

Statistical dependence measure for feature selection in microarray datasets

Verónica Bolón-Canedo¹, Sohan Seth², Noelia Sánchez-Maróño¹,
Amparo Alonso-Betanzos¹ and José C. Príncipe² *

1- Department of Computer Science
University of A Coruña - Spain

2- Department of Electrical and Computer Engineering
University of Florida, Gainesville, Florida, USA

Abstract. Feature selection is the domain of machine learning which studies data-driven methods to select, among a set of input variables, the ones that will lead to the most accurate predictive model. In this paper, a statistical dependence measure is presented for variable selection in the context of classification. Its performance is tested over DNA microarray data, a challenging dataset for machine learning researchers due to the high number of genes and relatively small number of measurements. This measure is compared against the so called mRMR approach, and is shown to obtain better or equal performance over the binary datasets.

1 Introduction

A typical classification task is to separate healthy patients from cancer patients based on their gene expression “profile” (DNA microarray) [1]. DNA microarray data *classification* is a challenging problem for machine learning researchers due to its high number of features (from 6000 to 60 000 genes) and the small number of samples (often less than 100 patients), and therefore, *feature selection* (FS) plays a crucial role in this domain [1]. Among the different FS methods [2], *filters* only rely on general characteristics of the data, and not on the learning machines; therefore, they are faster, and more suitable for large data sets. A common practice in this approach is to simply select the *top-ranked* genes where the ranks are determined by some *dependence criteria*, and the number of genes to retain is usually set by human intuition with trial-and-error. A deficiency of this ranking approach is that the selected features could be dependent among themselves. Therefore, a maximum relevance minimum redundancy approach is preferred in practice [3], that also minimizes the dependence among selected features.

In [4], the authors have followed this approach using a simple measure of *monotone dependence* to quantify both relevance and redundancy. This approach has been shown to perform equally well compared to other widely used measures of dependence (such as correlation coefficient, mutual information and Hilbert-Schmidt independence criterion). However, this method has not been applied in the context of classification. Therefore, in this paper, we extend this

*This work was supported by Spanish Ministerio de Ciencia e Innovación (under project TIN 2009-10748) partially supported by the European Union ERDF.

approach, and test it over 16 microarray datasets in a classification scenario. We compare its performance with the so called mRMR (minimum Redundancy Maximum Relevance) approach as proposed in [3], that uses crude binning based estimator of mutual information to capture relevance and redundancy. In [5] it has been demonstrated that mRMR is an appropriate tool for gene selection in microarray datasets, since genes selected by this framework provide a more balanced coverage of the space and capture broader characteristics of phenotypes. However, we observe that our approach outperforms mRMR on binary problems. Finally, we also compare our results over 10 out of the 16 datasets with those obtained in [6], where several widely-used filters were applied before the classification stage, including a discretization step.

2 Statistical dependence Measure (\mathcal{M}_d)

In the mRMR approach [3], the features are selected one at a time i.e. at a particular iteration an input variable, that is most relevant to the target and least redundant with respect to the already selected variables, is selected. The relevance and redundancy are often measured in terms of Mutual Information (MI) between the variable and the response and MI between the variable and the already selected variables, respectively. In [4], the authors have proposed to use a measure of monotone dependence (\mathcal{M}) to assess the relevance and the redundancy, since it is simpler to estimate from data (compared to MI), and it carries many desired properties of a measure of dependence (e.g. it is bounded, symmetric, and reaches maximum if and only two random variables share a monotonic relationship). However, the authors have only explored this concept in the context of continuous random variables by using,

$$\mathcal{M}_c(Y, X) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \int (P(Y \leq y_j, X \leq x_i) - P(Y \leq y_j)P(X \leq x_i))^2.$$

In this paper, we explore this method in a classification scenario where the class labels are categorical variables, and explore the following measure of dependence,

$$\mathcal{M}_d(Y, X) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \int (P(Y = y_j, X \leq x_i) - P(Y = y_j)P(X \leq x_i))^2.$$

In brief, the variable selection algorithm can be described as follows, given a set of features $X_{n-1} \subset X$ choose X_n in the following way,

$$X_n = \underset{X_i \in X/X_{n-1}}{\operatorname{argmin}} \left(\mathcal{M}(X_i, Y) - \frac{\lambda}{|X_{n-1}|} \sum_{X_j \in X_{n-1}} \mathcal{M}(X_i, X_j) \right),$$

where Y is the target and λ is a free parameter that controls the relative emphasis given on the relevance and the redundancy. Here, the relevance is evaluated by the dependence between a variable and the target, whereas the redundancy is evaluated by the average dependence between the new variable and the already selected variable.

3 Experimental results

Our goal is to test the previously described measurement (\mathcal{M}_d) in the classification of microarray datasets. This is a challenging problem since the number of features is relatively higher (order of 10^4) than the number of samples (order of 10^2). The experiments were carried out on 16 datasets of gene expression profiles; their properties (number of features, samples and classes) are listed in Table 1. Some of the datasets come originally with training and test samples that were drawn from different conditions. However, in this table we have presented the combined data for the purpose of performing a 10-fold cross-validation.

Table 1: Description for the datasets used (free to download in [7] and [8]).

Dataset	# feat	Samples	# cl	Dataset	# feat	Samples	# cl
Leukemia	7129	72	2	Myeloma	12 625	173	2
CNS	7129	60	2	Brain	12 625	21	2
DLBCL	4026	47	2	Gli	22 283	85	2
Colon	2000	62	2	SMK-CAN	19 993	187	2
Prostate	12 600	136	2	GCM	16 063	190	14
Lung	12 533	181	2	Lymphoma	4026	96	9
Ovarian	15 154	253	2	GLA-BRA	49 151	180	4
Breast	24 481	97	2	TOX	5748	171	4

Five well-known supervised classifiers, of different conceptual origin, were chosen for class prediction over the microarray datasets used in this work. All the classifiers (C4.5, naive Bayes, IB1, MLP and SVM) were executed using the Weka tool [9], with default values for their respective parameters. Over the 16 datasets, a 10-fold cross-validation was performed using 10 different numbers of features (1, 2, 3, 4, 5, 10, 15, 20, 25, 30). \mathcal{M}_d was tested with 5 different values for its parameter λ , which leads to a total of 4800 experiments. Therefore, in this paper, we only present a summary of the experiments. The results obtained show that \mathcal{M}_d performs equally well or even better than mRMR over the binary datasets, while mRMR is better over the multiclass datasets. In Table 2 we can see a summary of the experiments, where a method wins when obtaining the best result in terms of error (percentage of incorrectly classified instances) and, in case of same error, with less number of features. Regarding the classifiers, MLP obtained the best results in 50% of the datasets.

Table 2: Number of times \mathcal{M}_d wins, loses or ties compared to mRMR

\mathcal{M}_d vs mRMR	Wins	Loses	Ties
Binary datasets	9	2	1
Multiclass datasets	0	4	0

For the multiclass datasets, the number of tissue samples per class is small (e.g. in Lymphoma dataset, 6 out of the 9 classes have less than 10 samples) and unevenly distributed (e.g. in Lymphoma dataset, from 46 to 2). This fact, together with the large number of classes (see Table 1), makes the classification

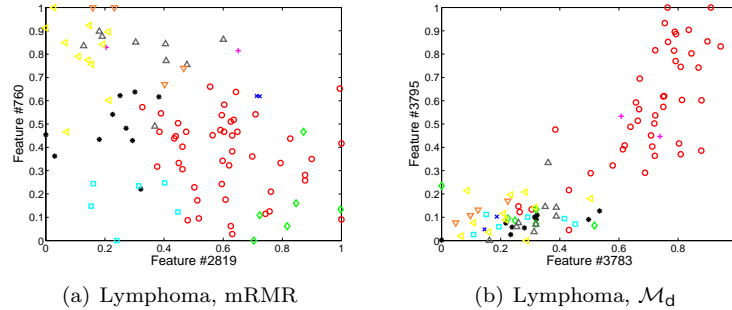


Fig. 1: Results for the dataset Lymphoma

task more complex than the two-class problem. \mathcal{M}_d seems to have difficulties when the number of classes is very large and selects features that would be more suitable for a two-class problem. Figure 1 shows the two first features selected by \mathcal{M}_d and mRMR over Lymphoma dataset (different shapes represent different classes). Note that Figure 1(b) can be almost linearly separated in two subsets, one formed by the class represented by red circles and the other one formed by the remaining classes, whereas this division is not appreciated in Figure 1(a).

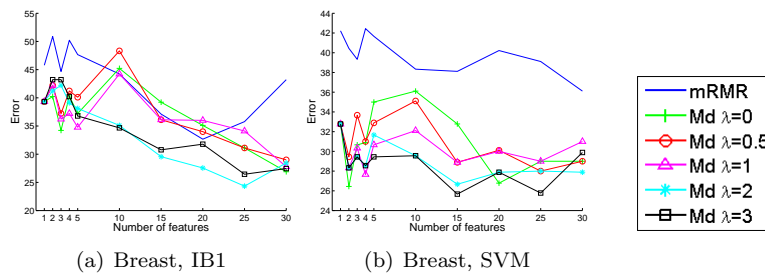


Fig. 2: Results for the dataset Breast with IB1 and SVM classifiers

Regarding the binary case, in Table 2 is shown that \mathcal{M}_d performs better or equal to mRMR in 10 out of the 12 datasets tested. However, in the two datasets where mRMR wins, \mathcal{M}_d achieves close results (for Ovarian dataset, \mathcal{M}_d obtains the same error but using 4 features instead of 3 and for Prostate dataset, there is a difference in error of 0.11% with the same number of features). When \mathcal{M}_d wins, it improves an average of 4.21%, whilst when it loses, its performance decays in 0.06%. As an illustrative example, Figure 2 presents the results for Breast dataset for classifiers IB1 and SVM. For this dataset, mRMR performs worse than \mathcal{M}_d . Specially remarkable is the difference between \mathcal{M}_d and mRMR for the SVM classifier, for any number of features tested. Focusing on the λ parameter, significant differences among these values were not found, although for the Breast dataset, high values of λ seem to be the more suitable.

An interesting result is observed with the Brain dataset. For the C4.5 classi-

fier, \mathcal{M}_d is able to select a *single* feature that is enough to get 0% classification error, where the best classification result obtained by mRMR is about 30%. This is due to the fact that \mathcal{M}_d always selects the #364 feature on the top of the ranking, while mRMR does not. This feature is crucial since it allows linear separation of the two classes when the boundary between the classes is at around 0.0001. However, this attribute of the feature is completely ignored by mRMR since it discretizes the feature in 3 states before using it in the algorithm. It can be easily understood from the contingency table 3 (left), where the rows are the two classes (C1 and C2) and the columns are the 3 states of the discretization (S1, S2 and S3). In the first table the realizations are discretized as suggested by [3] whereas in the second table the discretization involves the line through 0.0001. It is obvious that the second table has much more mutual information than the first, that the common scheme of mRMR overlooks.

	S1	S2	S3
C1	0	11	3
C2	0	7	0

	S1	S2	S3
C1	0	11	3
C2	7	0	0

Table 3: Contingency tables for mRMR (left) and \mathcal{M}_d (right)

4 Comparison with previous results

In a previous work [6], a combination of discretization and filter algorithms was proposed for the crucial task of accurate gene selection in class prediction problems over 10 DNA microarray datasets. Two discretizers (EMD, PKID) and two filters (INTERACT, Consistency-based) were used, and their performance was checked using three supervised classifiers (C4.5, NB, IB1). The datasets used were the 10 first ones in Table 1. For the sake of comparison, the datasets with only training set (marked with *) were randomly divided using the common rule 2/3 for training and 1/3 for testing.

Table 4: Comparison with previous results [6]

Dataset	Method	Test E	# g	Dataset	Method	Test E	# g
Colon*	EMD+Cons+NB	15.00	3	Prostate	PKID+INT+IB1	26.47	2
	Md+IB1	10.00	2		Md+C4.5	2.94	1
DLBCL*	EMD+INT+NB	6.67	36	Ovarian*	EMD+Cons+NB	0.00	3
	Md+MLP	0.00	33		Md+IB1	0.00	2
CNS*	PKID+INT+NB	25.00	4	Breast*	PKID+INT+C4.5	21.05	3
	Md+C4.5	25.00	3		Md+MLP	15.79	3
Leukemia	PKID+Cons+C4.5	5.88	2	Lymph.*	EMD+INT+NB	18.75	160
	Md+NB	0.00	21		Md+IB1	18.75	33
Lung	PKID+INT+IB1	0.00	40	GCM	EMD+Cons+NB	45.65	9
	Md+NB	1.34	37		Md+C4.5	52.17	26

In Table 4 a comparison with the results obtained in this work over the same datasets and conditions are shown. \mathcal{M}_d improves the test error in 5 of the datasets (especially remarkable in Prostate, with an improvement of 23%) . In 3 of them (CNS, Ovarian and Lymphoma) it achieves the same error but using less number of features, which is also considered as an improvement in performance.

Regarding the two remaining datasets, in Lung the error is slightly higher with \mathcal{M}_d but using less features, while GCM is a complex multiclass dataset and we have seen that \mathcal{M}_d does not perform adequately for this kind of datasets.

5 Conclusions

In this work, the adequacy of a measure of monotone dependence (\mathcal{M}_d) was tested in the classification domain, specifically on microarray data, which is a challenging problem for machine learning researchers where a feature selection step is a fundamental necessity. To check the capacity of this method to deal with this problem, it was compared to mRMR, a well-known and widely-used FS method, in terms of five different classifiers.

The results showed that for the binary datasets, \mathcal{M}_d selects better features, and outperforms mRMR in most of the data sets tested (10 out of 12). Moreover, it is remarkable that \mathcal{M}_d does not require a previous discretization of the data, so it can capture the importance of some features in cases where mRMR cannot.

For the multiclass datasets, the performance of \mathcal{M}_d decays. It has to be noted that the multiclass problem is much more complex than the binary one and the design of \mathcal{M}_d is more focused on the binary problem, which is the most common classification problem in the literature. However, as future work, it will be interesting to extend the method in order to handle this situation.

Finally, \mathcal{M}_d was also compared to previous results over several microarray datasets using a suite of well-known FS methods, leading to an improvement in 80% of them, both in number of features and in percentage of error (up to 23%).

References

- [1] Y. Saeys, I. Inza and P. Larrañaga, *A Review of Feature Selection Techniques in Bioinformatics*, Journal of Bioinformatics, vol. 23, no. 19, pp. 2507-2517, 2007
- [2] I. Guyon, S. Gunn, M. Nikravesh and L. Zadeh, *Feature Extraction. Foundations and Applications*, Springer, 2006
- [3] H. Peng, F. Long and C. Ding, *Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy*, IEEE TPAMI, vol. 27, no. 8, pp. 1226-1238, 2005
- [4] S. Seth and J.C. Principe, *Variable Selection: A Statistical Dependence Perspective*, In Proc. ICMLA, pp. 931-936, 2010
- [5] C. Ding and H. Peng, *Minimum Redundancy Feature Selection from Microarray Gene Expression Data*, JBCB, vol. 3, no. 2, pp. 185-206, 2005
- [6] V. Bolón-Canedo, N. Sánchez-Maróño and A. Alonso-Betanzos, *On the Effectiveness of Discretization on Gene Selection of Microarray Data*, In Proc. IJCNN, pp. 3167-3174, 2010
- [7] K. Ridge, *Kent Ridge Bio-Medical Dataset*, 2009 <http://datam.i2r.a-star.edu.sg/datasets/krbd> [Last access: October 2010]
- [8] Broad Institute. Cancer Program Data Sets <http://www.broadinstitute.org/cgi-bin/cancer/datasets.cgi> [Last access: October 2010]
- [9] I.H. Witten and E. Frank, *Data Mining: Practical machine learning tools and techniques*, 2nd Edition, Morgan Kaufmann, San Francisco, 2005 <http://www.cs.waikato.ac.nz/ml/weka/> [Last access: October 2010]