian variable selection for disease  
693 classification using gene expression data. *Bioinformatics*, **26**(2), 215–222.
- 694 Zhang, T. (2011). Adaptive forward-backward greedy algorithm for learning sparse  
695 representations. *Information Theory, IEEE Transactions on*, **57**(7), 4689–4708.
- 696 Zhu, J. and Hastie, T. (2004). Classification of Gene Microarrays by Penalized  
697 Logistic Regression. *Biostatistics*, **5**(3), 427–443.
- 698 Zou, H. and Hastie, T. (2005). Regularization and Variable Selection via the  
699 Elastic Net. *Journal of the Royal Statistical Society: Series B*, **67**(2), 301–320.

## 700 **A Adapting the algorithm to $p$**

701 In this subsection we provide two variants of the algorithm proposed in Section 3  
702 in order to adapt it to situations where  $p$  is either small or large.

### 703 **A.1 Adapting the algorithm to very large $p$**

704 In situations where  $p$  is extremely large and the initial step of the algorithm is not  
705 computationally feasible, this step can, for example, be replaced by the following  
706 modified initial step:

707 *A'. Large  $p$  Modified Initial Step:* We start by augmenting our initial variable  
708 set  $\mathcal{M}_0$  with  $d = 1$  variable in order to construct the set  $\mathcal{I}_1^*$ .

- 709 1. Augment  $\mathcal{M}_0$  with  $d = 1$  variable selected uniformly at random in  $\mathcal{J}_f$ .
- 710 2. Construct  $B$  models of dimension 1 by repeating Step A'.1  $B$  times.
- 711 3. From Steps A'.1 and A'.2, construct the set  $\mathcal{I}_1^*$  using (3). Go to Step  
712 B and let  $d = 2$ .

### 713 **A.2 Adapting the algorithm to small $p$**

714 On the other hand, when  $p$  is of reasonable size it may be possible to compute  
715 and evaluate all the  $\binom{p}{d'}$  models of dimension  $2 \leq d' \leq d_{\max}$ . In such cases, it may  
716 be feasible to also modify the initial step of the proposed algorithm to a different  
717 modified initial step. A possible modification is the following:

718 *A''. Small  $p$  Modified Initial Step:* We start by augmenting our initial variable set  
719  $\mathcal{M}_0$  with  $d$  ( $1 \leq d \leq d'$ ) variables in order to construct the sets  $\mathcal{I}_1^*, \dots, \mathcal{I}_{d'}^*$ .

- 720 1. We augment our initial variable set  $\mathcal{M}_0$  with 1 variable in order to  
721 construct the set  $\mathcal{I}_1^*$ .
  - 722 (i) Construct the  $p$  possible models obtained by augmenting  $\mathcal{M}_0$  with  
723 each of the  $p$  available variables.
  - 724 (ii) Compute  $\widehat{D}(\cdot, \cdot)$  for every model obtained in Step (i).
  - 725 (iii) From Steps (i) and (ii), construct the set  $\mathcal{I}_1^*$  using (3). Go to Step  
726 A''.2 and let  $d = 2$ .

- 727           2. We augment our initial model  $\mathcal{M}_0$  set by  $d$  variables in order to  
728           construct the set  $\mathcal{I}_d^*$ .
- 729           (i) Construct the  $\binom{p}{d}$  possible models and augment  $\mathcal{M}_0$  with all vari-  
730           ables of these constructed models.
- 731           (ii) Compute  $\widehat{D}$  for every model obtained in Step (i).
- 732           (iii) From Steps (i) and (ii), construct the set  $\mathcal{I}_d^*$  using (3) and let  
733            $d = d + 1$ . Go to Step A''.2 (if  $d < d'$ ) or Step B.1 (if  $d \geq d'$ ), with  
734           model dimension starting value  $d$ .

## 735   B   Complementary results on Acute Leukemia

736   Table 3 reports the main biomarker hubs and related biomarker networks for the  
737   *leukemia* data set analysed in Section 4.1.

738   Table 4 reports the performances of our implementation of the competing  
739   methods as described in Section 5. Unlike reported in Table 1, here the proposed  
740   method uses the classical tenfold-CV for  $D(\cdot, \cdot)$  ( $K = 1$ ). The other hyper-  
741   parameters are kept the same (i.e.  $\alpha = 0.01$ ,  $B = 20'000$  and  $\pi = 0.5$ ).

742

## 743   C   Breast Cancer

744   The second data-set we analyzed is the *breast cancer* data presented in Chin  
745   *et al.* (2006). The main goal behind analyzing this data is to identify the estrogen  
746   receptor expression on tumor cells which is a crucial step for the correct manage-  
747   ment of breast cancer. Similarly to Table 4 in Appendix B, Table 5 reports the  
748   performances of our implementation of the competing methods and the proposed  
749   approach on the *breast cancer* data. For the sake of this comparison, the data-set  
750   was randomly split into training (60) and test (58) sets. The hyper-parameters of  
751   the proposed method are  $\alpha = 0.01$ ,  $B = 30'000$ ,  $\pi = 0.5$  and  $D(\cdot, \cdot)$  is the classical  
752   tenfold-CV ( $K = 1$ ).

753

Affy ID	Gene ID	Gene Function	Biological Process	
<b>NETWORK 1</b>				
Position 1	M27891_at	ENSG00000101439 Cystatin C	AA	
Position 2	D80006_at	ENSG00000114978 MOB kinase activator 1A	AA	
	M20778_s_at	ENSG00000163359 Collagen, type VI, alpha 3	AA	
	U57316_at	ENSG00000108773 K(lysine) acetyltransferase 2A	TF	
	U90549_at	ENSG00000182952 High mobility group nucleosomal binding domain 4	TF	
	X66899_at	ENSG00000182944 Ewing Sarcoma region 1; RNA binding protein	TF	
	M74088_s_at	ENSG00000134982 Adenomatous polyposis coli, DP2, DP3, PPP1R46	TF	
	U51166_at	ENSG00000139372 thymine-DNA glycosylase	TF	
	Z69881_at	ENSG00000074370 ATPase, Ca++ transporting, ubiquitous	IPT	
	U49248_at	ENSG00000023839 ATP-binding cassette, sub-family C (CFTR/MRP), member 2	IPT	
	X89109_s_at	ENSG00000102879 Coronin, actin binding protein, 1A	IPT	
	HG2815-HT2931_at	ENSG00000092841 Myosin, Light Chain, Alkali, Smooth Muscle (Gb:U02629)	ACC	
	M94345_at	ENSG00000042493 Capping protein (actin filament), gelsolin-like	ACC	
	L33075_at	ENSG00000140575 IQ motif containing GTPase activating protein 1	ACC	
	L07633_at	ENSG00000092010 Proteasome (prosome, macropain) activator subunit 1 (PA28 alpha)	APC	
	J03589_at	ENSG00000102178 Ubiquitin-like 4A	APC	
	D83920_at	ENSG00000085265 FCN1, Ficolin-1	IR	
	X03934_at	ENSG00000167286 CD3d molecule, delta (CD3-TCR complex)	IR	
<b>NETWORK 2</b>				
Position 1	X95735_at	ENSG00000159840 Zyxin	ACC	
Position 2	X04526_at	ENSG00000185838 Guanine nucleotide binding protein (G protein), beta polypeptide 1	ST	
	D78577_s_at	ENSG00000128245 Tyrosine 3-monooxygenase/tryptophan 5-monooxygenase activation protein, eta	ST	
	U32645_at	ENSG00000102034 E74-like factor 4 (ets domain transcription factor)	TF	
	U93867_at	ENSG00000186141 Polymerase (RNA) III (DNA directed) polypeptide C (62kD)	TF	
	U29175_at	ENSG00000127616 SWI/SNF related, matrix associated, actin dependent regulator of chromatin, subfamily A, member 4	TF	
	Y00291_at	ENSG00000077092 Retinoic acid receptor, beta	TF	
	D17532_at	ENSG00000110367 DEAD (Asp-Glu-Ala-Asp) Box Helicase 6	TF	
	HG3521-HT3715_at	ENSG00000127314 Ras-Related Protein Rap1b	TF	
	M83233_at	ENSG00000140262 Transcription factor 12	TF	
	U94855_at	ENSG00000175390 Eukaryotic translation initiation factor 3, subunit F	TF	
	L07758_at	ENSG00000136045 PWP1 homolog	TF	
	D63506_at	ENSG00000116266 Syntaxin binding protein 3	IR	
	M33680_at	ENSG00000110651 CD81 molecule	IR	
	HG1612-HT1612_at	ENSG00000175130 Macmarcks	CG	
	M92287_at	ENSG00000112576 Cyclin D3	CG	
	M60483_rna1_s_at	ENSG00000113575 Protein Phosphatase 2 (formerly 2A), catalytic subunit, alpha isoform	CG	
	U84388_at	ENSG00000169372 CASP2 and RIPK1 domain containing adaptor with death domain	AA	
	S80437_s_at	ENSG00000169710 Fatty acid synthase		
	<b>NETWORK 3</b>			
	Position 1	M84526_at	ENSG00000197766 Complement factor D (adipsin)	IR
Position 2	M28130_rna1_s_at	ENSG00000169429 Interleukine-8	IR	
	Z32765_at	ENSG00000135218 CD36 - Thrombospondin receptor	IR	

TABLE 3: *Biomarker network organisation - leukemia data set - Lymphoblastic / Myeloblastic leukemia. TF = Transcription/translation factor activity, DNA repair and catabolism - AA = apoptotic activity - IR = immunity, inflammatory response (blood coagulation, antigen presentation and complement activation) - IPT = intracellular protein trafficking, transmembrane transport - ACC = actin activity, cytoskeleton organisation - APC = protein catabolism - ST = intracellular signal transduction - CG = cell growth, proliferation and division. Source: [www.ensembl.org](http://www.ensembl.org); [www.uniprot.org](http://www.uniprot.org)*

Method	Tenfold CV error	Test error	Number of genes
Support vector machine (with recursive feature elimination)	0/38	5/34	2/7129
Penalised logistic regression (with forward selection followed by backward deletion)	0/38*	4/34	3/7129
Nearest shrunken centroids	3/38	1/34	372/7129
Elastic net	3/38	2/34	74/7129
Panning Algorithm (131)			
Model a	0/38	1/34	2/7129
Model b	0/38	2/34	2/7129
Model c	0/38	2/34	2/7129
[...]			2/7129

TABLE 4: *Performances of our implementation of the competing methods on the leukemia data-set. For the Penalised logistic regression(\*), the in-sample error is reported instead of the tenfold-CV error. For the Panning Algorithm, models “a” to “c” are three examples out of the 131 models. All the 131 models have a tenfold-CV error of 0. The best test error is 1 and the worst is 20.*

Method	Tenfold CV error	Test error	Number of genes
Support vector machine (with recursive feature elimination)	0/60	10/58	3/22215
Penalised logistic regression (with forward selection followed by backward deletion)	0/60*	12/58	15/22215
Nearest shrunken centroids	2/60	11/58	5/22215
Elastic net	3/60	11/58	196/22215
Panning Algorithm (241)			
Model a	0/60	9/58	3/22215
Model b	0/60	9/58	3/22215
Model c	0/60	10/58	3/22215
[...]			3/22215

TABLE 5: *Performances of our implementation of the methods on the breast cancer data-set. For the Penalised logistic regression (\*), the in-sample error is reported instead of the tenfold-CV error. For the proposed method, models “a” to “c” are three examples out of 241 models. All the 241 models have a tenfold-CV error of 0. The best test error is 9 and the worst is 28.*

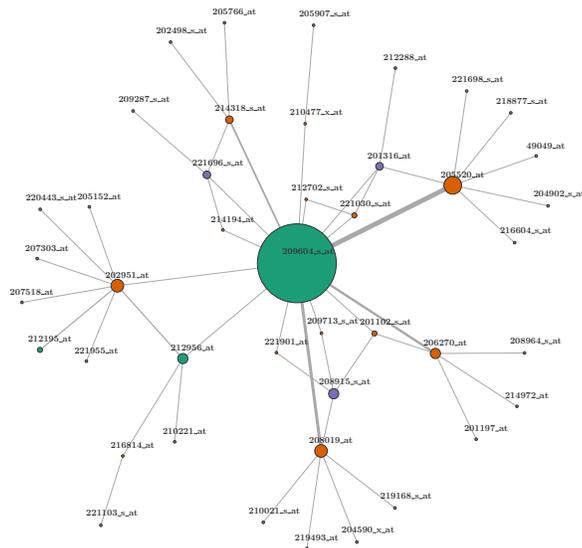


FIGURE 3: *Network representation of biomarkers selected from breast cancer data-set. Colors represent the position of covariates within the model: green for first position (hub), orange for second and purple for third. The width of the connecting lines is proportional to the frequency with which two biomarkers appear in the same model. The size of the circles is proportional to the frequency with which a biomarker is present within the selected set of models. (Note: biomarker “209602\_s\_at” is merged with biomarker “209604\_s\_at”).*

754 Figure 3 shows the paradigmatic network identified by our method for the  
 755 *breast cancer* data for which the selected model dimension is three (i.e. only three  
 756 biomarkers are needed in a model to well classify the breast cancer). We used the  
 757 hyper-parameters  $\alpha = 0.01$ ,  $B = 22'215$ ,  $\pi = 0.05$  and for  $D(\cdot, \cdot)$  the tenfold-CV  
 758 repeated  $K = 10$  times was used. Table 6 provides the details of the networks  
 759 based on the three main hubs and is to be interpreted as described in Section 4.1.

760 This figure is a clear example of the advantages of the proposed method since,  
 761 it not only selects a set of low-dimensional models with a high predictive power,  
 762 but also provides the basis for a more general biological interpretation which takes  
 763 into account interactions between different biomarkers as opposed to one single

764 model. The three main hubs identified through the proposed algorithm are:

- 765 1. GATA binding protein 3 (GATA3): a transcription factor regulating the  
766 differentiation of breast luminal epithelial cells;
- 767 2. IL6 Signal Transducer (IL6 ST): a pro-inflammatory cytokine signal trans-  
768 ducer;
- 769 3. TBC1 domain family, member 9 (TBC1D9): a GTPase-activating protein  
770 for Rab family protein involved in the expression of the ER in breast tumors.

771 GATA3 is known to regulate the differentiation of epithelial cells in mammary  
772 glands (see [Kouros-Mehr \*et al.\*, 2006](#)) and is required for luminal epithelial  
773 cell differentiation. Its expression is progressively lost during luminal breast  
774 cancer progression as cancer cells acquire a stem cell-like phenotype (see [Chou  
775 \*et al.\*, 2010](#)). IL6 ST has been linked to breast cancer epithelial-mesenchymal  
776 transition and cancer stem cell traits (see [Chung \*et al.\*, 2014](#)), cancer-promoting  
777 microenvironment (see [Bohrer \*et al.\*, 2014](#)) and resistance (see [Christer \*et al.\*,  
778 2013](#)). Moreover, this result supports the assertion by [Taniguchi and Karin \(2014\)](#)  
779 that IL6 ST and related cytokines are the critical lynchpins between inflammation  
780 and cancer. Finally, concerning the third biomarker, a recent publication by  
781 [Andres and Wittliff \(2012\)](#) has shown that the expression of the ER on the surface  
782 of breast tumor cells is highly correlated with the coordinate expression of different  
783 genes among which we can find TBC1D9 and GATA3. These two genes are not  
784 only considered as relevant genes according to the proposed method but as actual  
785 hubs of the “best” models which define the structure of the identified network.  
786 Instead of selecting a single model with many biomarkers whose interactions may  
787 be difficult to interpret, the proposed method selects a set of models with few  
788 biomarkers that allow them to be individually easy to interpret without losing  
789 the possibility of interpreting them within the larger network. This is what this  
790 paper intends with the expression “paradigmatic network” since by taking this  
791 approach it is possible to identify a set of biomarker families within which each  
792 biomarker is interchangeable with the others.

	Affy ID	Gene ID	Gene Function	Biological Process
<b>NETWORK 1</b>				
Position 1	209604_s.at	ENSG00000107485	GATA binding protein 3	TF
Position 2	205520_at	ENSG00000115808	Striatin, calmodulin binding protein	ER
Position 3	204902_s.at	ENSG00000168397	Autophagy related 4B, cysteine peptidase (APG4B, AUTL1, DKFZp586D1822, KIAA0943)	APC
	221698_s.at	ENSG00000172243	C-type lectin domain family 7, member A	IR
	49049_at	ENSG00000178498	Deltex 3, E3 ubiquitin ligase	APC
	209602_s.at	ENSG00000107485	GATA binding protein 3	TF
	216604_s.at	ENSG00000003989	Solute carrier family 7 (cationic amino acid transporter, y+ system), member 2	IPT
	218877_s.at	ENSG00000066651	TRNA methyltransferase 11 homolog	TF
	201316_at	ENSG00000106588	Proteasome (prosome, macropain) subunit, alpha type, 2	APC
Position 2	208019_at	ENSG00000147117	Zinc finger protein 157	TF
Position 3	219168_s.at	ENSG00000186654	PRR5 (Proline rich 5 (renal))	CG
	219493_at	ENSG00000171241	SHC SH2-domain binding protein 1	CG
	204590_x.at	ENSG00000139719	Vacuolar protein sorting 33 homolog A	APC
	210021_s.at	ENSG00000152669	Cyclin O	CG
	208915_s.at	ENSG00000103365	Golgi-associated, gamma adaptin ear containing, ARF binding protein 2	IPT
Position 2	214318_s.at	ENSG00000073910	Furry homolog	ACC
Position 3	205766_at	ENSG00000173991	Titin-cap (Telethonin)	ACC
	221696_s.at	ENSG00000060140	Serine/threonine/tyrosine kinase 1	CG
	202498_s.at	ENSG00000059804	Solute carrier family 2 (facilitated glucose transporter), member 3	STM
Position 2	201102_s.at	ENSG00000141959	Phosphofructokinase, liver	STM
Position 3	208915_s.at	ENSG00000103365	Golgi-associated, gamma adaptin ear containing, ARF binding protein 2	IPT
Position 2	201316_at	ENSG00000106588	Proteasome (prosome, macropain) subunit, alpha type, 2	APC
Position 3	212288_at	ENSG00000187239	Formin binding protein 1	ACC
Position 2	209713_s.at	ENSG00000116704	Solute carrier family 35 (UDP-GlcA/UDP-GalNAc transporter), member D1	STM
Position 3	208915_s.at	ENSG00000103365	Golgi-associated, gamma adaptin ear containing, ARF binding protein 2	IPT
Position 2	212702_s.at	ENSG00000185963	Bicaudal D homolog 2	ACC
Position 3	221030_s.at	ENSG00000138639	Rho GTPase activating protein 24	ACC
Position 2	212956_at	ENSG00000109436	TBC1 domain family, member 9 (with GRAM domain)	IPT
Position 3	210221_at	ENSG00000080644	Cholinergic receptor, nicotinic, alpha 3 (neuronal)	ITT
Position 2	214194_at	ENSG00000083520	DIS3 mitotic control homolog (Ribosomal RNA-processing protein 44)	TF
Position 3	221696_s.at	ENSG00000060140	Serine/threonine/tyrosine kinase 1	CG
Position 2	216814_at	ENSG00000232267	ACTR3 pseudogene 2	PUP
Position 3	221103_s.at	ENSG00000206530	Cilia and flagella associated protein 44	ACC
Position 2	221030_s.at	ENSG00000138639	Rho GTPase activating protein 24	ACC
Position 3	201316_at	ENSG00000106588	Proteasome (prosome, macropain) subunit, alpha type, 2	APC
Position 2	221696_s.at	ENSG00000060140	Serine/threonine/tyrosine kinase 1	CG
Position 3	209287_s.at	ENSG00000070831	Cell division control protein 42 homolog	ACC
Position 2	221901_at	ENSG00000138944	KIAA1644	PUP
Position 3	208915_s.at	ENSG00000103365	Golgi-associated, gamma adaptin ear containing, ARF binding protein 2	IPT
Position 1	209602_s.at	ENSG00000107485	GATA3	TF
Position 2	202951_at	ENSG00000112079	Serine/threonine kinase 38	CG
Position 3	220443_s.at	ENSG00000116035	VAX2 (ventral anterior homeobox 2)	TF
	221955_at	ENSG00000088256	Guanine nucleotide binding protein (G protein), alpha 11 (Gq class)	ITT
	207303_at	ENSG00000154678	Phosphodiesterase 1C, calmodulin-dependent 70kDa	ST

	205152_at	ENSG00000157103	Solute carrier family 6, member 1	ST
	207518_at	ENSG00000153933	Diacylglycerol kinase, epsilon 64kDa	ST
Position 2	206270_at	ENSG00000126583	Protein kinase C, gamma	ST
Position 3	208964_s.at	ENSG00000149485	Fatty acid desaturase 1	FAM
	201197_at	ENSG00000123505	Adenosylmethionine decarboxylase 1	CG
	201102_s.at	ENSG00000141959	ATP-dependent 6-phosphofructokinase, liver type	STM
	214972_at	ENSG00000198408	Protein O-GlcNAcase (Meningioma expressed antigen 5 (hyaluronidase))	ST
Position 2	210477_x.at	ENSG00000107643	Mitogen-activated protein kinase 8	CG
Position 3	205907_s.at	ENSG00000127083	Osteomodulin	STM
<b>NETWORK 2</b>				
Position 1	212195_at	ENSG00000134352	IL6 Signal Transducer	ICT
Position 2	202951_at	ENSG00000112079	Serine/threonine kinase 38	CG
Position 3	221955_at	ENSG00000088256	Guanine nucleotide binding protein (G protein), alpha 11 (Gq class)	ITT
	207303_at	ENSG00000154678	Phosphodiesterase 1C, calmodulin-dependent 70kDa	ICT
<b>NETWORK 3</b>				
Position 1	212956_at	ENSG00000109436	TBC1 domain family, member 9 (with GRAM domain)	IPT
Position 2	202951_at	ENSG00000112079	Serine/threonine kinase 38	CG
Position 3	205152_at	ENSG00000157103	Solute carrier family 6, member 1	ST
	207518_at	ENSG00000153933	Diacylglycerol kinase, epsilon 64kDa	ST
Position 2	216814_at	ENSG00000232267	ACTR3 pseudogene 2	PUP
Position 3	221103_s.at	ENSG00000206530	Cilia and flagella associated protein 44	ACC

TABLE 6: *Biomarker network organisation - breast cancer data set - Estrogen Receptor - Breast Cancer.*

*TF = Transcription/translation factor activity, DNA/RNA repair and catabolism - ER = estrogen receptor activity - APC = autophagy - protein catabolism - IR = immunity, inflammatory response (blood coagulation, antigen presentation and complement activation) - CC = cell/cell communication - ST = intracellular signal transduction, protein glycosylation - CG = cell growth and division - IPT = intracellular protein trafficking, transmembrane amino-acid transporter - ACC = actin activity, cytoskeleton organisation, cell projection - STM = sugar transport and metabolism - ITT = ion transmembrane transport, transmembrane signaling systems - PUP = pseudogene, uncharacterized protein - FAM = fatty acid metabolism. Source: [www.uniprot.org](http://www.uniprot.org); [www.ncbi.nlm.nih.gov/gene](http://www.ncbi.nlm.nih.gov/gene)*