

A Sparse-Modeling based approach for Class-Specific feature selection

Davide Nardone^{Corresp., 1}, **Angelo Ciaramella**¹, **Antonino Staiano**¹

¹ Dipartimento di Scienze e Tecnologie, Università degli Studi di Napoli "Parthenope", Naples, Italy

Corresponding Author: Davide Nardone
Email address: davide.nardone@live.it

In this work, we propose a novel Feature Selection framework, called Sparse-Modeling Based Approach for Class Specific Feature Selection (SMBA-CSFS), that simultaneously exploits the idea of Sparse Modeling and Class-Specific Feature Selection. Feature selection plays a key role in several fields (e.g., computational biology), making it possible to treat models with fewer variables which, in turn, are easier to explain, by providing valuable insights on the importance of their role, and possible speed-up of the experimental validation. Unfortunately, also corroborated by the no free lunch theorems, none of the approaches in literature is the most apt to detect the optimal feature subset for building a final model, thus it still represents a challenge. The proposed feature selection procedure conceives a two steps approach: (a) a sparse modeling-based learning technique is first used to find the best subset of features, for each class of a training set; (b) the discovered feature subsets are then fed to a class-specific feature selection scheme, in order to assess the effectiveness of the selected features in classification tasks. To this end, an ensemble of classifiers is built, where each classifier is trained on its own feature subset discovered in the previous phase, and a proper decision rule is adopted to compute the ensemble responses. In order to evaluate the performance of the proposed method, extensive experiments have been performed on publicly available datasets, in particular belonging to the computational biology field where feature selection is indispensable: the acute lymphoblastic leukemia and acute myeloid leukemia, the human carcinomas, the human lung carcinomas, the diffuse large B-cell lymphoma, and the malignant glioma. SMBA-CSFS is able to identify/retrieve the most representative features that maximize the classification accuracy. With top 20 and 80 features, SMBA-CSFS exhibits a promising performance when compared to its competitors from literature, on all considered datasets, especially those with a higher number of features. Experiments show that the proposed approach might outperform the state-of-the-art methods when the number of features is high. For this reason, the introduced approach proposes itself for selection and classification of data with a large number of features and classes.

A Sparse-Modeling Based Approach for Class-Specific Feature Selection

Davide Nardone¹, Angelo Ciaramella¹, and Antonino Staiano¹

¹Dipartimento di Scienze e Tecnologie, Università degli Studi di Napoli "Parthenope",
Centro Direzionale, Isola C4, 80143, Naples, Italy

Corresponding author:

Davide Nardone¹

Email address: davide.nardone@studenti.uniparthenope.it

ABSTRACT

In this work, we propose a novel Feature Selection framework, called Sparse-Modeling Based Approach for Class Specific Feature Selection (SMBA-CSFS), that simultaneously exploits the idea of *Sparse Modeling* and *Class-Specific Feature Selection*. Feature selection plays a key role in several fields (e.g., computational biology), making it possible to treat models with fewer variables which, in turn, are easier to explain, by providing valuable insights on the importance of their role, and possible speed-up of the experimental validation. Unfortunately, also corroborated by the no free lunch theorems, none of the approaches in literature is the most apt to detect the optimal feature subset for building a final model, thus it still represents a challenge. The proposed feature selection procedure conceives a two steps approach: (a) a sparse modeling-based learning technique is first used to find the best subset of features, for each class of a training set; (b) the discovered feature subsets are then fed to a class-specific feature selection scheme, in order to assess the effectiveness of the selected features in classification tasks. To this end, an ensemble of classifiers is built, where each classifier is trained on its own feature subset discovered in the previous phase, and a proper decision rule is adopted to compute the ensemble responses. In order to evaluate the performance of the proposed method, extensive experiments have been performed on publicly available data sets, in particular belonging to the computational biology field where feature selection is indispensable: the *acute lymphoblastic leukemia* and *acute myeloid leukemia*, the *human carcinomas*, the *human lung carcinomas*, the *diffuse large B-cell lymphoma*, and the *malignant glioma*. SMBA-CSFS is able to identify/retrieve the most representative features that maximize the classification accuracy. With top 20 and 80 features, SMBA-CSFS exhibits a promising performance when compared to its competitors from literature, on all considered data sets, especially those with a higher number of features. Experiments show that the proposed approach might outperform the state-of-the-art methods when the number of features is high. For this reason, the introduced approach proposes itself for selection and classification of data with a large number of features and classes.

INTRODUCTION

Data analysis is the process of evaluating data, that is often subject to high-dimensional feature spaces, i.e., where data are represented in, whatever the area of study, from biology to pattern recognition to computer vision. High dimensionality often translates into over-fitting, large computational costs and poor performance thus getting a learning task in trouble. Consequently, the high-dimensional feature spaces need to be lowered since its feature vectors are generally uninformative, redundant, correlated to each other and also noisy. In this paper, we focus on feature selection, which is undertaken to identify discriminative features by eliminating the ones with little or no predictive information, based on certain criteria, in order to treat with data in low dimensional spaces.

Feature Selection (FS) is the process of selecting a subset of relevant features for use in model construction. FS plays a key role in computational biology, for instance, microarray data analysis involves a huge number of genes w.r.t. a small number of samples, and effectively identifying the most significant differentially expressed genes under different conditions is prominent (Xiong et al., 2001). The selected genes are very useful in clinical applications such as recognizing diseased profiles (Calcagno et al., 2010; Staiano et al., 2013; Di Taranto et al., 2015; Camastra et al., 2015), nonetheless, because of its high costs, the number of experiments that can be used for classification purposes is usually limited so that the small number of samples, compared to the large number of genes in an experiment, gives rise to the *Curse of Dimensionality* problem (Friedman et al., 2001), which challenges the classification as well as other data analysis tasks (Staiano et al., 2004; Ciaramella et al., 2008). Furthermore, microarray data are usually not immune from several issues, such as sensitivity, accuracy, specificity, reproducibility of results, and noisy data (Draghici et al., 2006). For these reasons, it is unsuitable to use microarray data as they are, but, after several corrections, select the relevant genes by FS approaches and, for instance, validate the results using Real-Time PCR (Xiong et al., 2001).

Taking a look at the literature, by *googling* the keyword “*feature selection*”, one gets lost in an ocean of techniques (the reader might refer to classical reviews in (Saeys et al., 2007), (Guyon and Elisseeff, 2003) and (Hoque et al., 2014) on the topic), often designed to tackle a specific data set. The reasons for the abundance of techniques are in the heterogeneity of the available scientific data sets and also by the limitations dictated by *no free lunch theorems* (Wolpert and Macready, 1997), determining the existence of no general-purpose technique which well suites to a plethora of different kind of data. A typical taxonomy organizes FS techniques (Jović et al., 2015) in three main categories, namely *filter*, *wrapper* and *embedded* methods, whose belonging algorithms select a single feature subset from a complete list of features. Another perspective instead, divides FS techniques in two classes, namely, Traditional Feature Selection (TFS) for all classes (that includes filter, wrapper and embedded methods mentioned so far), and Class-Specific Feature Selection (CSFS) (Fu and Wang, 2002). Usually, a TFS algorithm selects one subset of features for all classes although it might be not the best one for some class, thus leading to undesirable results. Differently, a CSFS policy permits to select a distinct subset for each class, and it can use any traditional *feature selector*, for choosing, given the set of classes of a classification problem, one distinct grouping of features for every class. Depending on the type of the feature selector, the overall process may slightly change. Nevertheless, it is worth pointing out that a CSFS scheme heavily depends on the use of a specific classifier, while its use should be independent of both the classifier of the classification step and the feature selector strategy. To this end, a General Framework CSFS has been proposed in (Pineda-Bautista et al., 2011) which allows using any traditional feature selector as well as any classifier, consisting of four stages (the reader may refer to Methods section later on).

In this paper, on the basis of the general framework for CSFS, we propose a novel strategy to FS, namely a Sparse-Modeling based approach for Class-Specific Feature Selection, consisting of a two-steps procedure. Firstly, a sparse modeling based learning technique is used to find the best subset of features for each class of the training set. In so doing, it is assumed that a class is represented by using a subset of features, called *representatives*, such that each sample in a specific class, can be described as a linear combination of them. Secondly, the discovered feature subsets are fed to a class-specific feature selection scheme in order to assess the effectiveness of the selected features in classification task. To this end an ensemble of classifiers is built by training a given classifier, one for each class, on its own feature subset, i.e., the one discovered in the previous step, and a proper decision rule is adopted to compute the ensemble responses. In this way, the dilemma of choosing specific TFS strategy and classifiers in the CSFS framework is effectively mitigated.

87 METHODS

88 The sparse-modeling based approach for class-specific feature selection, is based on the concepts of sparse
89 modeling and class-specific feature selection that need to be properly introduced.

90 Sparse Modeling fundamentals

An active developing field of statistical learning is focused around the notion of sparsity (Tibshirani, 1994; Ciaramella and Giunta, 2016). A Sparse Model (SM) is a model that can be much easier to estimate and interpret than a dense model. The sparsity assumption allows extracting meaningful features from large data sets. Aim of the first phase of the proposed approach is to use a sparse modeling for finding data representatives without data transformation and to be performed directly in the data space. In other words, we wish to find a ranking of the most representatives features that best reconstruct the data collection. Most approaches are based on a l_1 -norm regularization (e.g, LASSO (Tibshirani, 1994), Sparse Dictionary Learning (Elhamifar et al., 2012)). Formally, given a set of features in \mathbb{R}^m arranged as columns of a data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, the task is to find representative features given a fixed feature space belonging to a collections of data points (see (Mairal et al., 2008; Aharon et al., 2006; Engan et al., 1999; Jolliffe, 1986; Ramirez et al., 2010)). That task can conveniently be described in the *Dictionary Learning* (DL) framework, where the aim is to simultaneously learn a compact dictionary $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_k] \in \mathbb{R}^{m \times k}$ and coefficients $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_n] \in \mathbb{R}^{k \times n}$, with $k \ll n$, that can well represent collections of data points (Ciaramella et al., 2016). The best representation of the data is obtained by minimizing the following objective function

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\mathbf{c}_i\|_2^2 = \|\mathbf{X} - \mathbf{D}\mathbf{C}\|_F^2 \quad (1)$$

91 w.r.t. the dictionary \mathbf{D} and the coefficient matrix \mathbf{C} , subject to appropriate constraints.

However, the dictionary learned atoms almost never correspond to the original feature space (Aharon et al., 2006; Ramirez et al., 2010; Mairal et al., 2009). In order to find a subset of features that best represent the entire feature space, the optimization problem in 1 is reformulated forcing the dictionary \mathbf{D} to be the data matrix \mathbf{X} (Elhamifar et al., 2012):

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{X}\mathbf{c}_i\|_2^2 = \|\mathbf{X} - \mathbf{X}\mathbf{C}\|_F^2, \quad (2)$$

where F is the Frobenius norm. Equation 2 is minimized w.r.t the coefficient matrix $\mathbf{C} \triangleq [\mathbf{c}_1, \dots, \mathbf{c}_n] \in \mathbb{R}^{n \times n}$, subject to additional constraints. In other words, the *reconstruction error* of each feature component is minimized by linearly combining all components of the feature space. To choose $k \ll n$ representatives involved in the linear reconstruction of the each component in (2), the following constraint is added to the model

$$\|\mathbf{C}\|_{0,q} \leq k, \quad (3)$$

where the mixed ℓ_0/ℓ_q norm is defined as $\|\mathbf{C}\|_{0,q} \triangleq \sum_{i=1}^N I(\|\mathbf{c}^i\|_q > 0)$, \mathbf{c}^i denotes the i -th row of \mathbf{C} , and $I(\cdot)$ denotes the indicator function. In a nutshell, $\|\mathbf{C}\|_{0,q}$ counts the number of nonzero rows of \mathbf{C} . The indices of the nonzero rows of \mathbf{C} correspond to the indices of the columns of \mathbf{X} which are chosen as the representative features. Since the aim is to select $k \ll n$ representatives features that can reconstruct each feature of the \mathbf{X} matrix up to a fixed error, the optimization problem to solve is

$$\begin{aligned} & \underset{\mathbf{C}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{X}\mathbf{C}\|_F^2 \\ & \text{subject to} \quad \|\mathbf{C}\|_{0,q} \leq k, \mathbf{1}^T \mathbf{C} = \mathbf{1}^T \end{aligned} \quad (4)$$

where $\mathbf{1}^T \mathbf{C} = \mathbf{1}^T$ is the affine constraint for selecting representatives that are invariant w.r.t. a global translation of the data (as requested by dimensionality reduction methods). This is an NP-hard problem as it implies a combinatorial calculation over every subset of the k columns of \mathbf{X} . Therefore, relaxing ℓ_0 to ℓ_1 norm, the problem becomes

$$\begin{aligned} & \underset{\mathbf{C}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{X}\mathbf{C}\|_F^2 \\ & \text{subject to} \quad \|\mathbf{C}\|_{1,q} \leq \tau, \mathbf{1}^T \mathbf{C} = \mathbf{1}^T \end{aligned} \quad (5)$$

Procedure SMBA

Input: \mathbf{X} , $N \times M$ matrix where N is the number observations and M is the number of features

$\theta = \{\alpha, \delta, \rho, \eta\}$, parameters vector

Output: I , set of features selected

1 Variables initialization

3 **while** $\varepsilon > \delta$ and $t > \rho$ **do**

4 $\beta^{t+1} \leftarrow (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1}$

5 $\theta^{t+1} \leftarrow (S_{\lambda/\rho}(\beta^{t+1} + \mu^t/\rho))$

6 $\mu^{t+1} \leftarrow \mu^t + \rho(\beta^{t+1} - \theta^{t+1})$

7 $\varepsilon \leftarrow \text{compute_error}(\beta, \theta)$

8 **end**

9 $I \leftarrow \text{find_representatives}(\theta, \eta)$

where $\|\mathbf{C}\|_{1,q} \triangleq \sum_{i=1}^N \|\mathbf{c}^i\|_q$ is the sum of the ℓ_q norms of the rows of \mathbf{C} and $\tau > 0$ is an appropriate chosen parameter. The solution of the optimization (5) not only provides the representative features as the nonzero rows of the \mathbf{C} , but also provides information about the ranking of the selected features. More precisely, a representative that has higher ranking takes part in the reconstruction process more than the others, hence, its corresponding row in the optimal coefficient matrix \mathbf{C} has many nonzero elements with large values. Conversely, a representative with lower ranking takes part in the reconstruction process less than the others, hence, its corresponding row in \mathbf{C} has a few nonzero elements with smaller values. Thus, the k representative features x_{i1}, \dots, x_{ik} are ranked as $i_1 \geq i_2 \geq \dots \geq i_k$, whenever for the corresponding rows of \mathbf{C} one gets

$$\|\mathbf{c}^{i_1}\|_q \geq \|\mathbf{c}^{i_2}\|_q \geq \dots \geq \|\mathbf{c}^{i_k}\|_q, \quad (6)$$

From a practical point of view, the optimization problem (5) can be expressed by using the Lagrange multipliers

$$\underset{\mathbf{C}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{X} - \mathbf{XC}\|_F^2 + \lambda \|\mathbf{C}\|_{1,q} \quad \text{subject to} \quad \mathbf{1}^T \mathbf{C} = \mathbf{1}^T. \quad (7)$$

In practice, the algorithm is implemented using an Alternating Direction Method of Multipliers (ADMM) optimization framework (Boyd et al., 2011). In particular, the features of a given data set are obtained considering representatives of small pairwise coherence features as in a sparse dictionary learning method. It is worth observing the resemblance with the Least Absolute Shrinkage and Selection Operator (LASSO) (Tibshirani, 1994). LASSO consists of an approach to regression analysis that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretation ability of the statistical model it produces. Recall that the objective of LASSO, in its basic form, is to solve

$$\underset{\beta}{\text{minimize}} \quad \frac{1}{N} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

$$\text{subject to} \quad \|\beta\|_1 \leq t, \quad (8)$$

92 where $\mathbf{y} = [y_1, \dots, y_N]$ is the N -dimensional vector of outcomes, \mathbf{X} the covariate matrix, t is a free
 93 parameter that determines the amount of regularization and β is the sparse vector to estimate.

94 From Equation 8, one can observe that a sparse matrix can be estimated as in equation 7 by considering
 95 \mathbf{X} itself as outcome and adding the affine constraint. In the following, the LASSO will be used for
 96 classification tasks, adopting a sigmoid function, as it will be described in the experimental setup.

97 **A Sparse-Modeling Based Approach for Class-Specific Feature Selection**

98 A General Framework for Class-Specific Feature Selection (GF-CSFS) is described in (Pineda-Bautista
 99 et al., 2011). The proposed Sparse-Modeling Based Approach for Class-Specific Feature Selection
 100 (SMBA-CSFS) tries to best represent each class-sample set of an input data set by only using few
 101 representatives features. More specifically, the method is made up of the following steps:

Algorithm 1: Sparse-Modeling Based Approach for Class-Specific Feature Selection

Input : $X = \{x_1, \dots, x_n\}$ data set
 y , class labels
 θ , SMBA parameters
 m , maximum number of features to select
 C , classifier model (e.g., SVM, KNN, etc)
 K , number of folds for performing K-Cross Validation

Output : \overline{ACM} , Averaged Classification Metrics on K folds

```

1 begin
2    $X \leftarrow$  Data standardization
3    $X \leftarrow$  Class balancing( $X$ ) by using SMOTE (Chawla et al., 2002)
4    $X \leftarrow$  Random shuffling( $X$ )
5   Divide  $X$  into  $K$  folds
6   foreach  $k_i \in K$  folds do
7     Set the  $k_i$  fold as the test set  $X_{test}$ 
8     Use the remaining  $K-2$  folds as the train set  $X_{train}$ 
9     Perform the Class-sample separation on the train set  $X_{train}$ 
10    (Note that  $I$  is the subset of features selected for each class  $c_i \in X_{train}$ )
11    foreach  $X_{c_i} \in X_{train}$  do
12       $I = \{I_{c_i} \dots I_{c_c}\} \leftarrow$  SMBA( $X_{c_i}, \theta$ )
13    end
14    for  $j \leftarrow 1$  to  $m$  do
15      Build an ensemble classifier  $E_j = \{e_{1,j}, \dots, e_{c,j}\}$  using the  $j$ -th selected feature  $\in I_{c_i}$ 
16      and the classifier  $C$ 
17      foreach  $O \in X_{test}$  do
18         $(ACM_j) \leftarrow$  Use  $E_j$  to classify the instance  $O$ 
19      end
20       $(ACM) \leftarrow (ACM_j)$ 
21    end
22   $(\overline{ACM}) \leftarrow$  Average( $ACM$ )
23 end
```

- 102 1. **Class-sample separation:** Unlike the GF-CSFS, SMBA-CSFS does not employ the *Class binarization* stage to transform a c -class problem into c binary problems, instead it just uses a simple
 103 *Class-sample separation*. It simply consists of differentiating the samples among all the classes of
 104 the training set for a given data set into several disjoint sets/configurations of samples, one for each
 105 class (See Fig. 1).
 106
- 107 2. **Class balancing:** Once the class sample set of the training set has been split apart (by applying the
 108 above *Class-sample separation* step), it may be possible that each class-subset results unbalanced.
 109 Therefore, the SMOTE (Chawla et al., 2002) re-sampling method is applied to balance each class-
 110 subset. Technically speaking, it is important to point out that steps 1-2 are interchangeable, meaning
 111 that there are no differences in doing the first one before the other.
- 112 3. **Intra-Class-Specific feature selection:** The *Sparse-Modeling Based Approach* is used for retriev-
 113 ing, minimizing equation 7, the most representative features for each class-sample set of the training
 114 set that best represent/reconstruct the whole class of objects. In doing so, the approach takes
 115 advantage of the intra-class properties for selecting the best feature subset (describing each class)
 116 which is used to improve the classification accuracy against TFS and GF-CSFS.
- 117 4. **Classification:** Since the training set gets split into different class-sample subsets, we embraced
 118 the idea of using a wise-ensemble procedure for training a classification model for discriminating
 119 new incoming instances. As in (Pineda-Bautista et al., 2011), given a class c_i , a classifier e_i is

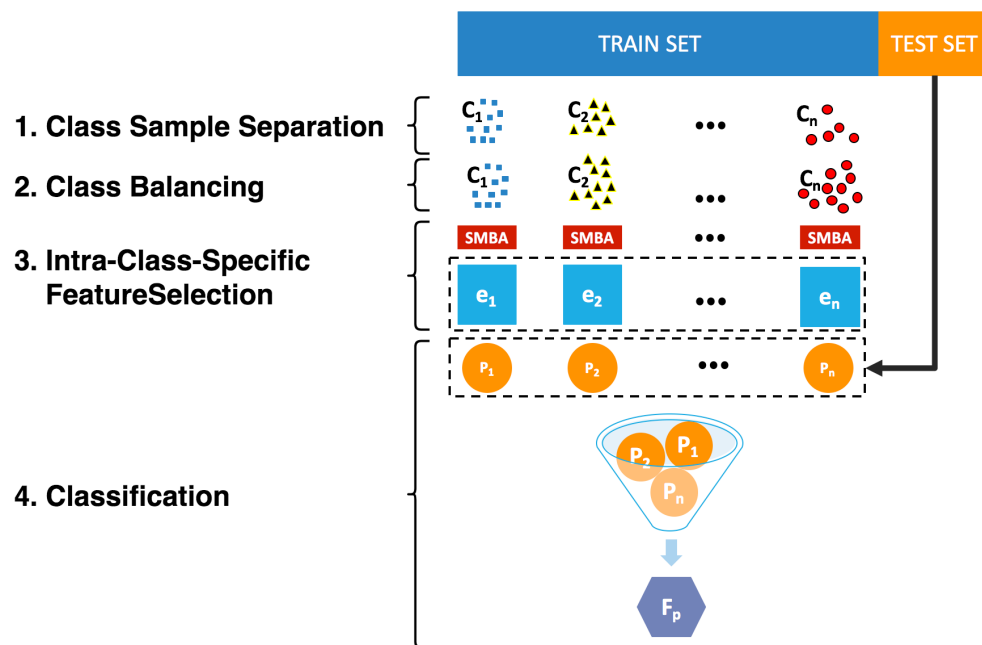


Figure 1. A Sparse-Modeling Based approach for Class-Specific Feature Selection.

trained on the original data set only using the selected features for c_i , for $i = 1, \dots, c$. Overall, a classifier ensemble $E = \{e_1, \dots, e_c\}$ is constructed. In order to classify a new instance O through the ensemble, the natural dimension of O needs to be lowered to the dimension d_i of the classifier e_i , $i = 1 \dots, c$. This way, for determining to which class O belongs to, an *ad-hoc majority rule* is used:

- If a classifier outputs the same class for which the features, used for e_i training, were selected, i.e., the e_i output is c_i , then O belongs to c_i . In case of a tie, i.e., when several classifiers respond c_i , a *majority vote* is needed among all classifiers to determine the class of O . If still a tie occurs, O will belong to the class that received more votes among the tied classes.
- If no classifier outputs the class whose selected features are used for e_i training, O belongs to the class winning the majority voting. If there is a tie, then O will belong to the class that received more votes among the tied classes.

Finally, since a recursive tie may occur, in that case, the instance O would be classified as c_i by randomly choosing a class among all the tied classes.

The algorithm in Fig. 1, illustrates the pseudo-code describing the CSFS-SMBA procedure. Basically, it first standardizes, *class-balances* and shuffles the data set X , then divides it into K -folds, assigning the k_i -th fold as test set X_{test} and the remaining $K - 1$ folders as train set X_{train} . The algorithm iteratively performs the task of *class-sample separation*, to split the sample belonging to different classes (X_{c_i}) on which the algorithm SMBA (illustrated in page 4) is performed and then output the m most representative features for each class (line 12). The selected features are then used, one at time, for training an ensemble classifier E_j , and later used for classifying each instance O belonging to the test set X_{test} . Finally, for all the ensemble models up to m selected features, the algorithm outputs the \overline{ACM} matrix, storing several model evaluation metrics.

EXPERIMENTAL RESULTS

In the experiments, the SMBA-CSFS performance have been assessed on nine publicly available microarray data sets. The classifier used to determine the goodness of the selected feature subsets are a Support Vector Machine (SVM) with a linear kernel and parameter $C = 1$, a Naive Bayes, a K-Nearest Neighbors (KNN) using $k = 5$, and a Decision Tree.

Data sets Description

In order to validate the introduced approach, a number of data sets exemplifying the typical data processing in the biological field are used in the experiments. In the following, a brief description of all data sets employed in the experiments.

1. The **ALLAML** data set (Golub et al., 1999) contains in total 72 samples in 2 classes: ALL and AML, which have 47 and 25 samples, respectively. Every sample contains 7,129 gene expression values.
2. The **LEUKEMIA** data set (Golub et al., 1999) contains in total 72 samples in 2 classes: acute lymphoblastic and acute myeloid. From 7,129 genes, the baseline genes were cut off before further analysis. The number of genes that are used in the binary classification task is 7,070.
3. The **CLL_SUB_111** data set (Haslinger et al., 2004) has gene expressions from high density oligonucleotide arrays containing genetically and clinically distinct subgroups of B-cell chronic lymphocytic leukemia (B-CLL). The data set consists of 11,340 attributes, 111 instances and 3 classes.
4. The **GLIOMA** data set (Nutt et al., 2003) contains in total 50 samples in 4 classes: cancer glioblastomas, non-cancer glioblastomas, cancer oligodendrogliomas and non-cancer oligodendrogliomas, which have 14, 14, 7, 15 samples, respectively. Each sample has 12,625 genes. After a preprocessing, the data set has been shrunk to 50 samples and 4,433 genes.
5. The **LUNG** data set (Bhattacharjee et al., 2001) contains in total 203 samples in 5 classes: adenocarcinomas, squamous cell lung carcinomas, pulmonary carcinoids, small-cell lung carcinomas and normal lung, with 139, 21, 20, 6, 17 samples, respectively. The genes with standard deviations smaller than 50 expression units were removed (the interested reader may refer to (Bhattacharjee et al., 2001) for details) getting a data set with 203 samples and 3,312 genes.
6. The **LUNG_DISCRETE** data set (Peng et al., 2005) contains 73 samples in 7 classes where, each sample consists of 325 gene expressions. The cardinalities of each sample in the LUNG_DISCRETE data set are 6, 5, 5, 16, 7, 13, 21, respectively.
7. The **DLBCL** data set (Alizadeh et al., 2000) is a modified version of the original DLBCL data set. It consists of 96 samples in 9 classes, where each sample is defined by the expression of 4,026 genes. The cardinalities of each sample in the DLBCL data set are 46, 10, 9, 11, 6, 6, 4, 2, 2, respectively.
8. The **CARCINOM** data set (Su et al., 2001) contains 174 samples in 11 classes: prostate, bladder/ureter, breast, colorectal, gastroesophagus, kidney, liver, ovary, pancreas, lung adenocarcinomas and lung squamous cell carcinoma, with 26, 8, 26, 23, 12, 11, 7, 27, 6, 14, 14 samples, respectively. After a preprocessing as described in (Yang et al., 2006), the data set has been shrunk to 174 samples and 9,182 genes.
9. The **GCM** data set (Ramaswamy et al., 2001) contains 190 samples in 14 classes: breast, prostate, lung, colorectal, lymphoma, bladder, melanoma, uterus, leukemia, renal, pancreas, ovary, mesothelioma and central nervous system, where each sample consist of 16,063 gene expression signatures. The cardinalities of each sample in the data set are 11, 11, 20, 11, 30, 11, 22, 10, 11, 11, 11, 10, 11, 10, respectively.

All data sets have been originally downloaded from the following source, migrated at later time at the following data repository (Nardone et al., 2019a). All the information about the data sets are summarized in Table 1.

Experiment Setup

To validate the effectiveness of the SMBA-CSFS model, it has been compared against several TFS and the GF-CSFS proposed in (Pineda-Bautista et al., 2011). SMBA-CSFS is firstly compared against TFS methods and, since the framework in (Pineda-Bautista et al., 2011) can use any TFS method as base for doing CSFS, some experiments using both filter and wrapper methods (injection process) were made. In

addition, the accuracy results were also compared against those obtained on the basis of all the features (BSL). The following TFS methods have been chosen for comparing purposes:

- **LASSO** (Tibshirani, 1994): It involves penalizing the absolute size of the regression coefficients and is usually used for creating parsimonious models in presence of a *large* number of features. The model implemented is a modified version of classical LASSO, adapted for classification purposes. In particular, in Equation 8, the product $\mathbf{X}\beta$ is transformed by a sigmoid function in order to address the classification problem.
- **EN** (Zou and Hastie, 2005): Elastic Net is a hybrid of ridge regression and LASSO regularization. Like lasso, Elastic Net can generate reduced models by generating zero-valued coefficients. Experimental studies have suggested that the Elastic Net technique can outperform LASSO on data with highly correlated features. As for LASSO, a modified version adapted for classification purposes has been implemented.
- **RFS** (Nie et al., 2010): Robust Feature Selection method is a sparse based-learning approach for feature selection which emphasizes the joint $\ell_{2,1}$ norm minimization on both loss and regularization function.
- **Is- $\ell_{2,1}$** (Tang et al., 2014): Is- $\ell_{2,1}$ is a supervised sparse feature selection method. It exploits the $\ell_{2,1}$ -norm regularized regression model for joint feature selection, from multiple tasks where the *classification objective function* is a quadratic loss.
- **Il- $\ell_{2,1}$** (Tang et al., 2014): Il- $\ell_{2,1}$ is a supervised sparse feature selection method which uses the same concept of Is- $\ell_{2,1}$ but instead uses a *logistic loss* as *classification objective function*.
- **Fisher** (Gu et al., 2012): Fisher is one of the most widely used supervised filter feature selection methods. It selects each feature as the ratio of inter-class separation and intraclass variance, where features are evaluated independently and, the final feature selection occurs by aggregating the m top ranked ones.
- **Relief-F** (Kira and Rendell, 1992; Kononenko, 1994): Relief-F is an iterative, randomized and supervised filter approach that estimates the quality of the features according to how well their values differentiate data samples that are near to each other; it does not discriminate among redundant features and performance decreases with few data.
- **mRmR** (Peng et al., 2005): Minimum-Redundancy-Maximum-Relevance is a mutual information filter based algorithm which selects features according to the maximal statistical dependency criterion.
- **MI** (Kraskov et al., 2004; Ross, 2014): Mutual Information is a non-negative value, which measures the dependency between the variables. Features are selected in a univariate way. The function relies on nonparametric methods based on entropy estimation from k-nearest neighbors distances.
- **SMBA**: Sparse-Modeling Based Approach is nothing else than our SMBA-CSFS model but that only takes into account the SDL strategy for selecting a subset of features considering all the classes in the feature selection process.

We pre-processed all the data sets by using the *Z-score* (Kreyszig, 2010) normalization. To fairly compare the considered supervised feature selection methods, we have firstly tuned the parameters for all methods by using a “grid-search” strategy (Tang et al., 2014) and finally, for evaluating the performance of all the methods, it has been considered a number of features ranging from 1 to 80, performing a 5-fold CV to report the average results along with the standard deviations (STD).

The performance of the classification algorithms among all the methods have been evaluated by using the metrics of Accuracy (ACC), Precision (P), Recall (R) and F-measure (F), which are computed as illustrated in (Sokolova and Lapalme, 2009). In addition, to give a better and summarized understanding between the performance of the models, we also computed the Area Under the Curve (AUC) and the Receiver Operating Characteristic (ROC) curves, where the former is a useful tool for evaluating the quality of class separation for a classifier while the latter makes it easier to compare the ROC curve of one model to another.

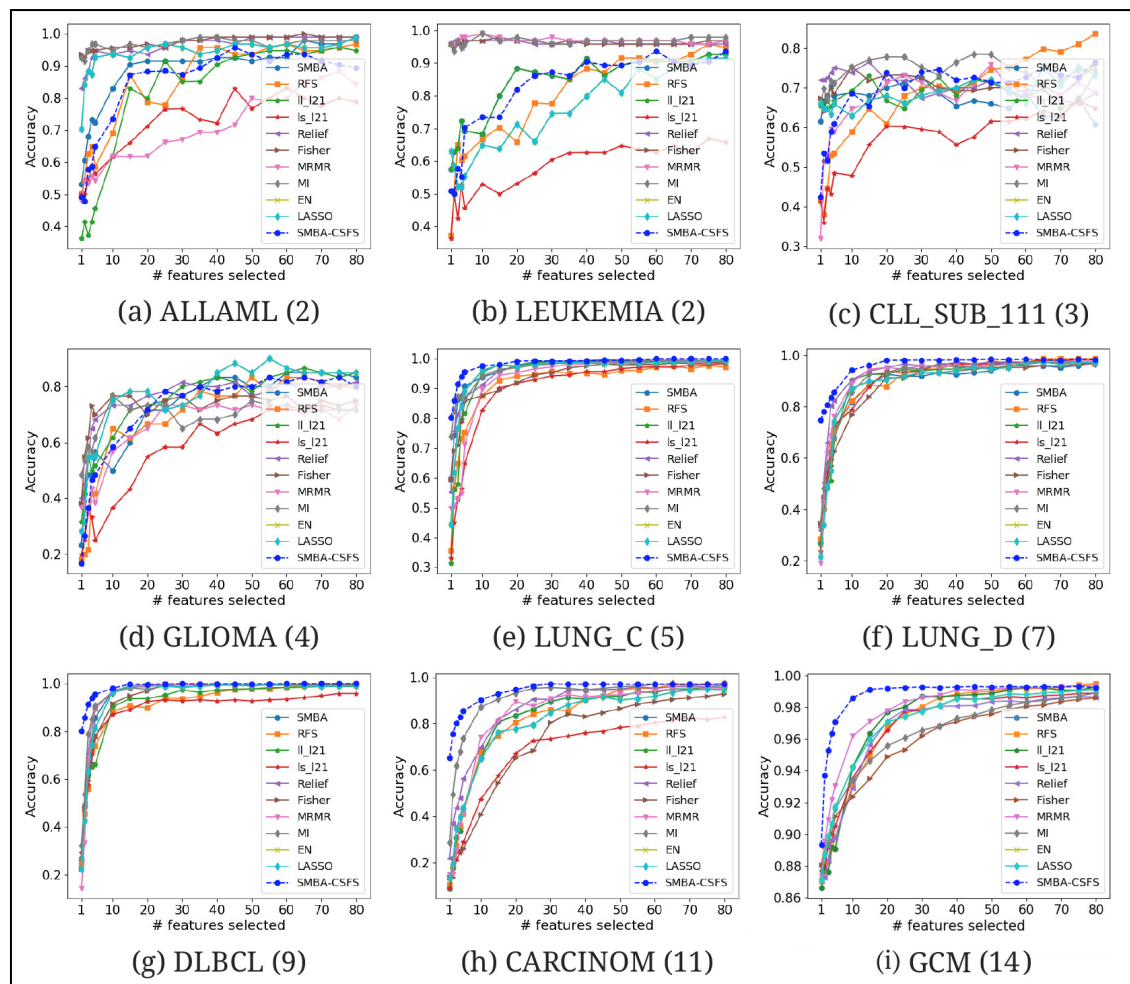


Figure 2. Comparison of several TFS accuracies against SMBA and SMBA-CSFS on nine data sets: (a) ALLAML(2), (b) LEUKEMIA(2), (c) CLL.SUB_111(3), (d) GLIOMA(4), (e) LUNG.C(5), (f) LUNG.D(7), (g) DLBCL(9), (h) CARCINOM(11), (i) GCM(14), when a varying number of features is selected. SVM classifier with 5-fold CV was used.

DISCUSSION

The experiments have been performed on a workstation with a dual Intel(R) Xeon(R) 2.40GHz and 64GB RAM. The developed code is available at (Nardone et al., 2019b). For the sake of readability, all the results presented here account only for the SVM classifier, since the performance proved that the proposed approach is a little sensitive to the choice of a specific classifier (indeed, the performance of each classifier are rather comparable). Nevertheless, the interested reader may refer to the Supplemental material for details on additional results concerning all the used classifiers. The experimental results on 5-CV for the SVM classifier are summarized in the Tables 2-5. Figures 2-5 show all the accounted model evaluation metrics for the ten feature selection methods on the nine considered data sets.

We compared the performance of our method against TFS methods (see Tables 2-3) and GF-CSFS framework (see Tables 4-5). By looking at precision, recall, F-measure and accuracy, SMBA-CSFS is able to better discriminate among the classes of the LUNG_C, LUNG_D, CAR, DLBLC, GCM data sets in most of the cases, when top 20 and 80 features are considered. In this latter case, when SMBA-CSFS performs worse than its competitors, the corresponding performance tend to be comparable. On the remaining data sets, each with a number of classes less than 5, namely, ALLAML, LEUKEMIA, CLL.SUB_111 and GLIOMA, SMBA-CSFS is instead outperformed by some of the competitors. Consequently, we can assert that SMBA-CSFS behaves better when working with data sets with many classes (at least 5). One possible reason is due to the sparse-modeling approach in selecting the features and the use of an

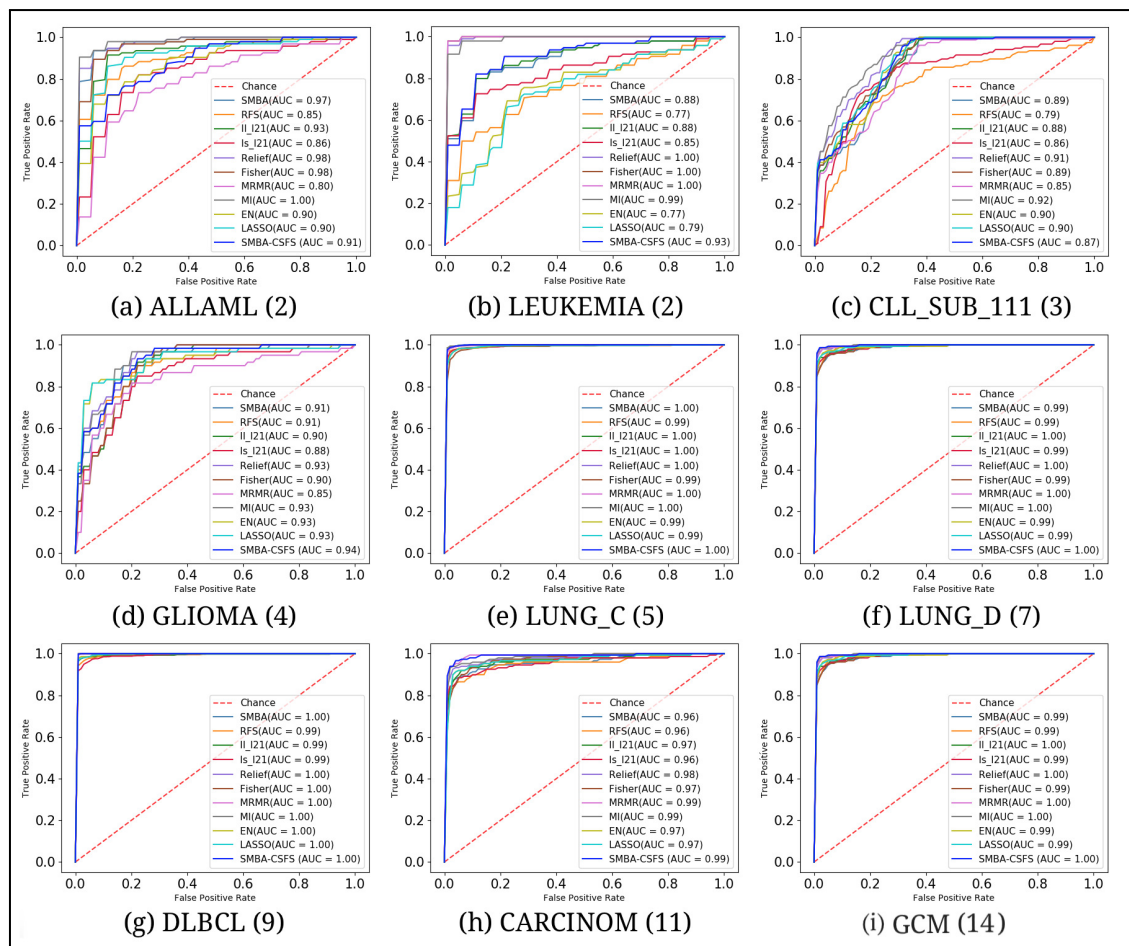


Figure 3. Averaged ROC curves on the first 20 features comparing the classification performance among SMBA-CSFS and TFS methods for nine data sets: (a) ALLAML(2), (b) LEUKEMIA(2), (c) CLL.SUB_111(3), (d) GLIOMA(4), (e) LUNG_C(5), (f) LUNG_D(7), (g) DLBCL(9), (h) CARCINOM(11), (i) GCM(14). SVM classifier with 5-fold CV was used.

ensemble classifier. Indeed, since the ensemble is based on a majority voting schema, SMBA-CSFS is able to guess, with higher probability, the belonging of samples coming from data sets with many classes. Just think that, whenever our method draws from a sample of a two-class data set, the probability of a right guess is proportional to a coin toss. Therefore if, on one hand, this leads to good performance when the data set consists of many classes, the probability of failure, on the other hand, increases in the case of data sets consisting of fewer classes. Anyhow, the local structure of data distribution which is crucial for feature selection, as stated in (He et al., 2005), may be a logical reason why the SBMA schema performs better on certain data set rather than others. In addition, as shown in Fig. 2, it is worth observing that SMBA-CSFS seems perform better w.r.t. TFS competitors on a fewer number of features. This would suggest that SMBA-CSFS is able to identify/retrieve the most representative features that maximize the classification accuracy. To assert the previous results achieved, we computed the averaged ROC curves between SMBA-CSFS and the other TFS methods on a subset of 20 and 80 features, respectively. Looking at the AUC values in Fig. 3, it would suggest SMBA-CSFS as the best model to choose for identifying the most representatives features in a classification task when dealing with data set with many classes. Concerning with the GF-CSFS competitors, as shown in Fig. 4, it would suggest that the *sparse modeling* process, underlying the proposed SMBA scheme for feature selection, is more suitable for retrieving the best features for the purpose of classification, often leading to get satisfactory results. Such statement is also proved by the good balance between precision and recall showed in Table 5 and the average ROC curves showed in Fig. 5, where SMBA-CSFS still hold a candle w.r.t. GF-CSF methods. The readers

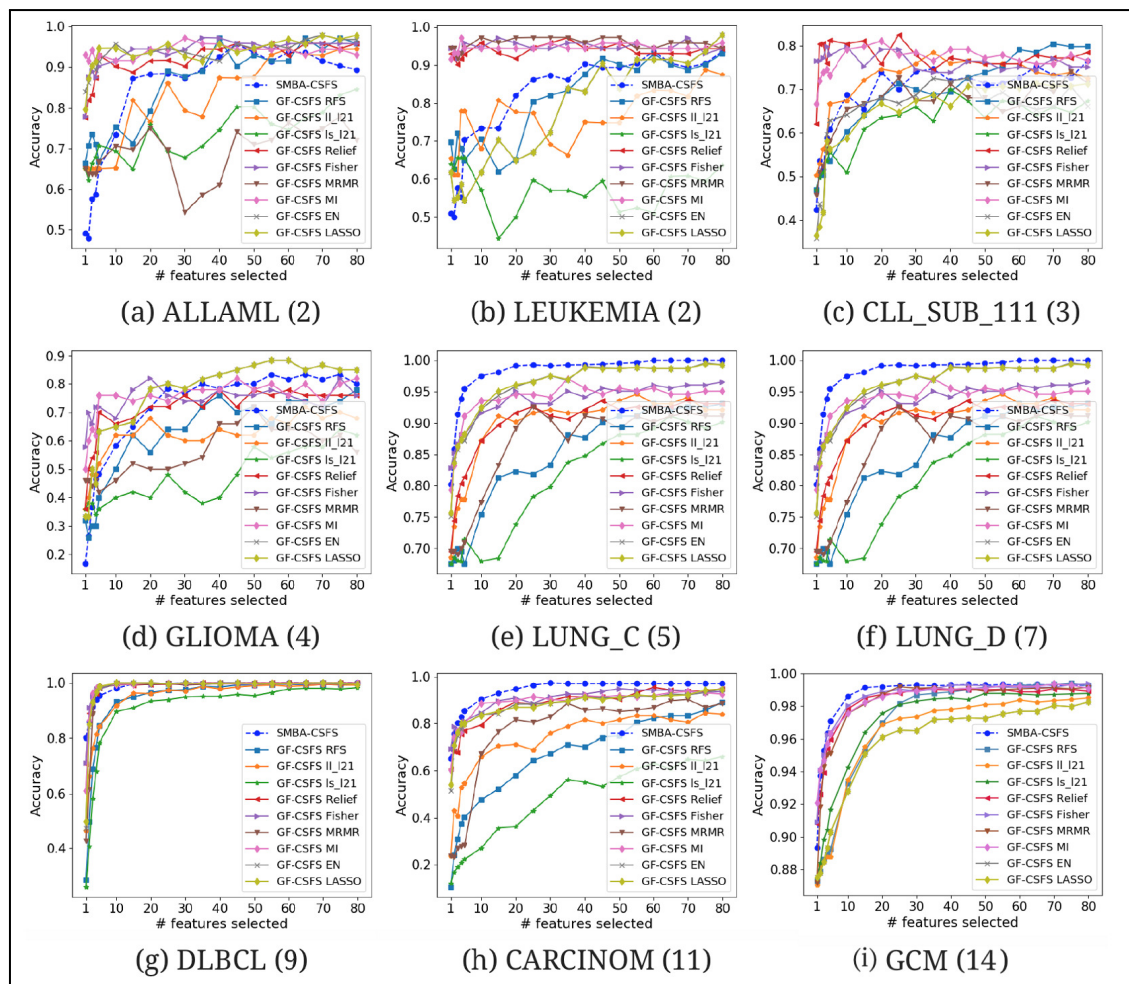


Figure 4. Comparison of several CSFS accuracies against SMBA-CSFS on nine data sets: (a) ALLAML(2), (b) LEUKEMIA(2), (c) CLL.SUB.111(3), (d) GLIOMA(4), (e) LUNG.C(5), (f) LUNG.D(7), (g) DLBCL(9), (h) CARCINOM(11), (i) GCM(14), when a varying number of features is selected. SVM classifier with 5-fold CV was used.

attention is drawn to the Supplemental material for all the experimental results and consideration arisen on the top 80 features.

To statistically validate the results and compare all the competing classifiers against the proposed SMBA-CSFS, on both 20 and 80 feature subsets, we ran *Non-Parametric multiple comparison tests (all vs all)* (Demšar, 2006; Rodríguez-Fdez et al., 2015) which sequentially performs a popular multi-class *Friedman nonparametric test* (Friedman, 1937) followed by a *Nemenyi Post-hoc multiple comparison* (Dunn, 1961). The ranking of the classifiers, when the top 20 and 80 features are selected, along with the corresponding p-values, are described in the supplementary material. Looking at the *Cumulative Rank (CR)* for each classifier, one notices how SMBA-CSFS achieves optimal results (e.g., always ranks within the first three places). However, it is worth emphasizing that our method ranks systematically on the top place when considering data sets consisting of five or more classes (named $CR_{\geq 5}$). These results prove again that SMBA-CSFS has good performance on data sets with many classes. Moreover, by using different classifiers we do not observe noteworthy differences in the results, meaning that the methodology is suitable for the classification of this kind of data, independently from the selected classifier. However, by looking at the p-values, corresponding to the single ranking method, one can better verify which algorithms have significantly different performance w.r.t. SMBA-CSFS. For detailed information regarding the results, see the Supplemental material. Concerning the computational complexity, from several conducted experiments we observed that the proposed methodology might be slower than other

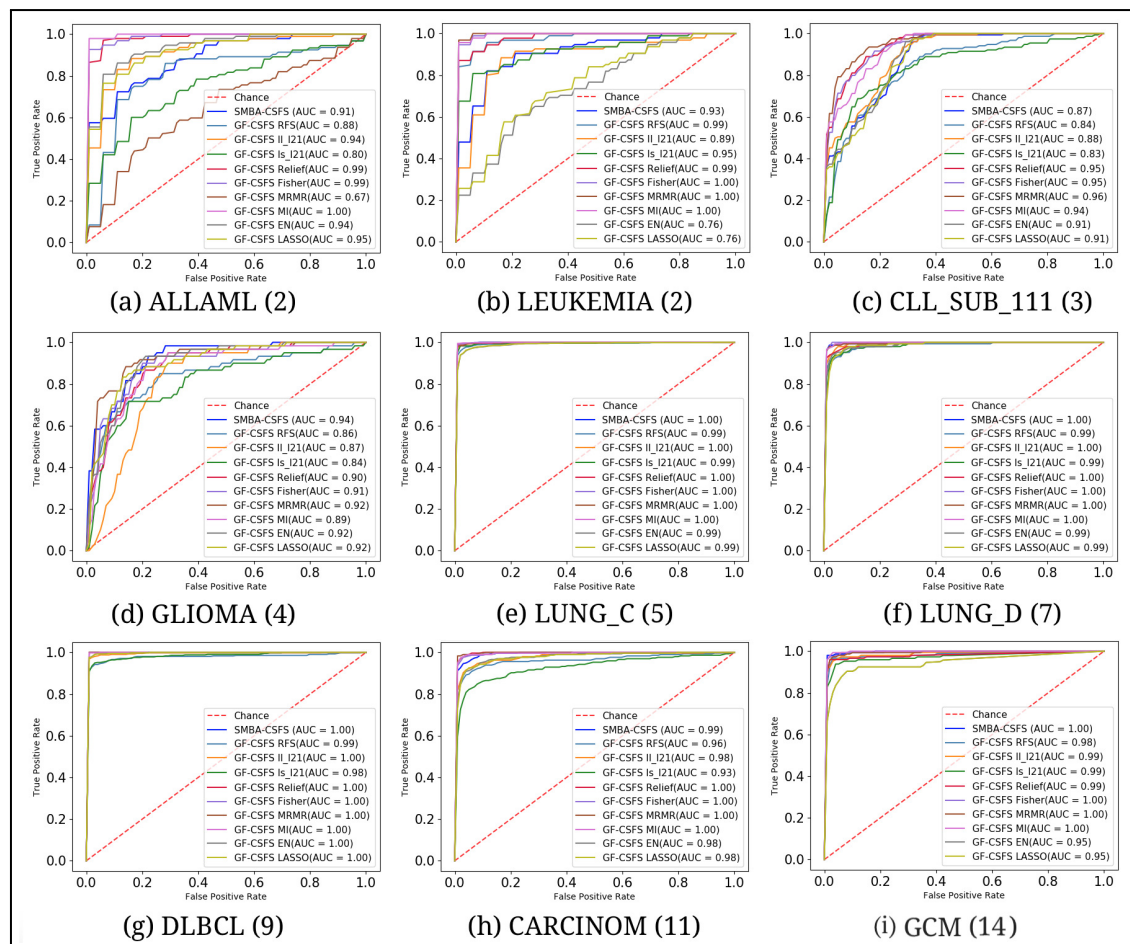


Figure 5. Averaged ROC curves on the first 20 features comparing the classification performance among SMBA-CSFS and several CSFS methods for nine data sets: (a) ALLAML(2), (b) LEUKEMIA(2), (c) CLL.SUB.111(3), (d) GLIOMA(4), (e) LUNG_C(5), (f) LUNG_D(7), (g) DLBCL(9), (h) CARCINOM(11), (i) GCM(14). SVM classifier with 5-fold CV was used.

techniques (e.g., FS and Relief whose running times are in term of few seconds) but comparable with SMBA. Its running time, depending on several parameters involved, especially in the size of the number of instances and classes of the data sets, might vary from a couple of hours to at most one day (see Table S9, in the Supplementary material, for details on the computational time). Nevertheless, SMBA-CSFS has appreciable performance when working on large data sets and number of classes, and sometimes, in the biological field, the accuracy in finding key features that are responsible for some biological processes would be preferred to the execution time. However, since most of the time consumed by the proposed approach is due to the solution of the optimization problem by using the ADMM method, and because the methodology is based on an ensemble of classifiers, a parallel computing approach could be adopted to obtain a faster computational time (Deng et al., 2017).

CONCLUSIONS

We proposed a Sparse-Modeling Based Approach for Feature Selection with emphasizing joint $\ell_{1,2}$ -norm minimization and the Class-Specific Feature Selection. Experimental results, on nine different data sets, validate the unique aspects of SMBA-CSFS and demonstrate the promising performance achieved against the-state-of-art methods. One of the main characteristics of our framework is that, by jointly exploiting the idea of Sparse Modeling and Class-Specific Feature Selection, it is able to identify/retrieve the most representative features that maximize the classification accuracy in those cases where a given data set is made up of many classes. Based on our experimental results, we can conclude that, usually applying TFS

allows achieving better results than using all the available features. However, in many cases, applying the proposed SMBA-CSFS method allows improving the performance of just TFS as well as GF-CSFS injected with several TFS methods. It has to be stressed, that SMBA-CSFS seems actually suitable for large data sets consisting of many classes, while on data sets with less than five classes other methods appear to be more effective. Although SMBA, SMBA-CSFS and TFS performance slightly differ on the whole, it is worth highlighting that SMBA-CSFS achieves its best performance when considering fewer features (i.e., from 1 to 20) on data sets with many classes, which is an important goal when certain biological tasks are taken into account. However, we do believe that these techniques might be effectively used in a systematic way after a microarray analysis. Indeed, a better gene selection step could avoid the waste of many resources in post-array wet analysis (e.g., Real Time-PCR) allowing researchers to focus their attention just on relevant features. Finally, we think this method demonstrated to be an interesting alternative among FS approaches on microarray data.

As future work, the focus will be moved towards the biologic interpretations of the SMBA framework behavior, by systematically studying the selected genes, especially taking into account the SMBA-CSFS approach which, as proved by the experimental results, is more effective in selecting genes of interest than the standard SMBA. Furthermore, we are planning to test our approach on EPIC data set (Demetriou et al., 2013), after a thorough analysis of pre-filtering, and a parallel implementation to substantially reduce its computational time.

1 ACKNOWLEDGEMENTS

The research was entirely developed when Davide Nardone was a Master Degree student in Applied Computer Science at University of Naples Parthenope.

Table 1. Data sets Description.

	Size	# of Features	# of Classes
ALLAML	72	7,129	2
LEUKEMIA	72	7,070	2
CLL_SUB_111	111	11,340	3
GLIOMA	50	4,434	4
LUNG_C	203	3,312	5
LUNG_D	73	325	7
DLBCL	96	4,026	9
CAR	174	9,182	11
GCM	190	16,063	14

Table 2. SVM accuracy results (ACC±STD) on top 20 features using 5-fold CV on different data sets.

TFS methods are compared against our methods (SMBA and SMBA-CSFS). FS: Fisher Score, mRmR: Minimum-Redundancy-Maximum-Relevance, MI: Mutual Information, RFS: Robust Feature Selector, EN: Elastic Net, BSL: all features. The best results are highlighted in bold. The number in parentheses is the number of features when the performance is achieved.

	Average Accuracy of top 20 features (%)															
	ALLAML	LEUKEMIA	CLL.SUB.111	GLIOMA	LUNG.C	LUNG.D	DLBCL	CAR	GCM							
Fisher	96.84±0.04(19)	98.95±0.02(16)	75.20±0.1(19)	80±0.04(13)	91.94±0.02(19)	91.24±0.1(20)	97.11±0.02(19)	65.33±0.05(20)	94.9±0.00(20)							
Relief	95.78±0.04(8)	97.89±0.03(12)	76.45±0.03(15)	80±0.07(19)	97.12±0.01(20)	95.2±0.03(14)	99.76±0.00(20)	86.52±0.03(18)	97.14±0.01(20)							
mRmR	66.14±0.13(12)	98.95±0.02(9)	71.27±0.1(20)	66.67±0.1(17)	95.68±0.013(19)	95.22±0.02(20)	99.03±0.01(16)	89.57±0.04(20)	97.79±0.01(20)							
MI	96.84±0.04(215)	98.95±0.02(10)	81.03±0.06(17)	78.33±0.04(12)	97.41±0.014(17)	94.53±0.03(18)	98.79±0.01(19)	93.25±0.05(20)	95.58±0.01(20)							
Is-21	71.34±0.14(19)	59.42±0.2(12)	60.30±0.14(19)	55±0.07(20)	92.66±0.05(19)	93.86±0.04(20)	92.52±0.01(20)	66.99±0.03(20)	96.56±0.01(20)							
l1-21	83±0.11(15)	88.36±0.06(20)	73.12±0.06(15)	0.75±0.12(17)	98.27±0.015(16)	93.24±0.04(16)	94.44±0.02(19)	83.49±0.03(20)	97.69±0.01(20)							
RFS	87±0.01(15)	74.33±0.1(18)	64.73±0.09(15)	66.67±0.07(17)	94.10±0.022(20)	89.77±0.02(19)	91.06±0.03(18)	81.85±0.07(18)	96.77±0.01(20)							
LASSO	98.95±0.02(17)	71.3±0.08(21)	68.02±0.06(20)	83.33±0.05(17)	97.99±0.012(16)	92.51±0.03(12)	99.52±0.01(16)	82.14±0.05(18)	97.07±0.01(20)							
EN	98.95±0.02(17)	71.3±0.08(21)	68.02±0.06(20)	83.33±0.05(17)	97.99±0.012(16)	92.51±0.03(12)	99.52±0.01(16)	82.14±0.05(18)	97.07±0.01(20)							
SMBA	93.68±0.084(16)	88.36±0.06(20)	70.60±0.10(19)	71.67±0.134(17)	97.84±0.00(20)	92.55±0.03(20)	99.28±0.01(20)	83.49±0.03(20)	97.69±0.01(20)							
SMBA-CSFS	88.24±0.04(20)	81.93±0.02(20)	75.53±0.06(20)	73.34±0.18(16)	98.41±0.014(19)	97.93±0.03(19)	98.30±0.02(13)	94.95±0.02(19)	99.2±0.01(20)							
BSL	97.89±0.04	98.95±0.021	84.26±0.06	85±0.1	99.57±0.00	98.62±0.02	100±0.00	98.65±0.01	100±0.00							

Table 3. SVM Precision(P), Recall(R) and F-measure(F) on top 20 features using 5-fold CV on different data sets.

TFS methods are compared against our methods (SMBA and SMBA-CSFS). FS: Fisher Score, mRmR: Minimum-Redundancy-Maximum-Relevance, MI: Mutual Information, RFS: Robust Feature Selector, EN: Elastic Net, BSL: all features. The best results are highlighted in bold.

	ALLAML				LEUKEMIA				CLL.SUB.111				GLIOMA				LUNG.C				LUNG.D				DLBCL				CARCINOM				GCM(14)			
	P	R	F		P	R	F		P	R	F		P	R	F		P	R	F		P	R	F		P	R	F		P	R	F		P	R	F	
Fisher	0.98(18)	0.98(18)	0.98		0.99(15)	0.99(15)	0.99		0.75(11)	0.75(11)	0.75		0.68(20)	0.67(14)	0.67		0.92(19)	0.92(19)	0.92		0.89(20)	0.88(15)	0.88		0.9(17)	0.90(20)	0.93		0.9(19)	0.89(20)	0.89		0.64(20)	0.64(20)	0.64	
Relief	0.96(12)	0.96(12)	0.96		0.99(4)	0.99(4)	0.99		0.75(17)	0.75(17)	0.75		0.77(19)	0.77(19)	0.77		0.97(20)	0.97(20)	0.97		0.95(20)	0.95(15)	0.95		0.89(18)	0.88(18)	0.88		0.89(18)	0.88(18)	0.88		0.8(20)	0.8(20)	0.8	
MRMR	0.8(19)	0.8(19)	0.8		0.98(6)	0.98(17)	0.98		0.64(14)	0.64(14)	0.65		0.71(12)	0.71(12)	0.7		0.97(20)	0.97(20)	0.97		0.96(19)	0.95(19)	0.95		0.95(20)	0.94(14)	0.92		0.88(20)	0.91(20)	0.89		0.85(20)	0.85(20)	0.85	
MI	0.98(12)	0.98(12)	0.98		0.98(2)	0.98(2)	0.98		0.76(16)	0.76(16)	0.76		0.74(20)	0.73(17)	0.73		0.97(20)	0.97(20)	0.97		0.95(20)	0.95(20)	0.95		0.95(17)	0.95(17)	0.93		0.69(20)	0.69(20)	0.69		0.69(20)	0.69(20)	0.69	
l1-21	0.92(15)	0.91(15)	0.91		0.83(20)	0.83(20)	0.83		0.72(20)	0.72(20)	0.7		0.71(16)	0.71(17)	0.7		0.97(20)	0.97(20)	0.97		0.88(19)	0.88(19)	0.88		0.81(19)	0.81(20)	0.81		0.76(20)	0.76(20)	0.76		0.84(20)	0.84(20)	0.84	
RFS	0.86(18)	0.84(19)	0.85		0.84(20)	0.76(20)	0.8		0.63(12)	0.64(12)	0.63		0.71(12)	0.71(12)	0.7		0.96(19)	0.96(19)	0.96		0.88(18)	0.86(18)	0.87		0.89(19)	0.93(16)	0.84		0.89(18)	0.84(19)	0.86		0.77(20)	0.77(20)	0.77	
LASSO	0.84(20)	0.84(13)	0.84		0.77(20)	0.77(20)	0.77		0.71(16)	0.71(10)	0.71		0.79(14)	0.78(14)	0.78		0.94(20)	0.94(19)	0.94		0.93(19)	0.92(20)	0.91		0.84(18)	0.84(18)	0.84		0.84(18)	0.84(18)	0.84		0.8(20)	0.8(20)	0.8	
EN	0.84(20)	0.84(13)	0.84		0.77(20)	0.77(20)	0.77		0.71(16)	0.71(10)	0.71		0.79(14)	0.78(14)	0.78		0.94(20)	0.94(19)	0.94		0.93(19)	0.92(20)	0.91		0.84(18)	0.84(18)	0.84		0.84(18)	0.84(18)	0.84		0.8(20)	0.8(20)	0.8	
SMBA	0.9(13)	0.89(16)	0.89		0.83(20)	0.83(20)	0.83		0.71(11)	0.71(11)	0.7		0.68(15)	0.68(15)	0.68		0.97(18)	0.97(18)	0.97		0.91(19)	0.91(19)	0.9		0.92(19)	0.99(17)	0.92		0.9(19)	0.86(20)	0.88		0.84(20)	0.84(20)	0.84	
SMBA-CSFS	0.83(16)	0.83(16)	0.83		0.86(20)	0.86(20)	0.86		0.67(20)	0.68(20)	0.67		0.8(20)	0.77(20)	0.78		0.98(15)	0.98(15)	0.98		0.99(19)	0.99(19)	0.99		1.0(20)	1.0(20)	1.0		0.99(20)	0.98(20)	0.98		0.97(20)	0.97(20)	0.97	
BSL	1	1	1		1	1	1		0.74	0.74	0.74		0.92	0.92	0.92		0.93	0.93	0.93		0.8	0.8	0.8		1	1	1		0.98	0.98	0.98		1	1	1	

Table 4. SVM accuracy results (ACC±STD) on top 20 features using 5-fold CV on different data sets.

GF-CSFS (Pineda – Bautista et al., 2011) framework is compared against our SMBA-CSFS. FS: Fisher Score, mRmR: Minimum-Redundancy-Maximum-Relevance, MI: Mutual Information, RFS: Robust Feature Selector, EN: Elastic Net, BSL: all features. The best results are highlighted in bold. The number in parentheses is the number of features when the performance is achieved.

	Average Accuracy of top 20 features (%)															
	ALLAML	LEUKEMIA	CLL.SUB.111	GLIOMA	LUNG.C	LUNG.D	DLBCL	CAR	GCM							
Fisher	95.90±0.03(13)	98.57±0.03(18)	80.41±0.02(7)	82±0.16(17)	95.09±0.03(20)	86.38±0.14(16)	100±0.00(14)	90.86±0.08(20)	98.98±0.0(18)							
Relief	92.95±0.04(5)	95.81±0.03(10)	82.41±0.05(12)	80±0.19(12)	91.63±0.02(20)	86.39±0.07(20)	100±0.00(11)	89.68±0.03(17)	98.71±0.0(20)							
mRmR	75.14±0.09(16)	98.57±0.03(11)	70.69±0.07(12)	62±0.12(14)	89.16±0.03(20)	86.48±0.09(17)	99.52±0.01(15)	81.61±0.07(20)	98.71±0.0(20)							
MI	94.38±0.03(18)	97.14±0.03(4)	81.03±0.05(20)	82±0.21(19)	95.07±0.015(11)	79.90±0.18(14)	100±0.00(19)	90.86±0.06(11)	98.67±0.0(19)							
Is-21	76.47±0.13(6)	65.52±0.08(3)	63.44±0.03(20)	46±0.21(7)	73.88±0.04(19)	75.43±0.07(18)	93.46±0.03(20)	39.68±0.04(19)	97.59±0.0(19)							
l1-21	82.1±0.05(16)	80.67±0.09(15)	74.58±0.07(20)	68±0.13(18)	91.15±0.02(15)	67.24±0.12(15)	96.38±0.02(17)	72.40±0.05(17)	96.87±0.0(20)							
RFS	79.24±0.168(17)	74.95±0.09(6)	71.94±0.10(19)	68±0.21(13)	82.79±0.05(17)	68.67±0.07(18)	96.62±0.01(20)	58.03±0.18(20)	96.97±0.01(20)							
LASSO	95.73±0.02(6)	70.3±0.08(15)	71.29±0.05(18)	81.67±0.08(19)	96.26±0.00(18)	93.22±0.02(10)	100±0.00(10)	87.88±0.03(18)	96.09±0.0(20)							
EN	95.73±0.04(10)	70.3±0.08(15)	68.73±0.10(19)	81.67±0.08(19)	95.97±0.012(18)	93.22±0.02(10)	100±0.00(10)	88.56±0.03(19)	96.09±0.0(20)							
SMBA-CSFS	88.24±0.04(20)	81.93±0.02(20)	75.53±0.06(20)	73.34±0.18(16)	98.41±0.014(19)	97.93±0.03(19)	98.30±0.02(13)	94.95±0.02(19)	99.2±0.01(20)							
BSL	97.89±0.04	98.95±0.021	84.26±0.06	85±0.1	99.57±0.00	98.62±0.02	100±0.00	98.65±0.01	100±0.00							

Table 5. SVM Precision(P), Recall(R) and F-measure(F) on top 20 features using 5-fold CV on different data sets.

GF-CSFS (Pineda – Bautista et al., 2011) framework is compared against our SMBA-CSFS. FS: Fisher Score, mRmR: Minimum-Redundancy-Maximum-Relevance, MI: Mutual Information, RFS: Robust Feature Selector, EN: Elastic Net, BSL: all features. The best results are highlighted in bold. The number in parentheses is the number of features when the performance is achieved.

	ALLAML			LEUKEMIA			CLL.SUB.111			GLIOMA			LUNG.C			LUNG.D			DLBCL			CARCINOM			GCM(14)		
	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F
Fisher	0.96(15)	0.96(14)	0.96	0.97(12)	0.97(12)	0.97	0.84(4)	0.84(4)	0.84	0.76(8)	0.75(8)	0.75	0.96(18)	0.96(18)	0.96	0.97(16)	0.97(16)	0.97	1.0(17)	1.0(17)	1.0	0.95(13)	0.95(13)	0.95	0.93(18)	0.93(18)	0.93
Relief	0.98(16)	0.98(16)	0.98	0.97(18)	0.97(18)	0.97	0.82(5)	0.82(5)	0.82	0.72(19)	0.71(5)	0.71	0.95(19)	0.95(19)	0.95	0.96(9)	0.95(9)	0.95	1.0(10)	1.0(10)	1.0	0.96(17)	0.96(17)	0.96	0.91(20)	0.91(20)	0.91
MRMR	0.69(8)	0.69(8)	0.69	0.97(13)	0.97(14)	0.97	0.84(15)	0.84(15)	0.84	0.72(19)	0.72(19)	0.72	0.95(19)	0.95(19)	0.95	0.97(17)	0.97(17)	0.97	1.0(11)	1.0(11)	1.0	0.96(17)	0.97(15)	0.97	0.91(20)	0.91(20)	0.91
JFS	0.82(18)	0.82(18)	0.82	0.97(17)	0.97(17)	0.97	0.84(15)	0.84(15)	0.84	0.72(19)	0.72(19)	0.72	0.95(19)	0.95(19)	0.95	0.97(17)	0.97(17)	0.97	1.0(11)	1.0(11)	1.0	0.96(17)	0.97(15)	0.97	0.91(20)	0.91(20)	0.91
RFS	0.82(18)	0.78(18)	0.8	0.92(17)	0.91(17)	0.91	0.71(4)	0.69(4)	0.69	0.67(20)	0.67(20)	0.67	0.96(20)	0.96(20)	0.96	0.96(20)	0.96(20)	0.96	0.91(6)	0.91(6)	0.91	0.91(19)	0.91(19)	0.91	0.77(18)	0.77(18)	0.77
JFS	0.82(18)	0.78(18)	0.8	0.92(17)	0.91(17)	0.91	0.71(4)	0.69(4)	0.69	0.67(20)	0.67(20)	0.67	0.96(20)	0.96(20)	0.96	0.96(20)	0.96(20)	0.96	0.91(6)	0.91(6)	0.91	0.91(19)	0.91(19)	0.91	0.77(18)	0.77(18)	0.77
JFS	0.87(14)	0.85(14)	0.86	0.96(19)	0.96(19)	0.96	0.68(12)	0.68(12)	0.68	0.69(20)	0.67(20)	0.68	0.95(20)	0.95(20)	0.95	0.93(19)	0.93(19)	0.92	0.94(20)	0.93(20)	0.93	0.85(19)	0.85(19)	0.85	0.78(19)	0.78(19)	0.78
LASSO	0.87(16)	0.87(16)	0.87	0.72(16)	0.71(16)	0.71	0.78(18)	0.78(18)	0.78	0.88(18)	0.78(18)	0.78	0.94(17)	0.94(17)	0.94	0.98(20)	0.98(20)	0.98	0.97(19)	0.97(19)	0.97	0.94(20)	0.94(20)	0.94	0.73(20)	0.73(20)	0.73
JFS	0.87(16)	0.87(16)	0.87	0.72(16)	0.71(16)	0.71	0.78(18)	0.78(18)	0.78	0.88(18)	0.78(18)	0.78	0.94(17)	0.94(17)	0.94	0.98(20)	0.98(20)	0.98	0.97(19)	0.97(19)	0.97	0.94(20)	0.94(20)	0.94	0.73(20)	0.73(20)	0.73
MA-SFCS	0.83(16)	0.83(16)	0.83	0.86(20)	0.86(20)	0.86	0.67(20)	0.67(20)	0.67	0.82(20)	0.77(20)	0.78	0.94(17)	0.94(17)	0.94	0.99(19)	0.99(19)	0.99	1.0(20)	1.0(20)	1.0	0.99(20)	0.99(20)	0.99	0.73(20)	0.73(20)	0.73

REFERENCES

- Aharon, M., Elad, M., and Bruckstein, A. (2006). K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11):4311–4322.
- Alizadeh, A. A., Eisen, M. B., Davis, R. E., Ma, C., Lossos, I. S., Rosenwald, A., Boldrick, J. C., Sabet, H., Tran, T., Yu, X., Powell, J. I., Yang, L., Marti, G. E., Moore, T., Hudson, J. J., Lu, L., Lewis, D. B., Tibshirani, R., Sherlock, G., Chan, W. C., Greiner, T. C., Weisenburger, D. D., Armitage, J. O., Warnke, R., Levy, R., Wilson, W., Grever, M. R., Byrd, J. C., Botstein, D., Brown, P., and Staudt, L. M. (2000). Distinct types of diffuse large b-cell lymphoma identified by gene expression profiling. *Nature*, 403(6769):503–511.
- Bhattacharjee, A., Richards, W. G., Staunton, J., Li, C., Monti, S., Vasa, P., Ladd, C., Beheshti, J., Bueno, R., Gillette, M., Loda, M., Weber, G., Mark, E. J., Lander, E. S., Wong, W., Johnson, B. E., Golub, T. R., Sugarbaker, D. J., and Meyerson, M. (2001). Classification of human lung carcinomas by mrna expression profiling reveals distinct adenocarcinoma subclasses. *Proceedings of the National Academy of Sciences*, 98(24):13790–13795.
- Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.*, 3(1):1–122.
- Calcagno, G., Staiano, A., Fortunato, G., Brescia-Morra, V., Salvatore, E., Liguori, R., Capone, S., Filla, A., Longo, G., and Sacchetti, L. (2010). A multilayer perceptron neural network-based approach for the identification of responsiveness to interferon therapy in multiple sclerosis patients. *Information Sciences*, 180(21):4153–4163.
- Camasta, F., Di Taranto, M., and Staiano, A. (2015). Statistical and computational methods for genetic diseases: An overview. *Computational and Mathematical Methods in Medicine*, 2015(Article ID 954598).
- Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P. (2002). Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357.
- Ciaramella, A., Cocozza, S., Iorio, F., Miele, G., Napolitano, F., Pinelli, M., Raiconi, G., and Tagliaferri, R. (2008). Interactive data analysis and clustering of genomic data. *Neural Networks*, 21(2-3):368–378.
- Ciaramella, A., Gianfico, M., and Giunta, G. (2016). Compressive sampling and adaptive dictionary learning for the packet loss recovery in audio multimedia streaming. *Multimedia Tools and Applications*, 75(24):17375–17392.
- Ciaramella, A. and Giunta, G. (2016). Packet loss recovery in audio multimedia streaming by using compressive sensing. *IET Communications*, 10(4):387–392.
- Demetriou, C. A., Chen, J., Polidoro, S., Van Veldhoven, K., Cuenin, C., Campanella, G., Brennan, K., Clavel-Chapelon, F., Dossus, L., Kvaskoff, M., Drohan, D., Boeing, H., Kaaks, R., Risch, A., Trichopoulos, D., Lagiou, P., Masala, G., Sieri, S., Tumino, R., Panico, S., Quirós, J. R., Sánchez Perez, M. J., Amiano, P., Huerta Castaño, J. M., Ardanaz, E., Onland-Moret, C., Peeters, P., Khaw, K. T., Wareham, N., Key, T. J., Travis, R. C., Romieu, I., Gallo, V., Gunter, M., Herceg, Z., Kyriacou, K., Riboli, E., Flanagan, J. M., and Vineis, P. (2013). Methylome analysis and epigenetic changes associated with menarcheal age. *PloS one*, 8(11):e79391.
- Demšar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research*, 7(Jan):1–30.
- Deng, W., Lai, M.-J., Peng, Z., and Yin, W. (2017). Parallel multi-block admm with o (1/k) convergence. *Journal of Scientific Computing*, 71(2):712–736.
- Di Taranto, M. D., Staiano, A., D’Agostino, M. N., D’Angelo, A., Bloise, E., Morgante, A., Marotta, G., Gentile, M., Rubba, P., and Fortunato, G. (2015). Association of usf1 and apoa5 polymorphisms with familial combined hyperlipidemia in an italian population. *Molecular and cellular probes*, 29(1):19–24.
- Draghici, S., Khatri, P., Eklund, A., and Szallasi, Z. (2006). Reliability and reproducibility issues in dna microarray measurements. *Trends Genet.*, 22(2):101–109.
- Dunn, O. J. (1961). Multiple comparisons among means. *Journal of the American Statistical Association*, 56(293):52–64.
- Elhamifar, E., Sapiro, G., and Vidal, R. (2012). See all by looking at a few: Sparse modeling for finding representative objects. In *IEEE Conf. Comput. Vis. Pattern Recognit.*, pages 1600–1607. IEEE.
- Engan, K., Aase, S. O., and Husoy, J. H. (1999). Method of optimal directions for frame design. In *Acoustics, Speech, and Signal Processing, 1999. Proceedings., 1999 IEEE International Conference on*, volume 5, pages 2443–2446. IEEE.

- 394 Friedman, J., Hastie, T., and Tibshirani, R. (2001). *The elements of statistical learning*, volume 1. Springer,
395 New-York.
- 396 Friedman, M. (1937). The use of ranks to avoid the assumption of normality implicit in the analysis of
397 variance. *Journal of the american statistical association*, 32(200):675–701.
- 398 Fu, X. and Wang, L. (2002). A ga-based rbf classifier with class-dependent features. In *Evolutionary*
399 *Computation, 2002. CEC'02. Proceedings of the 2002 Congress on*, volume 2, pages 1890–1894. IEEE.
- 400 Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., Coller, H., Loh,
401 M. L., Downing, J. R., Caligiuri, M. A., Bloomfield, C. D., and Lander, E. S. (1999). Molecular
402 classification of cancer: class discovery and class prediction by gene expression monitoring. *science*,
403 286(5439):531–537.
- 404 Gu, Q., Li, Z., and Han, J. (2012). Generalized fisher score for feature selection. *arXiv preprint*
405 *arXiv:1202.3a725*.
- 406 Guyon, I. and Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of Machine*
407 *Learning Research*, 3:1157–1182.
- 408 Haslinger, C., Schweifer, N., Stilgenbauer, S., Döhner, H., Lichter, P., Kraut, N., Stratowa, C., and
409 Abseher, R. (2004). Microarray gene expression profiling of b-cell chronic lymphocytic leukemia
410 subgroups defined by genomic aberrations and vh mutation status. *Journal of Clinical Oncology*,
411 22(19):3937–3949.
- 412 He, X., Cai, D., and Niyogi, P. (2005). Laplacian score for feature selection, advances in neural information
413 processing systems.
- 414 Hoque, N., Bhattacharyya, D. K., and Kalita, J. K. (2014). Mifs-nd: A mutual information-based feature
415 selection method. *Expert Systems with Applications*, 41(14):6371–6385.
- 416 Jolliffe, I. T. (1986). Principal component analysis and factor analysis. In *Principal component analysis*,
417 pages 115–128. Springer, New York, NY.
- 418 Jović, A., Brkić, K., and Bogunović, N. (2015). A review of feature selection methods with applications.
419 In *Information and Communication Technology, Electronics and Microelectronics (MIPRO), 2015 38th*
420 *International Convention on*, pages 1200–1205. IEEE.
- 421 Kira, K. and Rendell, L. A. (1992). A practical approach to feature selection. In *Proceedings of the ninth*
422 *international workshop on Machine learning*, pages 249–256.
- 423 Kononenko, I. (1994). Estimating attributes: analysis and extensions of relief. In *European conference on*
424 *machine learning*, pages 171–182. Springer.
- 425 Kraskov, A., Stögbauer, H., and Grassberger, P. (2004). Estimating mutual information. *Physical review*
426 *E*, 69(6):066138.
- 427 Kreyszig, E. (2010). *Advanced engineering mathematics*. John Wiley & Sons, Great Britain.
- 428 Mairal, J., Bach, F., Ponce, J., Sapiro, G., and Zisserman, A. (2008). Discriminative learned dictionaries
429 for local image analysis. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE*
430 *Conference on*, pages 1–8. IEEE.
- 431 Mairal, J., Bach, F., Ponce, J., Sapiro, G., and Zisserman, A. (2009). Non-local sparse models for image
432 restoration. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 2272–2279.
433 IEEE.
- 434 Nardone, D., Ciaramella, A., and Staiano, A. (2019a). Biological datasets. [https://zenodo.org/
435 record/3405292#.XXkAtugzaUk](https://zenodo.org/record/3405292#.XXkAtugzaUk).
- 436 Nardone, D., Ciaramella, A., and Staiano, A. (2019b). Source code. [https://github.com/
437 DavideNardone/A-Sparse-Coding-Based-Approach-for-Class-Specific-Feature-Sele](https://github.com/DavideNardone/A-Sparse-Coding-Based-Approach-for-Class-Specific-Feature-Sele)
- 438 Nie, F., Huang, H., Cai, X., and Ding, C. H. (2010). Efficient and robust feature selection via joint
439 $\ell_{2,1}$ -norms minimization. In *Advances in neural information processing systems*, pages 1813–1821.
- 440 Nutt, C. L., Mani, D., Betensky, R. A., Tamayo, P., Cairncross, J. G., Ladd, C., Pohl, U., Hartmann, C.,
441 McLaughlin, M. E., Batchelor, T. T., Black, P. M., von Deimling, A., Pomeroy, S. L., Golub, S. L., and
442 Louis, D. N. (2003). Gene expression-based classification of malignant gliomas correlates better with
443 survival than histological classification. *Cancer research*, 63(7):1602–1607.
- 444 Peng, H., Long, F., and Ding, C. (2005). Feature selection based on mutual information criteria of
445 max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on pattern analysis and*
446 *machine intelligence*, 27(8):1226–1238.
- 447 Pineda-Bautista, B. B., Carrasco-Ochoa, J. A., and Martinez-Trinidad, J. F. (2011). General framework
448 for class-specific feature selection. *Expert Systems with Applications*, 38(8):10018–10024.

- 449 Ramaswamy, S., Tamayo, P., Rifkin, R., Mukherjee, S., Yeang, C.-H., Angelo, M., Ladd, C., Reich,
450 M., Latulippe, E., Mesirov, J. P., Poggio, T., Gerald, W. L., Loda, M. F., Lander, E. S., and Golub,
451 T. R. (2001). Multiclass cancer diagnosis using tumor gene expression signatures. *Proceedings of the*
452 *National Academy of Sciences*, 98(26):15149–15154.
- 453 Ramirez, I., Sprechmann, P., and Sapiro, G. (2010). Classification and clustering via dictionary learning
454 with structured incoherence and shared features. In *Computer Vision and Pattern Recognition (CVPR),*
455 *2010 IEEE Conference on*, pages 3501–3508. IEEE.
- 456 Rodríguez-Fdez, I., Canosa, A., Mucientes, M., and Bugarín, A. (2015). Stac: a web platform for
457 the comparison of algorithms using statistical tests. In *Fuzzy Systems (FUZZ-IEEE), 2015 IEEE*
458 *International Conference on*, pages 1–8. IEEE.
- 459 Ross, B. C. (2014). Mutual information between discrete and continuous data sets. *PloS one*, 9(2):e87357.
- 460 Saeys, Y., Inza, I., and Larrañaga, P. (2007). A review of feature selection techniques in bioinformatics.
461 *Bioinformatics*, 23(19):2507–2517.
- 462 Sokolova, M. and Lapalme, G. (2009). A systematic analysis of performance measures for classification
463 tasks. *Information processing & management*, 45(4):427–437.
- 464 Staiano, A., De Vinco, L., Ciaramella, A., Raiconi, G., Tagliaferri, R., Amato, R., Longo, G., Donalek, C.,
465 Miele, G., and Di Bernardo, D. (2004). Probabilistic principal surfaces for yeast gene microarray data
466 mining. In *Proceedings of the Fourth IEEE International Conference on Data Mining, ICDM 2004*,
467 pages 202–208. IEEE.
- 468 Staiano, A., Di Taranto, M. D., Bloise, E., D’Agostino, M. N., D’Angelo, A., Marotta, G., Gentile,
469 M., Jossa, F., Iannuzzi, A., Rubba, P., and Fortunato, G. (2013). Investigation of single nucleotide
470 polymorphisms associated to familial combined hyperlipidemia with random forests. In *Neural Nets*
471 *and Surroundings*, volume 19, pages 169–178. Springer, Berlin, Heidelberg.
- 472 Su, A. I., Welsh, J. B., Sapinoso, L. M., Kern, S. G., Dimitrov, P., Lapp, H., Schultz, P. G., Powell,
473 S. M., Moskaluk, C. A., Frierson, H., and Hampton, G. M. (2001). Molecular classification of human
474 carcinomas by use of gene expression signatures. *Cancer research*, 61(20):7388–7393.
- 475 Tang, J., Alelyani, S., and Liu, H. (2014). Feature selection for classification: A review. *Data Classifica-*
476 *tion: Algorithms and Applications*, page 37.
- 477 Tibshirani, R. (1994). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*
478 *Society, Series B*, 58:267–288.
- 479 Wolpert, D. H. and Macready, W. G. (1997). No free lunch theorems for optimization. *IEEE Transactions*
480 *on Evolutionary Computation*, 1(1):67–82.
- 481 Xiong, M., Fang, X., and Zhao, J. (2001). Biomarker identification by feature wrappers. *Genome*
482 *Research*, 11(11):1878–1887.
- 483 Yang, K., Cai, Z., Li, J., and Lin, G. (2006). A stable gene selection in microarray data analysis. *BMC*
484 *bioinformatics*, 7(1):228.
- 485 Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the*
486 *Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320.